
QEMU Documentation

Release 5.2.50

The QEMU Project Developers

Feb 20, 2021

Contents:

1	QEMU System Emulation User's Guide	1
1.1	Quick Start	1
1.2	Invocation	1
1.3	Keys in the graphical frontends	51
1.4	Keys in the character backend multiplexer	52
1.5	QEMU Monitor	52
1.6	Disk Images	61
1.7	Network emulation	74
1.8	QEMU virtio-net standby (net_failover)	75
1.9	USB emulation	76
1.10	Inter-VM Shared Memory device	78
1.11	Direct Linux Boot	79
1.12	VNC security	79
1.13	TLS setup for network services	81
1.14	GDB usage	86
1.15	Managed start up options	87
1.16	Virtual CPU hotplug	88
1.17	virtio pmem	90
1.18	Persistent reservation managers	91
1.19	QEMU System Emulator Targets	92
1.20	Security	158
1.21	Multi-process QEMU	160
1.22	Deprecated features	161
1.23	Removed features	169
1.24	Supported build platforms	173
1.25	License	174
2	QEMU User Mode Emulation User's Guide	175
2.1	QEMU User space emulator	175
3	QEMU Tools Guide	179
3.1	QEMU disk image utility	179
3.2	QEMU Storage Daemon	190
3.3	QEMU Disk Network Block Device Server	192
3.4	QEMU persistent reservation helper	196
3.5	QEMU SystemTap trace tool	197
3.6	QEMU 9p virtfs proxy filesystem helper	199

3.7	QEMU virtio-fs shared file system daemon	200
4	QEMU System Emulation Management and Interoperability Guide	205
4.1	Dirty Bitmaps and Incremental Backup	205
4.2	D-Bus	229
4.3	D-Bus VMState	230
4.4	Live Block Device Operations	231
4.5	Persistent reservation helper protocol	248
4.6	QEMU Guest Agent	249
4.7	QEMU Guest Agent Protocol Reference	251
4.8	QEMU QMP Reference Manual	280
4.9	QEMU Storage Daemon QMP Reference Manual	629
4.10	Vhost-user Protocol	782
4.11	Vhost-user-gpu Protocol	807
4.12	Vhost-vdpa Protocol	811
5	QEMU System Emulation Guest Hardware Specifications	813
5.1	POWER9 XIVE interrupt controller	813
5.2	XIVE for sPAPR (pseries machines)	816
5.3	NUMA mechanics for sPAPR (pseries machines)	820
5.4	How the pseries Linux guest calculates NUMA distances	822
5.5	pseries NUMA mechanics	822
5.6	Legacy (5.1 and older) pseries NUMA mechanics	825
5.7	QEMU and ACPI BIOS Generic Event Device interface	826
5.8	QEMU TPM Device	827
5.9	APEI tables generating and CPER record	834
6	QEMU Developer's Guide	837
6.1	The QEMU build system architecture	837
6.2	QEMU and Kconfig	843
6.3	Testing in QEMU	847
6.4	Fuzzing	863
6.5	Control-Flow Integrity (CFI)	866
6.6	Load and Store APIs	868
6.7	The memory API	877
6.8	Migration	906
6.9	Atomic operations in QEMU	919
6.10	QEMU and the stable process	925
6.11	QTest Device Emulation Testing Framework	926
6.12	Decodetree Specification	942
6.13	Secure Coding Practices	946
6.14	Translator Internals	947
6.15	TCG Instruction Counting	949
6.16	Tracing	950
6.17	Introduction	958
6.18	vCPU Scheduling	958
6.19	Shared Data Structures	958
6.20	Memory Consistency	961
6.21	QEMU TCG Plugins	963
6.22	Bitwise operations	967
6.23	Reset in QEMU: the Resettable interface	974
6.24	Bootimg from real channel-attached devices on s390x	978
6.25	Modelling a clock tree in QEMU	980
6.26	The QEMU Object Model (QOM)	986

6.27	block-coroutine-wrapper	1022
6.28	Multi-process QEMU	1023
Index		1037

QEMU System Emulation User's Guide

This manual is the overall guide for users using QEMU for full system emulation (as opposed to user-mode emulation). This includes working with hypervisors such as KVM, Xen, Hax or Hypervisor.Framework.

Contents:

1.1 Quick Start

Download and uncompress a PC hard disk image with Linux installed (e.g. `linux.img`) and type:

```
qemu-system-x86_64 linux.img
```

Linux should boot and give you a prompt.

1.2 Invocation

```
qemu-system-x86_64 [options] [disk_image]
```

`disk_image` is a raw hard disk image for IDE hard disk 0. Some targets do not need a disk image.

1.2.1 Standard options

-h Display help and exit

-version Display version information and exit

-machine [**type=**]**name**[,**prop=value**[, ...]] Select the emulated machine by name. Use `-machine help` to list available machines.

For architectures which aim to support live migration compatibility across releases, each release will introduce a new versioned machine type. For example, the 2.8.0 release introduced machine types “pc-i440fx-2.8” and “pc-q35-2.8” for the x86_64/i686 architectures.

To allow live migration of guests from QEMU version 2.8.0, to QEMU version 2.9.0, the 2.9.0 version must support the “pc-i440fx-2.8” and “pc-q35-2.8” machines too. To allow users live migrating VMs to skip multiple intermediate releases when upgrading, new releases of QEMU will support machine types from many previous versions.

Supported machine properties are:

accel=accels1[:accels2[:...]] This is used to enable an accelerator. Depending on the target architecture, kvm, xen, hax, hvf, whpx or tcg can be available. By default, tcg is used. If there is more than one accelerator specified, the next one is used if the previous one fails to initialize.

vmport=on|off|auto Enables emulation of VMWare IO port, for vmmouse etc. auto says to select the value based on accel. For accel=xen the default is off otherwise the default is on.

dump-guest-core=on|off Include guest memory in a core dump. The default is on.

mem-merge=on|off Enables or disables memory merge support. This feature, when supported by the host, de-duplicates identical memory pages among VMs instances (enabled by default).

aes-key-wrap=on|off Enables or disables AES key wrapping support on s390-ccw hosts. This feature controls whether AES wrapping keys will be created to allow execution of AES cryptographic functions. The default is on.

dea-key-wrap=on|off Enables or disables DEA key wrapping support on s390-ccw hosts. This feature controls whether DEA wrapping keys will be created to allow execution of DEA cryptographic functions. The default is on.

nvdimm=on|off Enables or disables NVDIMM support. The default is off.

memory-encryption= Memory encryption object to use. The default is none.

hmat=on|off

Enables or disables ACPI Heterogeneous Memory Attribute Table (HMAT) support. The default is off.

memory-backend='id' An alternative to legacy `-mem-path` and `mem-prealloc` options. Allows to use a memory backend as main RAM.

For example: `:: -object memory-backend-file,id=pc.ram,size=512M,mem-path=/hugetlbfs,prealloc=on,share=on -machine memory-backend=pc.ram -m 512M`

Migration compatibility note: a) as backend id one shall use value of ‘default-ram-id’, advertised by machine type (available via `query-machines QMP` command), if migration to/from old QEMU (<5.0) is expected. b) for machine types 4.0 and older, user shall use `x-use-canonical-path-for-ramblock-id=off` backend option if migration to/from old QEMU (<5.0) is expected. For example: `:: -object memory-backend-ram,id=pc.ram,size=512M,x-use-canonical-path-for-ramblock-id=off -machine memory-backend=pc.ram -m 512M`

-cpu model Select CPU model (`-cpu help` for list and additional feature selection)

-accel name[,prop=value[,...]] This is used to enable an accelerator. Depending on the target architecture, kvm, xen, hax, hvf, whpx or tcg can be available. By default, tcg is used. If there is more than one accelerator specified, the next one is used if the previous one fails to initialize.

igd-passthru=on|off When Xen is in use, this option controls whether Intel integrated graphics devices can be passed through to the guest (default=off)

kernel-irqchip=on|off|split Controls KVM in-kernel irqchip support. The default is full acceleration of the interrupt controllers. On x86, split irqchip reduces the kernel attack surface, at a performance cost for non-MSI interrupts. Disabling the in-kernel irqchip completely is not recommended except for debugging purposes.

kvm-shadow-mem=size Defines the size of the KVM shadow MMU.

split-wx=on|off Controls the use of split w^x mapping for the TCG code generation buffer. Some operating systems require this to be enabled, and in such a case this will default on. On other operating systems, this will default off, but one may enable this for testing or debugging.

tb-size=n Controls the size (in MiB) of the TCG translation block cache.

thread=single|multi Controls number of TCG threads. When the TCG is multi-threaded there will be one thread per vCPU therefor taking advantage of additional host cores. The default is to enable multi-threading where both the back-end and front-ends support it and no incompatible TCG features have been enabled (e.g. icount/replay).

-smp [cpus=n[,cores=cores][,threads=threads][,dies=dies][,sockets=sockets][,maxcpus=maxcpus]

Simulate an SMP system with n CPUs. On the PC target, up to 255 CPUs are supported. On Sparc32 target, Linux limits the number of usable CPUs to 4. For the PC target, the number of cores per die, the number of threads per cores, the number of dies per packages and the total number of sockets can be specified. Missing values will be computed. If any on the three values is given, the total number of CPUs n can be omitted. maxcpus specifies the maximum number of hotpluggable CPUs.

-numa node[,mem=size][,cpus=firstcpu[-lastcpu]][,nodeid=node][,initiator=initiator]

-numa node[,memdev=id][,cpus=firstcpu[-lastcpu]][,nodeid=node][,initiator=initiator]

-numa dist,src=source,dst=destination,val=distance

-numa cpu,node-id=node[,socket-id=x][,core-id=y][,thread-id=z]

-numa hmat-lb,initiator=node,target=node,hierarchy=hierarchy,data-type=tpye[,latency=lat][,lines=lines]

-numa hmat-cache,node-id=node,size=size,level=level[,associativity=str][,policy=str][,lines=lines]

Define a NUMA node and assign RAM and VCPUs to it. Set the NUMA distance from a source node to a destination node. Set the ACPI Heterogeneous Memory Attributes for the given nodes.

Legacy VCPU assignment uses ‘cpus’ option where firstcpu and lastcpu are CPU indexes. Each ‘cpus’ option represent a contiguous range of CPU indexes (or a single VCPU if lastcpu is omitted). A non-contiguous set of VCPUs can be represented by providing multiple ‘cpus’ options. If ‘cpus’ is omitted on all nodes, VCPUs are automatically split between them.

For example, the following option assigns VCPUs 0, 1, 2 and 5 to a NUMA node:

```
-numa node,cpus=0-2,cpus=5
```

‘cpu’ option is a new alternative to ‘cpus’ option which uses ‘socket-id|core-id|thread-id’ properties to assign CPU objects to a node using topology layout properties of CPU. The set of properties is machine specific, and depends on used machine type/‘smp’ options. It could be queried with ‘hotpluggable-cpus’ monitor command. ‘node-id’ property specifies node to which CPU object will be assigned, it’s required for node to be declared with ‘node’ option before it’s used with ‘cpu’ option.

For example:

```
-M pc \
-smp 1,sockets=2,maxcpus=2 \
-numa node,nodeid=0 -numa node,nodeid=1 \
-numa cpu,node-id=0,socket-id=0 -numa cpu,node-id=1,socket-id=1
```

Legacy ‘mem’ assigns a given RAM amount to a node (not supported for 5.1 and newer machine types). ‘memdev’ assigns RAM from a given memory backend device to a node. If ‘mem’ and ‘memdev’ are omitted in all nodes, RAM is split equally between them.

‘mem’ and ‘memdev’ are mutually exclusive. Furthermore, if one node uses ‘memdev’, all of them have to use it.

‘initiator’ is an additional option that points to an initiator NUMA node that has best performance (the lowest latency or largest bandwidth) to this NUMA node. Note that this option can be set only when the machine property ‘hmat’ is set to ‘on’.

Following example creates a machine with 2 NUMA nodes, node 0 has CPU, node 1 has only memory, and its initiator is node 0. Note that because node 0 has CPU, by default the initiator of node 0 is itself and must be itself.

```
-machine hmat=on \
-m 2G,slots=2,maxmem=4G \
-object memory-backend-ram,size=1G,id=m0 \
-object memory-backend-ram,size=1G,id=m1 \
-numa node,nodeid=0,memdev=m0 \
-numa node,nodeid=1,memdev=m1,initiator=0 \
-smp 2,sockets=2,maxcpus=2 \
-numa cpu,node-id=0,socket-id=0 \
-numa cpu,node-id=0,socket-id=1
```

source and destination are NUMA node IDs. distance is the NUMA distance from source to destination. The distance from a node to itself is always 10. If any pair of nodes is given a distance, then all pairs must be given distances. Although, when distances are only given in one direction for each pair of nodes, then the distances in the opposite directions are assumed to be the same. If, however, an asymmetrical pair of distances is given for even one node pair, then all node pairs must be provided distance values for both directions, even when they are symmetrical. When a node is unreachable from another node, set the pair’s distance to 255.

Note that the -numa option doesn’t allocate any of the specified resources, it just assigns existing resources to NUMA nodes. This means that one still has to use the -m, -smp options to allocate RAM and VCPUs respectively.

Use ‘hmat-lb’ to set System Locality Latency and Bandwidth Information between initiator and target NUMA nodes in ACPI Heterogeneous Attribute Memory Table (HMAT). Initiator NUMA node can create memory requests, usually it has one or more processors. Target NUMA node contains addressable memory.

In ‘hmat-lb’ option, node are NUMA node IDs. hierarchy is the memory hierarchy of the target NUMA node: if hierarchy is ‘memory’, the structure represents the memory performance; if hierarchy is ‘first-levelsecond-levelthird-level’, this structure represents aggregated performance of memory side caches for each domain. type of ‘data-type’ is type of data represented by this structure instance: if ‘hierarchy’ is ‘memory’, ‘data-type’ is ‘accessreadlwrite’ latency or ‘accessreadlwrite’ bandwidth of the target memory; if ‘hierarchy’ is ‘first-levelsecond-levelthird-level’, ‘data-type’ is ‘accessreadlwrite’ hit latency or ‘accessreadlwrite’ hit bandwidth of the target memory side cache.

lat is latency value in nanoseconds. bw is bandwidth value, the possible value and units are NUM[M|G|T], mean that the bandwidth value are NUM byte per second (or MB/s, GB/s or TB/s depending on used suffix). Note that if latency or bandwidth value is 0, means the corresponding latency or bandwidth information is not provided.

In ‘hmat-cache’ option, node-id is the NUMA-id of the memory belongs. size is the size of memory side cache in bytes. level is the cache level described in this structure, note that the cache level 0 should not be used with ‘hmat-cache’ option. associativity is the cache associativity, the possible value is ‘none/direct(direct-mapped)/complex(complex cache indexing)’. policy is the write policy. line is the cache Line size in bytes.

For example, the following options describe 2 NUMA nodes. Node 0 has 2 cpus and a ram, node 1 has only a ram. The processors in node 0 access memory in node 0 with access-latency 5 nanoseconds, access-bandwidth is

200 MB/s; The processors in NUMA node 0 access memory in NUMA node 1 with access-latency 10 nanoseconds, access-bandwidth is 100 MB/s. And for memory side cache information, NUMA node 0 and 1 both have 1 level memory cache, size is 10KB, policy is write-back, the cache Line size is 8 bytes:

```
-machine hmat=on \
-m 2G \
-object memory-backend-ram,size=1G,id=m0 \
-object memory-backend-ram,size=1G,id=m1 \
-smp 2 \
-numa node,nodeid=0,memdev=m0 \
-numa node,nodeid=1,memdev=m1,initiator=0 \
-numa cpu,node-id=0,socket-id=0 \
-numa cpu,node-id=0,socket-id=1 \
-numa hmat-lb,initiator=0,target=0,hierarchy=memory,data-type=access-latency,
↪latency=5 \
-numa hmat-lb,initiator=0,target=0,hierarchy=memory,data-type=access-bandwidth,
↪bandwidth=200M \
-numa hmat-lb,initiator=0,target=1,hierarchy=memory,data-type=access-latency,
↪latency=10 \
-numa hmat-lb,initiator=0,target=1,hierarchy=memory,data-type=access-bandwidth,
↪bandwidth=100M \
-numa hmat-cache,node-id=0,size=10K,level=1,associativity=direct,policy=write-
↪back,line=8 \
-numa hmat-cache,node-id=1,size=10K,level=1,associativity=direct,policy=write-
↪back,line=8
```

-add-fd fd=fd,set=set[,opaque=opaque] Add a file descriptor to an fd set. Valid options are:

fd=fd This option defines the file descriptor of which a duplicate is added to fd set. The file descriptor cannot be stdin, stdout, or stderr.

set=set This option defines the ID of the fd set to add the file descriptor to.

opaque=opaque This option defines a free-form string that can be used to describe fd.

You can open an image using pre-opened file descriptors from an fd set:

```
qemu-system-x86_64 \
-add-fd fd=3,set=2,opaque="rdwr:/path/to/file" \
-add-fd fd=4,set=2,opaque="ronly:/path/to/file" \
-drive file=/dev/fdset/2,index=0,media=disk
```

-set group.id.arg=value Set parameter arg for item id of type group

-global driver.prop=value

-global driver=driver,property=property,value=value Set default value of driver's property prop to value, e.g.:

```
qemu_system-x86_64 -global ide-hd.physical_block_size=4096 disk-image.img
```

In particular, you can use this to set driver properties for devices which are created automatically by the machine model. To create a device which is not created automatically and set properties on it, use `-device`.

`-global driver.prop=value` is shorthand for `-global driver=driver,property=prop,value=value`. The longhand syntax works even when driver contains a dot.

-boot [order=drives][,once=drives][,menu=on|off][,splash=sp_name][,splash-time=sp_time][,reboot=on|off]

Specify boot order drives as a string of drive letters. Valid drive letters depend on the target architecture. The x86 PC uses: a, b (floppy 1 and 2), c (first hard disk), d (first CD-ROM), n-p (Etherboot from network adapter 1-4), hard disk boot is the default. To apply a particular boot order only on the first startup, specify it via

once. Note that the `order` or `once` parameter should not be used together with the `bootindex` property of devices, since the firmware implementations normally do not support both at the same time.

Interactive boot menus/prompts can be enabled via `menu=on` as far as firmware/BIOS supports them. The default is non-interactive boot.

A splash picture could be passed to bios, enabling user to show it as logo, when option `splash=sp_name` is given and `menu=on`. If firmware/BIOS supports them. Currently Seabios for X86 system support it. limitation: The splash file could be a jpeg file or a BMP file in 24 BPP format(true color). The resolution should be supported by the SVGA mode, so the recommended is 320x240, 640x480, 800x640.

A timeout could be passed to bios, guest will pause for `rb_timeout` ms when boot failed, then reboot. If `rb_timeout` is `-1`, guest will not reboot, qemu passes `-1` to bios by default. Currently Seabios for X86 system support it.

Do strict boot via `strict=on` as far as firmware/BIOS supports it. This only effects when boot priority is changed by `bootindex` options. The default is non-strict boot.

```
# try to boot from network first, then from hard disk
qemu_system-x86_64 -boot order=nc
# boot from CD-ROM first, switch back to default order after reboot
qemu_system-x86_64 -boot once=d
# boot with a splash picture for 5 seconds.
qemu_system-x86_64 -boot menu=on,splash=/root/boot.bmp,splash-time=5000
```

Note: The legacy format `-boot drives` is still supported but its use is discouraged as it may be removed from future versions.

-m [size=]megs[, slots=n,maxmem=size] Sets guest startup RAM size to megs megabytes. Default is 128 MiB. Optionally, a suffix of “M” or “G” can be used to signify a value in megabytes or gigabytes respectively. Optional pair slots, maxmem could be used to set amount of hotpluggable memory slots and maximum amount of memory. Note that maxmem must be aligned to the page size.

For example, the following command-line sets the guest startup RAM size to 1GB, creates 3 slots to hotplug additional memory and sets the maximum memory the guest can reach to 4GB:

```
qemu-system-x86_64 -m 1G,slots=3,maxmem=4G
```

If slots and maxmem are not specified, memory hotplug won’t be enabled and the guest startup RAM will never increase.

-mem-path path Allocate guest RAM from a temporarily created file in path.

-mem-prealloc Preallocate memory when using -mem-path.

-k language Use keyboard layout language (for example `fr` for French). This option is only needed where it is not easy to get raw PC keycodes (e.g. on Macs, with some X11 servers or with a VNC or curses display). You don’t normally need to use it on PC/Linux or PC/Windows hosts.

The available layouts are:

ar	de-ch	es	fo	fr-ca	hu	ja	mk	no	pt-br	sv
da	en-gb	et	fr	fr-ch	is	lt	nl	pl	ru	th
de	en-us	fi	fr-be	hr	it	lv	nl-be	pt	sl	tr

The default is `en-us`.

-audio-help Will show the -audiodev equivalent of the currently specified (deprecated) environment variables.

-audiodev [driver=]driver,id=id[,prop[=value][, ...]] Adds a new audio backend driver identified by id. There are global and driver specific properties. Some values can be set differently for input

and output, they're marked with `in|out..` You can set the input's property with `in.prop` and the output's property with `out.prop`. For example:

```
-audiodev alsa,id=example,in.frequency=44110,out.frequency=8000
-audiodev alsa,id=example,out.channels=1 # leaves in.channels unspecified
```

NOTE: parameter validation is known to be incomplete, in many cases specifying an invalid option causes QEMU to print an error message and continue emulation without sound.

Valid global options are:

id=identifier Identifies the audio backend.

timer-period=period Sets the timer period used by the audio subsystem in microseconds. Default is 10000 (10 ms).

in|out.mixing-engine=on|off Use QEMU's mixing engine to mix all streams inside QEMU and convert audio formats when not supported by the backend. When off, fixed-settings must be off too. Note that disabling this option means that the selected backend must support multiple streams and the audio formats used by the virtual cards, otherwise you'll get no sound. It's not recommended to disable this option unless you want to use 5.1 or 7.1 audio, as mixing engine only supports mono and stereo audio. Default is on.

in|out.fixed-settings=on|off Use fixed settings for host audio. When off, it will change based on how the guest opens the sound card. In this case you must not specify frequency, channels or format. Default is on.

in|out.frequency=frequency Specify the frequency to use when using fixed-settings. Default is 44100Hz.

in|out.channels=channels Specify the number of channels to use when using fixed-settings. Default is 2 (stereo).

in|out.format=format Specify the sample format to use when using fixed-settings. Valid values are: s8, s16, s32, u8, u16, u32, f32. Default is s16.

in|out.voices=voices Specify the number of voices to use. Default is 1.

in|out.buffer-length=usecs Sets the size of the buffer in microseconds.

-audiodev none,id=id[,prop[=value][, ...]] Creates a dummy backend that discards all outputs. This backend has no backend specific properties.

-audiodev alsa,id=id[,prop[=value][, ...]] Creates backend using the ALSA. This backend is only available on Linux.

ALSA specific options are:

in|out.dev=device Specify the ALSA device to use for input and/or output. Default is `default`.

in|out.period-length=usecs Sets the period length in microseconds.

in|out.try-poll=on|off Attempt to use poll mode with the device. Default is on.

threshold=threshold Threshold (in microseconds) when playback starts. Default is 0.

-audiodev coreaudio,id=id[,prop[=value][, ...]] Creates a backend using Apple's Core Audio. This backend is only available on Mac OS and only supports playback.

Core Audio specific options are:

in|out.buffer-count=count Sets the count of the buffers.

-audiodev dsound, **id=id**[, **prop**[=**value**][, ...]] Creates a backend using Microsoft's DirectSound. This backend is only available on Windows and only supports playback.

DirectSound specific options are:

latency=usecs Add extra usecs microseconds latency to playback. Default is 10000 (10 ms).

-audiodev oss, **id=id**[, **prop**[=**value**][, ...]] Creates a backend using OSS. This backend is available on most Unix-like systems.

OSS specific options are:

in|out.dev=device Specify the file name of the OSS device to use. Default is `/dev/dsp`.

in|out.buffer-count=count Sets the count of the buffers.

in|out.try-poll=on|off Attempt to use poll mode with the device. Default is on.

try-mmap=on|off Try using memory mapped device access. Default is off.

exclusive=on|off Open the device in exclusive mode (vmix won't work in this case). Default is off.

dsp-policy=policy Sets the timing policy (between 0 and 10, where smaller number means smaller latency but higher CPU usage). Use -1 to use buffer sizes specified by `buffer` and `buffer-count`. This option is ignored if you do not have OSS 4. Default is 5.

-audiodev pa, **id=id**[, **prop**[=**value**][, ...]] Creates a backend using PulseAudio. This backend is available on most systems.

PulseAudio specific options are:

server=server Sets the PulseAudio server to connect to.

in|out.name=sink Use the specified source/sink for recording/playback.

in|out.latency=usecs Desired latency in microseconds. The PulseAudio server will try to honor this value but actual latencies may be lower or higher.

-audiodev sdl, **id=id**[, **prop**[=**value**][, ...]] Creates a backend using SDL. This backend is available on most systems, but you should use your platform's native backend if possible.

SDL specific options are:

in|out.buffer-count=count Sets the count of the buffers.

-audiodev spice, **id=id**[, **prop**[=**value**][, ...]] Creates a backend that sends audio through SPICE. This backend requires `-spice` and automatically selected in that case, so usually you can ignore this option. This backend has no backend specific properties.

-audiodev wav, **id=id**[, **prop**[=**value**][, ...]] Creates a backend that writes audio to a WAV file.

Backend specific options are:

path=path Write recorded audio into the specified file. Default is `qemu.wav`.

-soundhw card1[, **card2**, ...] or **-soundhw all** Enable audio and selected sound hardware. Use 'help' to print all available sound hardware. For example:

```
qemu_system-x86_64 -soundhw sb16,adlib disk.img
qemu_system-x86_64 -soundhw es1370 disk.img
qemu_system-x86_64 -soundhw ac97 disk.img
qemu_system-x86_64 -soundhw hda disk.img
qemu_system-x86_64 -soundhw all disk.img
qemu_system-x86_64 -soundhw help
```

Note that Linux's i810_audio OSS kernel (for AC97) module might require manually specifying clocking.

```
modprobe i810_audio clocking=48000
```

-device driver[,prop=[value][, ...]] Add device driver. `prop=value` sets driver properties. Valid properties depend on the driver. To get help on possible drivers and properties, use `-device help` and `-device driver,help`.

Some drivers are:

-device ipmi-bmc-sim,id=id[,prop=[value][, ...]] Add an IPMI BMC. This is a simulation of a hardware management interface processor that normally sits on a system. It provides a watchdog and the ability to reset and power control the system. You need to connect this to an IPMI interface to make it useful

The IPMI slave address to use for the BMC. The default is 0x20. This address is the BMC's address on the I2C network of management controllers. If you don't know what this means, it is safe to ignore it.

id=id The BMC id for interfaces to use this device.

slave_addr=val Define slave address to use for the BMC. The default is 0x20.

sdrfile=file file containing raw Sensor Data Records (SDR) data. The default is none.

fruareasize=val size of a Field Replaceable Unit (FRU) area. The default is 1024.

frudatafile=file file containing raw Field Replaceable Unit (FRU) inventory data. The default is none.

guid=uuid value for the GUID for the BMC, in standard UUID format. If this is set, get "Get GUID" command to the BMC will return it. Otherwise "Get GUID" will return an error.

-device ipmi-bmc-extern,id=id,chardev=id[,slave_addr=val] Add a connection to an external IPMI BMC simulator. Instead of locally emulating the BMC like the above item, instead connect to an external entity that provides the IPMI services.

A connection is made to an external BMC simulator. If you do this, it is strongly recommended that you use the "reconnect=" chardev option to reconnect to the simulator if the connection is lost. Note that if this is not used carefully, it can be a security issue, as the interface has the ability to send resets, NMIs, and power off the VM. It's best if QEMU makes a connection to an external simulator running on a secure port on localhost, so neither the simulator nor QEMU is exposed to any outside network.

See the "lanserv/README.vm" file in the OpenIPMI library for more details on the external interface.

-device isa-ipmi-kcs,bmc=id[,iport=val][,irq=val] Add a KCS IPMI interface on the ISA bus. This also adds a corresponding ACPI and SMBIOS entries, if appropriate.

bmc=id The BMC to connect to, one of ipmi-bmc-sim or ipmi-bmc-extern above.

iport=val Define the I/O address of the interface. The default is 0xca0 for KCS.

irq=val Define the interrupt to use. The default is 5. To disable interrupts, set this to 0.

-device isa-ipmi-bt,bmc=id[,iport=val][,irq=val] Like the KCS interface, but defines a BT interface. The default port is 0xe4 and the default interrupt is 5.

-device pci-ipmi-kcs,bmc=id Add a KCS IPMI interface on the PCI bus.

bmc=id The BMC to connect to, one of ipmi-bmc-sim or ipmi-bmc-extern above.

-device pci-ipmi-bt,bmc=id Like the KCS interface, but defines a BT interface on the PCI bus.

-name name Sets the name of the guest. This name will be displayed in the SDL window caption. The name will also be used for the VNC server. Also optionally set the top visible process name in Linux. Naming of individual threads can also be enabled on Linux to aid debugging.

-uuid uuid Set system UUID.

1.2.2 Block device options

-fda file

-fdb file Use file as floppy disk 0/1 image (see the *Disk Images* chapter in the System Emulation Users Guide).

-hda file

-hdb file

-hdc file

-hdd file Use file as hard disk 0, 1, 2 or 3 image (see the *Disk Images* chapter in the System Emulation Users Guide).

-cdrom file Use file as CD-ROM image (you cannot use **-hdc** and **-cdrom** at the same time). You can use the host CD-ROM by using `/dev/cdrom` as filename.

-blockdev option[,option[,option[,...]]] Define a new block driver node. Some of the options apply to all block drivers, other options are only accepted for a specific block driver. See below for a list of generic options and options for the most common block drivers.

Options that expect a reference to another node (e.g. `file`) can be given in two ways. Either you specify the node name of an already existing node (`file=node-name`), or you define a new node inline, adding options for the referenced node after a dot (`file.filename=path,file.aio=native`).

A block driver node created with **-blockdev** can be used for a guest device by specifying its node name for the `drive` property in a **-device** argument that defines a block device.

Valid options for any block driver node:

driver Specifies the block driver to use for the given node.

node-name This defines the name of the block driver node by which it will be referenced later. The name must be unique, i.e. it must not match the name of a different block driver node, or (if you use **-drive** as well) the ID of a drive.

If no node name is specified, it is automatically generated. The generated node name is not intended to be predictable and changes between QEMU invocations. For the top level, an explicit node name must be specified.

read-only Open the node read-only. Guest write attempts will fail.

Note that some block drivers support only read-only access, either generally or in certain configurations. In this case, the default value `read-only=off` does not work and the option must be specified explicitly.

auto-read-only If `auto-read-only=on` is set, QEMU may fall back to read-only usage even when `read-only=off` is requested, or even switch between modes as needed, e.g. depending on whether the image file is writable or whether a writing user is attached to the node.

force-share Override the image locking system of QEMU by forcing the node to utilize weaker shared access for permissions where it would normally request exclusive access. When there is the potential for multiple instances to have the same file open (whether this invocation of QEMU is the first or the second instance), both instances must permit shared access for the second instance to succeed at opening the file.

Enabling `force-share=on` requires `read-only=on`.

cache.direct The host page cache can be avoided with `cache.direct=on`. This will attempt to do disk IO directly to the guest's memory. QEMU may still perform an internal copy of the data.

cache.no-flush In case you don't care about data integrity over host failures, you can use `cache.no-flush=on`. This option tells QEMU that it never needs to write any data to the disk but can

instead keep things in cache. If anything goes wrong, like your host losing power, the disk storage getting disconnected accidentally, etc. your image will most probably be rendered unusable.

discard=discard discard is one of “ignore” (or “off”) or “unmap” (or “on”) and controls whether discard (also known as trim or unmap) requests are ignored or passed to the filesystem. Some machine types may not support discard requests.

detect-zeroes=detect-zeroes detect-zeroes is “off”, “on” or “unmap” and enables the automatic conversion of plain zero writes by the OS to driver specific optimized zero write commands. You may even choose “unmap” if discard is set to “unmap” to allow a zero write to be converted to an unmap operation.

Driver-specific options for file This is the protocol-level block driver for accessing regular files.

filename The path to the image file in the local filesystem

aio Specifies the AIO backend (threads/native/io_uring, default: threads)

locking Specifies whether the image file is protected with Linux OFD / POSIX locks. The default is to use the Linux Open File Descriptor API if available, otherwise no lock is applied. (auto/on/off, default: auto)

Example:

```
-blockdev driver=file,node-name=disk,filename=disk.img
```

Driver-specific options for raw This is the image format block driver for raw images. It is usually stacked on top of a protocol level block driver such as file.

file Reference to or definition of the data source block driver node (e.g. a file driver node)

Example 1:

```
-blockdev driver=file,node-name=disk_file,filename=disk.img
-blockdev driver=raw,node-name=disk,file=disk_file
```

Example 2:

```
-blockdev driver=raw,node-name=disk,file.driver=file,file.filename=disk.img
```

Driver-specific options for qcow2 This is the image format block driver for qcow2 images. It is usually stacked on top of a protocol level block driver such as file.

file Reference to or definition of the data source block driver node (e.g. a file driver node)

backing Reference to or definition of the backing file block device (default is taken from the image file). It is allowed to pass null here in order to disable the default backing file.

lazy-refcounts Whether to enable the lazy refcounts feature (on/off; default is taken from the image file)

cache-size The maximum total size of the L2 table and refcount block caches in bytes (default: the sum of l2-cache-size and refcount-cache-size)

l2-cache-size The maximum size of the L2 table cache in bytes (default: if cache-size is not specified - 32M on Linux platforms, and 8M on non-Linux platforms; otherwise, as large as possible within the cache-size, while permitting the requested or the minimal refcount cache size)

refcount-cache-size The maximum size of the refcount block cache in bytes (default: 4 times the cluster size; or if cache-size is specified, the part of it which is not used for the L2 cache)

cache-clean-interval Clean unused entries in the L2 and refcount caches. The interval is in seconds. The default value is 600 on supporting platforms, and 0 on other platforms. Setting it to 0 disables this feature.

pass-discard-request Whether discard requests to the qcow2 device should be forwarded to the data source (on/off; default: on if discard=unmap is specified, off otherwise)

pass-discard-snapshot Whether discard requests for the data source should be issued when a snapshot operation (e.g. deleting a snapshot) frees clusters in the qcow2 file (on/off; default: on)

pass-discard-other Whether discard requests for the data source should be issued on other occasions where a cluster gets freed (on/off; default: off)

overlap-check Which overlap checks to perform for writes to the image (none/constant/cached/all; default: cached). For details or finer granularity control refer to the QAPI documentation of `blockdev-add`.

Example 1:

```
-blockdev driver=file,node-name=my_file,filename=/tmp/disk.qcow2
-blockdev driver=qcow2,node-name=hda,file=my_file,overlap-check=none,cache-
↪size=16777216
```

Example 2:

```
-blockdev driver=qcow2,node-name=disk,file.driver=http,file.filename=http://
↪example.com/image.qcow2
```

Driver-specific options for other drivers Please refer to the QAPI documentation of the `blockdev-add QMP` command.

-drive option[,option[,option[,...]]] Define a new drive. This includes creating a block driver node (the backend) as well as a guest device, and is mostly a shortcut for defining the corresponding `-blockdev` and `-device` options.

`-drive` accepts all options that are accepted by `-blockdev`. In addition, it knows the following options:

file=file This option defines which disk image (see the [Disk Images](#) chapter in the System Emulation Users Guide) to use with this drive. If the filename contains comma, you must double it (for instance, “file=my,file” to use file “my,file”).

Special files such as iSCSI devices can be specified using protocol specific URLs. See the section for “Device URL Syntax” for more information.

if=interface This option defines on which type on interface the drive is connected. Available types are: ide, scsi, sd, mtd, floppy, pflash, virtio, none.

bus=bus,unit=unit These options define where is connected the drive by defining the bus number and the unit id.

index=index This option defines where is connected the drive by using an index in the list of available connectors of a given interface type.

media=media This option defines the type of the media: disk or cdrom.

snapshot=snapshot snapshot is “on” or “off” and controls snapshot mode for the given drive (see `-snapshot`).

cache=cache cache is “none”, “writeback”, “unsafe”, “directsync” or “writethrough” and controls how the host cache is used to access block data. This is a shortcut that sets the `cache.direct` and `cache.no-flush` options (as in `-blockdev`), and additionally `cache.writeback`, which provides a de-

fault for the `write-cache` option of block guest devices (as in `-device`). The modes correspond to the following settings:

	cache.writeback	cache.direct	cache.no-flush
writeback	on	off	off
none	on	on	off
writethrough	off	off	off
directsync	off	on	off
unsafe	on	off	on

The default mode is `cache=writeback`.

aio=aio `aio` is “threads”, “native”, or “io_uring” and selects between pthread based disk I/O, native Linux AIO, or Linux io_uring API.

format=format Specify which disk format will be used rather than detecting the format. Can be used to specify `format=raw` to avoid interpreting an untrusted format header.

werror=action, rerror=action Specify which action to take on write and read errors. Valid actions are: “ignore” (ignore the error and try to continue), “stop” (pause QEMU), “report” (report the error to the guest), “enospc” (pause QEMU only if the host disk is full; report the error to the guest otherwise). The default setting is `werror=enospc` and `rerror=report`.

copy-on-read=copy-on-read `copy-on-read` is “on” or “off” and enables whether to copy read backing file sectors into the image file.

bps=b, bps_rd=r, bps_wr=w Specify bandwidth throttling limits in bytes per second, either for all request types or for reads or writes only. Small values can lead to timeouts or hangs inside the guest. A safe minimum for disks is 2 MB/s.

bps_max=bm, bps_rd_max=rm, bps_wr_max=wm Specify bursts in bytes per second, either for all request types or for reads or writes only. Bursts allow the guest I/O to spike above the limit temporarily.

iops=i, iops_rd=r, iops_wr=w Specify request rate limits in requests per second, either for all request types or for reads or writes only.

iops_max=bm, iops_rd_max=rm, iops_wr_max=wm Specify bursts in requests per second, either for all request types or for reads or writes only. Bursts allow the guest I/O to spike above the limit temporarily.

iops_size=is Let every `is` bytes of a request count as a new request for iops throttling purposes. Use this option to prevent guests from circumventing iops limits by sending fewer but larger requests.

group=g Join a throttling quota group with given name `g`. All drives that are members of the same group are accounted for together. Use this option to prevent guests from circumventing throttling limits by using many small disks instead of a single larger disk.

By default, the `cache.writeback=on` mode is used. It will report data writes as completed as soon as the data is present in the host page cache. This is safe as long as your guest OS makes sure to correctly flush disk caches where needed. If your guest OS does not handle volatile disk write caches correctly and your host crashes or loses power, then the guest may experience data corruption.

For such guests, you should consider using `cache.writeback=off`. This means that the host page cache will be used to read and write data, but write notification will be sent to the guest only after QEMU has made sure to flush each write to the disk. Be aware that this has a major impact on performance.

When using the `-snapshot` option, unsafe caching is always used.

Copy-on-read avoids accessing the same backing file sectors repeatedly and is useful when the backing file is over a slow network. By default `copy-on-read` is off.

Instead of `-cdrom` you can use:

```
qemu-system-x86_64 -drive file=file,index=2,media=cdrom
```

Instead of `-hda`, `-hdb`, `-hdc`, `-hdd`, you can use:

```
qemu-system-x86_64 -drive file=file,index=0,media=disk
qemu-system-x86_64 -drive file=file,index=1,media=disk
qemu-system-x86_64 -drive file=file,index=2,media=disk
qemu-system-x86_64 -drive file=file,index=3,media=disk
```

You can open an image using pre-opened file descriptors from an fd set:

```
qemu-system-x86_64 \
-add-fd fd=3,set=2,opaque="rdwr:/path/to/file" \
-add-fd fd=4,set=2,opaque="ronly:/path/to/file" \
-drive file=/dev/fdset/2,index=0,media=disk
```

You can connect a CDROM to the slave of ide0:

```
qemu_system-x86_64 -drive file=file,if=ide,index=1,media=cdrom
```

If you don't specify the "file=" argument, you define an empty drive:

```
qemu_system-x86_64 -drive if=ide,index=1,media=cdrom
```

Instead of `-fda`, `-fdb`, you can use:

```
qemu_system-x86_64 -drive file=file,index=0,if=floppy
qemu_system-x86_64 -drive file=file,index=1,if=floppy
```

By default, interface is "ide" and index is automatically incremented:

```
qemu_system-x86_64 -drive file=a -drive file=b"
```

is interpreted like:

```
qemu_system-x86_64 -hda a -hdb b
```

-mtdblock file Use file as on-board Flash memory image.

-sd file Use file as SecureDigital card image.

-pflash file Use file as a parallel flash image.

-snapshot Write to temporary files instead of disk image files. In this case, the raw disk image you use is not written back. You can however force the write back by pressing C-a s (see the *Disk Images* chapter in the System Emulation Users Guide).

-fsdev local,id=id,path=path,security_model=security_model [,writeout=writeout][,readonly=on]

-fsdev proxy,id=id,socket=socket [,writeout=writeout][,readonly=on]

-fsdev proxy,id=id,sock_fd=sock_fd[,writeout=writeout][,readonly=on]

-fsdev synth,id=id[,readonly=on] Define a new file system device. Valid options are:

local Accesses to the filesystem are done by QEMU.

proxy Accesses to the filesystem are done by virtfs-proxy-helper(1).

synth Synthetic filesystem, only used by QTests.

id=id Specifies identifier for this device.

path=path Specifies the export path for the file system device. Files under this path will be available to the 9p client on the guest.

security_model=security_model Specifies the security model to be used for this export path. Supported security models are “passthrough”, “mapped-xattr”, “mapped-file” and “none”. In “passthrough” security model, files are stored using the same credentials as they are created on the guest. This requires QEMU to run as root. In “mapped-xattr” security model, some of the file attributes like uid, gid, mode bits and link target are stored as file attributes. For “mapped-file” these attributes are stored in the hidden .virtfs_metadata directory. Directories exported by this security model cannot interact with other unix tools. “none” security model is same as passthrough except the server won’t report failures if it fails to set file attributes like ownership. Security model is mandatory only for local fsdriver. Other fsdrivers (like proxy) don’t take security model as a parameter.

writeout=writeout This is an optional argument. The only supported value is “immediate”. This means that host page cache will be used to read and write data but write notification will be sent to the guest only when the data has been reported as written by the storage subsystem.

readonly=on Enables exporting 9p share as a readonly mount for guests. By default read-write access is given.

socket=socket Enables proxy filesystem driver to use passed socket file for communicating with virtfs-proxy-helper(1).

sock_fd=sock_fd Enables proxy filesystem driver to use passed socket descriptor for communicating with virtfs-proxy-helper(1). Usually a helper like libvirt will create socketpair and pass one of the fds as sock_fd.

fmode=fmode Specifies the default mode for newly created files on the host. Works only with security models “mapped-xattr” and “mapped-file”.

dmode=dmode Specifies the default mode for newly created directories on the host. Works only with security models “mapped-xattr” and “mapped-file”.

throttling.bps-total=b, throttling.bps-read=r, throttling.bps-write=w Specify bandwidth throttling limits in bytes per second, either for all request types or for reads or writes only.

throttling.bps-total-max=bm, bps-read-max=rm, bps-write-max=wm Specify bursts in bytes per second, either for all request types or for reads or writes only. Bursts allow the guest I/O to spike above the limit temporarily.

throttling.iops-total=i, throttling.iops-read=r, throttling.iops-write=w Specify request rate limits in requests per second, either for all request types or for reads or writes only.

throttling.iops-total-max=im, throttling.iops-read-max=irm, throttling.iops-write-max=iw Specify bursts in requests per second, either for all request types or for reads or writes only. Bursts allow the guest I/O to spike above the limit temporarily.

throttling.iops-size=is Let every is bytes of a request count as a new request for iops throttling purposes.

-fsdev option is used along with -device driver “virtio-9p-...”.

-device virtio-9p-type, fsdev=id, mount_tag=mount_tag Options for virtio-9p-... driver are:

type Specifies the variant to be used. Supported values are “pci”, “ccw” or “device”, depending on the machine type.

fsdev=id Specifies the id value specified along with -fsdev option.

mount_tag=mount_tag Specifies the tag name to be used by the guest to mount this export point.

-virtfs local, path=path, mount_tag=mount_tag , security_model=security_model [, writeout=writeout]

-virtfs proxy, socket=socket, mount_tag=mount_tag [, writeout=writeout] [, readonly=on]

-virtfs proxy,sock_fd=sock_fd,mount_tag=mount_tag [,writeout=writeout][,readonly=on]

-virtfs synth,mount_tag=mount_tag Define a new virtual filesystem device and expose it to the guest using a virtio-9p-device (a.k.a. 9pfs), which essentially means that a certain directory on host is made directly accessible by guest as a pass-through file system by using the 9P network protocol for communication between host and guests, if desired even accessible, shared by several guests simultaneously.

Note that **-virtfs** is actually just a convenience shortcut for its generalized form **-fsdev -device virtio-9p-pci**.

The general form of pass-through file system options are:

local Accesses to the filesystem are done by QEMU.

proxy Accesses to the filesystem are done by virtfs-proxy-helper(1).

synth Synthetic filesystem, only used by QTests.

id=id Specifies identifier for the filesystem device

path=path Specifies the export path for the file system device. Files under this path will be available to the 9p client on the guest.

security_model=security_model Specifies the security model to be used for this export path. Supported security models are “passthrough”, “mapped-xattr”, “mapped-file” and “none”. In “passthrough” security model, files are stored using the same credentials as they are created on the guest. This requires QEMU to run as root. In “mapped-xattr” security model, some of the file attributes like uid, gid, mode bits and link target are stored as file attributes. For “mapped-file” these attributes are stored in the hidden .virtfs_metadata directory. Directories exported by this security model cannot interact with other unix tools. “none” security model is same as passthrough except the server won’t report failures if it fails to set file attributes like ownership. Security model is mandatory only for local fsdriver. Other fsdrivers (like proxy) don’t take security model as a parameter.

writeout=writeout This is an optional argument. The only supported value is “immediate”. This means that host page cache will be used to read and write data but write notification will be sent to the guest only when the data has been reported as written by the storage subsystem.

readonly=on Enables exporting 9p share as a readonly mount for guests. By default read-write access is given.

socket=socket Enables proxy filesystem driver to use passed socket file for communicating with virtfs-proxy-helper(1). Usually a helper like libvirt will create socketpair and pass one of the fds as sock_fd.

sock_fd Enables proxy filesystem driver to use passed ‘sock_fd’ as the socket descriptor for interfacing with virtfs-proxy-helper(1).

fmode=fmode Specifies the default mode for newly created files on the host. Works only with security models “mapped-xattr” and “mapped-file”.

dmode=dmode Specifies the default mode for newly created directories on the host. Works only with security models “mapped-xattr” and “mapped-file”.

mount_tag=mount_tag Specifies the tag name to be used by the guest to mount this export point.

multidevs=multidevs Specifies how to deal with multiple devices being shared with a 9p export. Supported behaviours are either “remap”, “forbid” or “warn”. The latter is the default behaviour on which virtfs 9p expects only one device to be shared with the same export, and if more than one device is shared and accessed via the same 9p export then only a warning message is logged (once) by qemu on host side. In order to avoid file ID collisions on guest you should either create a separate virtfs export for each device to be shared with guests (recommended way) or you might use “remap” instead which allows you to share multiple devices with only one export instead, which is achieved by remapping the original inode numbers

from host to guest in a way that would prevent such collisions. Remapping inodes in such use cases is required because the original device IDs from host are never passed and exposed on guest. Instead all files of an export shared with virtfs always share the same device id on guest. So two files with identical inode numbers but from actually different devices on host would otherwise cause a file ID collision and hence potential misbehaviours on guest. “forbid” on the other hand assumes like “warn” that only one device is shared by the same export, however it will not only log a warning message but also deny access to additional devices on guest. Note though that “forbid” does currently not block all possible file access operations (e.g. `readdir()` would still return entries from other devices).

-iscsi Configure iSCSI session parameters.

1.2.3 USB options

-usb Enable USB emulation on machine types with an on-board USB host controller (if not enabled by default). Note that on-board USB host controllers may not support USB 3.0. In this case `-device qemu-xhci` can be used instead on machines with PCI.

-usbdevice devname Add the USB device devname. Note that this option is deprecated, please use `-device usb-...` instead. See the chapter about *Connecting USB devices* in the System Emulation Users Guide.

mouse Virtual Mouse. This will override the PS/2 mouse emulation when activated.

tablet Pointer device that uses absolute coordinates (like a touchscreen). This means QEMU is able to report the mouse position without having to grab the mouse. Also overrides the PS/2 mouse emulation when activated.

braille Braille device. This will use BrlAPI to display the braille output on a real or fake device.

1.2.4 Display options

-display type Select type of display to use. This option is a replacement for the old style `-sdl/-curses/...` options. Use `-display help` to list the available display types. Valid values for type are

sdl Display video output via SDL (usually in a separate graphics window; see the SDL documentation for other possibilities).

curses Display video output via curses. For graphics device models which support a text mode, QEMU can display this output using a curses/ncurses interface. Nothing is displayed when the graphics device is in graphical mode or if the graphics device does not support a text mode. Generally only the VGA device models support text mode. The font charset used by the guest can be specified with the `charset` option, for example `charset=CP850` for IBM CP850 encoding. The default is CP437.

none Do not display video output. The guest will still see an emulated graphics card, but its output will not be displayed to the QEMU user. This option differs from the `-nographic` option in that it only affects what is done with video output; `-nographic` also changes the destination of the serial and parallel port data.

gtk Display video output in a GTK window. This interface provides drop-down menus and other UI elements to configure and control the VM during runtime.

vnc Start a VNC server on display <arg>

egl-headless Offload all OpenGL operations to a local DRI device. For any graphical display, this display needs to be paired with either VNC or SPICE displays.

spice-app Start QEMU as a Spice server and launch the default Spice client application. The Spice server will redirect the serial consoles and QEMU monitors. (Since 4.0)

-nographic Normally, if QEMU is compiled with graphical window support, it displays output such as guest graphics, guest console, and the QEMU monitor in a window. With this option, you can totally disable graphical output so that QEMU is a simple command line application. The emulated serial port is redirected on the console and muxed with the monitor (unless redirected elsewhere explicitly). Therefore, you can still use QEMU to debug a Linux kernel with a serial console. Use C-a h for help on switching between the console and monitor.

-curses Normally, if QEMU is compiled with graphical window support, it displays output such as guest graphics, guest console, and the QEMU monitor in a window. With this option, QEMU can display the VGA output when in text mode using a curses/ncurses interface. Nothing is displayed in graphical mode.

-alt-grab Use Ctrl-Alt-Shift to grab mouse (instead of Ctrl-Alt). Note that this also affects the special keys (for fullscreen, monitor-mode switching, etc).

-ctrl-grab Use Right-Ctrl to grab mouse (instead of Ctrl-Alt). Note that this also affects the special keys (for fullscreen, monitor-mode switching, etc).

-no-quit Disable SDL window close capability.

-sdl Enable SDL.

-spice option[,option[,...]] Enable the spice remote desktop protocol. Valid options are

port=<nr> Set the TCP port spice is listening on for plaintext channels.

addr=<addr> Set the IP address spice is listening on. Default is any address.

ipv4; ipv6; unix Force using the specified IP version.

password=<secret> Set the password you need to authenticate.

sasl Require that the client use SASL to authenticate with the spice. The exact choice of authentication method used is controlled from the system / user's SASL configuration file for the 'qemu' service. This is typically found in /etc/sasl2/qemu.conf. If running QEMU as an unprivileged user, an environment variable SASL_CONF_PATH can be used to make it search alternate locations for the service config. While some SASL auth methods can also provide data encryption (eg GSSAPI), it is recommended that SASL always be combined with the 'tls' and 'x509' settings to enable use of SSL and server certificates. This ensures a data encryption preventing compromise of authentication credentials.

disable-ticketing Allow client connects without authentication.

disable-copy-paste Disable copy paste between the client and the guest.

disable-agent-file-xfer Disable spice-vdagent based file-xfer between the client and the guest.

tls-port=<nr> Set the TCP port spice is listening on for encrypted channels.

x509-dir=<dir> Set the x509 file directory. Expects same filenames as -vnc \$display,x509=\$dir

x509-key-file=<file>; x509-key-password=<file>; x509-cert-file=<file>; x509-cacert-file=<file>
The x509 file names can also be configured individually.

tls-ciphers=<list> Specify which ciphers to use.

tls-channel=[main|display|cursor|inputs|record|playback]; plaintext-channel=[main|display|cursor|inputs|record|playback]
Force specific channel to be used with or without TLS encryption. The options can be specified multiple times to configure multiple channels. The special name "default" can be used to set the default mode. For channels which are not explicitly forced into one mode the spice client is allowed to pick tls/plaintext as he pleases.

image-compression=[auto_glz|auto_lz|quic|glz|lz|off] Configure image compression (lossless). Default is auto_glz.

jpeg-wan-compression=[auto|never|always]; zlib-glz-wan-compression=[auto|never|always]
Configure wan image compression (lossy for slow links). Default is auto.

streaming-video=[off|all|filter] Configure video stream detection. Default is off.

agent-mouse=[on|off] Enable/disable passing mouse events via vdaagent. Default is on.

playback-compression=[on|off] Enable/disable audio stream compression (using celt 0.5.1). Default is on.

seamless-migration=[on|off] Enable/disable spice seamless migration. Default is off.

gl=[on|off] Enable/disable OpenGL context. Default is off.

rendernode=<file> DRM render node for OpenGL rendering. If not specified, it will pick the first available. (Since 2.9)

-portrait Rotate graphical output 90 deg left (only PXA LCD).

-rotate deg Rotate graphical output some deg left (only PXA LCD).

-vga type Select type of VGA card to emulate. Valid values for type are

cirrus Cirrus Logic GD5446 Video card. All Windows versions starting from Windows 95 should recognize and use this graphic card. For optimal performances, use 16 bit color depth in the guest and the host OS. (This card was the default before QEMU 2.2)

std Standard VGA card with Bochs VBE extensions. If your guest OS supports the VESA 2.0 VBE extensions (e.g. Windows XP) and if you want to use high resolution modes ($\geq 1280 \times 1024 \times 16$) then you should use this option. (This card is the default since QEMU 2.2)

vmware VMWare SVGA-II compatible adapter. Use it if you have sufficiently recent XFree86/XOrg server or Windows guest with a driver for this card.

qxl QXL paravirtual graphic card. It is VGA compatible (including VESA 2.0 VBE support). Works best with qxl guest drivers installed though. Recommended choice when using the spice protocol.

tcx (sun4m only) Sun TCX framebuffer. This is the default framebuffer for sun4m machines and offers both 8-bit and 24-bit colour depths at a fixed resolution of 1024×768 .

cg3 (sun4m only) Sun cgthree framebuffer. This is a simple 8-bit framebuffer for sun4m machines available in both 1024×768 (OpenBIOS) and 1152×900 (OBP) resolutions aimed at people wishing to run older Solaris versions.

virtio Virtio VGA card.

none Disable VGA card.

-full-screen Start in full screen.

-g widthxheight[xdepth] Set the initial graphical resolution and depth (PPC, SPARC only).

For PPC the default is $800 \times 600 \times 32$.

For SPARC with the TCX graphics device, the default is $1024 \times 768 \times 8$ with the option of $1024 \times 768 \times 24$. For cgthree, the default is $1024 \times 768 \times 8$ with the option of $1152 \times 900 \times 8$ for people who wish to use OBP.

-vnc display[,option[,option[,...]]] Normally, if QEMU is compiled with graphical window support, it displays output such as guest graphics, guest console, and the QEMU monitor in a window. With this option, you can have QEMU listen on VNC display and redirect the VGA display over the VNC session. It is very useful to enable the usb tablet device when using this option (option **-device usb-tablet**). When using the VNC display, you must use the **-k** parameter to set the keyboard layout if you are not using en-us. Valid syntax for the display is

to=L With this option, QEMU will try next available VNC displays, until the number L, if the originally defined “-vnc display” is not available, e.g. port 5900+display is already used by another application. By default, to=0.

host:d TCP connections will only be allowed from host on display d. By convention the TCP port is 5900+d. Optionally, host can be omitted in which case the server will accept connections from any host.

unix:path Connections will be allowed over UNIX domain sockets where path is the location of a unix socket to listen for connections on.

none VNC is initialized but not started. The monitor `change` command can be used to later start the VNC server.

Following the display value there may be one or more option flags separated by commas. Valid options are

reverse Connect to a listening VNC client via a “reverse” connection. The client is specified by the display. For reverse network connections (host:d, “reverse”), the d argument is a TCP port number, not a display number.

websocket Opens an additional TCP listening port dedicated to VNC Websocket connections. If a bare websocket option is given, the Websocket port is 5700+display. An alternative port can be specified with the syntax `websocket=port`.

If host is specified connections will only be allowed from this host. It is possible to control the websocket listen address independently, using the syntax `websocket=host:port`.

If no TLS credentials are provided, the websocket connection runs in unencrypted mode. If TLS credentials are provided, the websocket connection requires encrypted client connections.

password Require that password based authentication is used for client connections.

The password must be set separately using the `set_password` command in the *QEMU Monitor*. The syntax to change your password is: `set_password <protocol> <password>` where <protocol> could be either “vnc” or “spice”.

If you would like to change <protocol> password expiration, you should use `expire_password <protocol> <expiration-time>` where expiration time could be one of the following options: now, never, +seconds or UNIX time of expiration, e.g. +60 to make password expire in 60 seconds, or 1335196800 to make password expire on “Mon Apr 23 12:00:00 EDT 2012” (UNIX time for this date and time).

You can also use keywords “now” or “never” for the expiration time to allow <protocol> password to expire immediately or never expire.

tls-creds=ID Provides the ID of a set of TLS credentials to use to secure the VNC server. They will apply to both the normal VNC server socket and the websocket socket (if enabled). Setting TLS credentials will cause the VNC server socket to enable the VeNCrypt auth mechanism. The credentials should have been previously created using the `-object tls-creds` argument.

tls-authz=ID Provides the ID of the QAuthZ authorization object against which the client’s x509 distinguished name will be validated. This object is only resolved at time of use, so can be deleted and recreated on the fly while the VNC server is active. If missing, it will default to denying access.

sasl Require that the client use SASL to authenticate with the VNC server. The exact choice of authentication method used is controlled from the system / user’s SASL configuration file for the ‘qemu’ service. This is typically found in `/etc/sasl2/qemu.conf`. If running QEMU as an unprivileged user, an environment variable `SASL_CONF_PATH` can be used to make it search alternate locations for the service config. While some SASL auth methods can also provide data encryption (eg GSSAPI), it is recommended that SASL always be combined with the ‘tls’ and ‘x509’ settings to enable use of SSL and server certificates. This ensures a data encryption preventing compromise of authentication credentials. See the *VNC security* section in the System Emulation Users Guide for details on using SASL authentication.

sasl-authz=ID Provides the ID of the QAuthZ authorization object against which the client’s SASL username will be validated. This object is only resolved at time of use, so can be deleted and recreated on the fly while the VNC server is active. If missing, it will default to denying access.

acl Legacy method for enabling authorization of clients against the x509 distinguished name and SASL username. It results in the creation of two `authz-list` objects with IDs of `vnc.username` and `vnc.x509dname`. The rules for these objects must be configured with the HMP ACL commands.

This option is deprecated and should no longer be used. The new `sasl-authz` and `tls-authz` options are a replacement.

lossy Enable lossy compression methods (gradient, JPEG, ...). If this option is set, VNC client may receive lossy framebuffer updates depending on its encoding settings. Enabling this option can save a lot of bandwidth at the expense of quality.

non-adaptive Disable adaptive encodings. Adaptive encodings are enabled by default. An adaptive encoding will try to detect frequently updated screen regions, and send updates in these regions using a lossy encoding (like JPEG). This can be really helpful to save bandwidth when playing videos. Disabling adaptive encodings restores the original static behavior of encodings like Tight.

share=[allow-exclusive|force-shared|ignore] Set display sharing policy. ‘allow-exclusive’ allows clients to ask for exclusive access. As suggested by the rfb spec this is implemented by dropping other connections. Connecting multiple clients in parallel requires all clients asking for a shared session (`vncviewer: -shared` switch). This is the default. ‘force-shared’ disables exclusive client access. Useful for shared desktop sessions, where you don’t want someone forgetting specify `-shared` disconnect everybody else. ‘ignore’ completely ignores the shared flag and allows everybody connect unconditionally. Doesn’t conform to the rfb spec but is traditional QEMU behavior.

key-delay-ms Set keyboard delay, for key down and key up events, in milliseconds. Default is 10. Keyboards are low-bandwidth devices, so this slowdown can help the device and guest to keep up and not lose events in case events are arriving in bulk. Possible causes for the latter are flaky network connections, or scripts for automated testing.

audiodev=audiodev Use the specified audiodev when the VNC client requests audio transmission. When not using an `-audiodev` argument, this option must be omitted, otherwise it must be present and specify a valid audiodev.

power-control Permit the remote client to issue shutdown, reboot or reset power control requests.

1.2.5 i386 target only

-win2k-hack Use it when installing Windows 2000 to avoid a disk full bug. After Windows 2000 is installed, you no longer need this option (this option slows down the IDE transfers).

-no-fd-bootchk Disable boot signature checking for floppy disks in BIOS. May be needed to boot from old floppy disks.

-no-acpi Disable ACPI (Advanced Configuration and Power Interface) support. Use it if your guest OS complains about ACPI problems (PC target machine only).

-no-hpet Disable HPET support.

-acpitable [sig=str] [, rev=n] [, oem_id=str] [, oem_table_id=str] [, oem_rev=n] [,asl_compiler_id=...]

Add ACPI table with specified header fields and context from specified files. For `file=`, take whole ACPI table from the specified files, including all ACPI headers (possible overridden by other options). For `data=`, only data portion of the table is used, all header information is specified in the command line. If a SLIC table is supplied to QEMU, then the SLIC’s `oem_id` and `oem_table_id` fields will override the same in the RSDT and the FADT (a.k.a. FACP), in order to ensure the field matches required by the Microsoft SLIC spec and the ACPI spec.

-smbios file=binary Load SMBIOS entry from binary file.

-smbios type=0 [, vendor=str] [, version=str] [, date=str] [, release=%d.%d] [, uefi=on|off]
Specify SMBIOS type 0 fields

-smbios type=1[,**manufacturer=**str][,**product=**str][,**version=**str][,**serial=**str][,**uuid=**uuid][,**sku=**str]
Specify SMBIOS type 1 fields

-smbios type=2[,**manufacturer=**str][,**product=**str][,**version=**str][,**serial=**str][,**asset=**str][,**location=**str]
Specify SMBIOS type 2 fields

-smbios type=3[,**manufacturer=**str][,**version=**str][,**serial=**str][,**asset=**str][,**sku=**str]
Specify SMBIOS type 3 fields

-smbios type=4[,**sock_pfx=**str][,**manufacturer=**str][,**version=**str][,**serial=**str][,**asset=**str][,**part=**str]
Specify SMBIOS type 4 fields

-smbios type=11[,**value=**str][,**path=**filename] Specify SMBIOS type 11 fields

This argument can be repeated multiple times, and values are added in the order they are parsed. Applications intending to use OEM strings data are encouraged to use their application name as a prefix for the value string. This facilitates passing information for multiple applications concurrently.

The `value=`str syntax provides the string data inline, while the `path=`filename syntax loads data from a file on disk. Note that the file is not permitted to contain any NUL bytes.

Both the `value` and `path` options can be repeated multiple times and will be added to the SMBIOS table in the order in which they appear.

Note that on the x86 architecture, the total size of all SMBIOS tables is limited to 65535 bytes. Thus the OEM strings data is not suitable for passing large amounts of data into the guest. Instead it should be used as an indicator to inform the guest where to locate the real data set, for example, by specifying the serial ID of a block device.

An example passing three strings is

```
-smbios type=11,value=cloud-init:ds=nocloud-net;s=http://10.10.0.1:8000/,\  
value=anaconda:method=http://dl.fedoraproject.org/pub/  
↪fedora/linux/releases/25/x86_64/os,\  
path=/some/file/with/oemstringsdata.txt
```

In the guest OS this is visible with the `dmidecode` command

```
$ dmidecode -t 11  
Handle 0x0E00, DMI type 11, 5 bytes  
OEM Strings  
String 1: cloud-init:ds=nocloud-net;s=http://10.10.0.1:8000/  
String 2: anaconda:method=http://dl.fedoraproject.org/pub/fedora/  
↪linux/releases/25/x86_64/os  
String 3: myapp:some extra data
```

-smbios type=17[,**loc_pfx=**str][,**bank=**str][,**manufacturer=**str][,**serial=**str][,**asset=**str][,**part=**str]
Specify SMBIOS type 17 fields

1.2.6 Network options

-nic [tap|bridge|user|l2tpv3|vde|netmap|vhost-user|socket][,**...**][,**mac=**macaddr][,**model=**mn]

This option is a shortcut for configuring both the on-board (default) guest NIC hardware and the host network backend in one go. The host backend options are the same as with the corresponding `-netdev` options below. The guest NIC model can be set with `model=modelname`. Use `model=help` to list the available device types. The hardware MAC address can be set with `mac=macaddr`.

The following two example do exactly the same, to show how `-nic` can be used to shorten the command line length:

```
qemu-system-x86_64 -netdev user,id=n1,ipv6=off -device e1000,netdev=n1,
↪mac=52:54:98:76:54:32
qemu-system-x86_64 -nic user,ipv6=off,model=e1000,mac=52:54:98:76:54:32
```

-nic none Indicate that no network devices should be configured. It is used to override the default configuration (default NIC with “user” host network backend) which is activated if no other networking options are provided.

-netdev user,id=id[,option][,option][, ...] Configure user mode host network backend which requires no administrator privilege to run. Valid options are:

id=id Assign symbolic name for use in monitor commands.

ipv4=on|off and ipv6=on|off Specify that either IPv4 or IPv6 must be enabled. If neither is specified both protocols are enabled.

net=addr[/mask] Set IP network address the guest will see. Optionally specify the netmask, either in the form a.b.c.d or as number of valid top-most bits. Default is 10.0.2.0/24.

host=addr Specify the guest-visible address of the host. Default is the 2nd IP in the guest network, i.e. x.x.x.2.

ipv6-net=addr[/int] Set IPv6 network address the guest will see (default is fec0::/64). The network prefix is given in the usual hexadecimal IPv6 address notation. The prefix size is optional, and is given as the number of valid top-most bits (default is 64).

ipv6-host=addr Specify the guest-visible IPv6 address of the host. Default is the 2nd IPv6 in the guest network, i.e. xxxx::2.

restrict=on|off If this option is enabled, the guest will be isolated, i.e. it will not be able to contact the host and no guest IP packets will be routed over the host to the outside. This option does not affect any explicitly set forwarding rules.

hostname=name Specifies the client hostname reported by the built-in DHCP server.

dhcpstart=addr Specify the first of the 16 IPs the built-in DHCP server can assign. Default is the 15th to 31st IP in the guest network, i.e. x.x.x.15 to x.x.x.31.

dns=addr Specify the guest-visible address of the virtual nameserver. The address must be different from the host address. Default is the 3rd IP in the guest network, i.e. x.x.x.3.

ipv6-dns=addr Specify the guest-visible address of the IPv6 virtual nameserver. The address must be different from the host address. Default is the 3rd IP in the guest network, i.e. xxxx::3.

dnssearch=domain Provides an entry for the domain-search list sent by the built-in DHCP server. More than one domain suffix can be transmitted by specifying this option multiple times. If supported, this will cause the guest to automatically try to append the given domain suffix(es) in case a domain name can not be resolved.

Example:

```
qemu-system-x86_64 -nic user,dnssearch=mgmt.example.org,
↪dnssearch=example.org
```

domainname=domain Specifies the client domain name reported by the built-in DHCP server.

tftp=dir When using the user mode network stack, activate a built-in TFTP server. The files in dir will be exposed as the root of a TFTP server. The TFTP client on the guest must be configured in binary mode (use the command `bin` of the Unix TFTP client).

tftp-server-name=name In BOOTP reply, broadcast name as the “TFTP server name” (RFC2132 option 66). This can be used to advise the guest to load boot files or configurations from a different server than the host address.

bootfile=file When using the user mode network stack, broadcast file as the BOOTP filename. In conjunction with `tftp`, this can be used to network boot a guest from a local directory.

Example (using `pxelinux`):

```
qemu-system-x86_64 -hda linux.img -boot n -device e1000,netdev=n1 \
    -netdev user,id=n1,tftp=/path/to/tftp/files,bootfile=/pxelinux.0
```

smb=dir[,smbserver=addr] When using the user mode network stack, activate a built-in SMB server so that Windows OSes can access to the host files in `dir` transparently. The IP address of the SMB server can be set to `addr`. By default the 4th IP in the guest network is used, i.e. `x.x.x.4`.

In the guest Windows OS, the line:

```
10.0.2.4 smbserver
```

must be added in the file `C:\WINDOWS\LMHOSTS` (for windows 9x/Me) or `C:\WINNT\SYSTEM32\DRIVERS\ETC\LMHOSTS` (Windows NT/2000).

Then `dir` can be accessed in `\\smbserver\qemu`.

Note that a SAMBA server must be installed on the host OS.

hostfwd=[tcp|udp]:[hostaddr]:hostport-[guestaddr]:guestport Redirect incoming TCP or UDP connections to the host port `hostport` to the guest IP address `guestaddr` on guest port `guestport`. If `guestaddr` is not specified, its value is `x.x.x.15` (default first address given by the built-in DHCP server). By specifying `hostaddr`, the rule can be bound to a specific host interface. If no connection type is set, TCP is used. This option can be given multiple times.

For example, to redirect host X11 connection from screen 1 to guest screen 0, use the following:

```
# on the host
qemu-system-x86_64 -nic user,hostfwd=tcp:127.0.0.1:6001-:6000
# this host xterm should open in the guest X11 server
xterm -display :1
```

To redirect telnet connections from host port 5555 to telnet port on the guest, use the following:

```
# on the host
qemu-system-x86_64 -nic user,hostfwd=tcp::5555-:23
telnet localhost 5555
```

Then when you use on the host `telnet localhost 5555`, you connect to the guest telnet server.

guestfwd=[tcp]:server:port-dev; guestfwd=[tcp]:server:port-cmd:command

Forward guest TCP connections to the IP address `server` on port `port` to the character device `dev` or to a program executed by `cmd:command` which gets spawned for each connection. This option can be given multiple times.

You can either use a chardev directly and have that one used throughout QEMU's lifetime, like in the following example:

```
# open 10.10.1.1:4321 on bootup, connect 10.0.2.100:1234 to it_
↪ whenever
# the guest accesses it
qemu-system-x86_64 -nic user,guestfwd=tcp:10.0.2.100:1234-tcp:10.10.1.1:4321
```

Or you can execute a command on every TCP connection established by the guest, so that QEMU behaves similar to an `inetd` process for that virtual server:

```
# call "netcat 10.10.1.1 4321" on every TCP connection to 10.0.2.
↪ 100:1234
```

```
# and connect the TCP stream to its stdin/stdout
qemu-system-x86_64 -nic 'user,id=n1,guestfwd=tcp:10.0.2.100:1234-
↳cmd:netcat 10.10.1.1 4321'
```

-netdev tap,id=id[,fd=h][,ifname=name][,script=file][,downscript=dfile][,br=bridge][,helper=helper] Configure a host TAP network backend with ID id.

Use the network script file to configure it and the network script dfile to deconfigure it. If name is not provided, the OS automatically provides one. The default network configure script is /etc/qemu-ifup and the default network deconfigure script is /etc/qemu-ifdown. Use script=no or downscript=no to disable script execution.

If running QEMU as an unprivileged user, use the network helper to configure the TAP interface and attach it to the bridge. The default network helper executable is /path/to/qemu-bridge-helper and the default bridge device is br0.

fd=h can be used to specify the handle of an already opened host TAP interface.

Examples:

```
#launch a QEMU instance with the default network script
qemu-system-x86_64 linux.img -nic tap

#launch a QEMU instance with two NICs, each one connected
#to a TAP device
qemu-system-x86_64 linux.img \
    -netdev tap,id=nd0,ifname=tap0 -device e1000,netdev=nd0 \
    -netdev tap,id=nd1,ifname=tap1 -device rtl8139,netdev=nd1

#launch a QEMU instance with the default network helper to
#connect a TAP device to bridge br0
qemu-system-x86_64 linux.img -device virtio-net-pci,netdev=n1 \
    -netdev tap,id=n1,"helper=/path/to/qemu-bridge-helper"
```

-netdev bridge,id=id[,br=bridge][,helper=helper] Connect a host TAP network interface to a host bridge device.

Use the network helper helper to configure the TAP interface and attach it to the bridge. The default network helper executable is /path/to/qemu-bridge-helper and the default bridge device is br0.

Examples:

```
#launch a QEMU instance with the default network helper to
#connect a TAP device to bridge br0
qemu-system-x86_64 linux.img -netdev bridge,id=n1 -device virtio-net,
↳netdev=n1

#launch a QEMU instance with the default network helper to
#connect a TAP device to bridge qemubr0
qemu-system-x86_64 linux.img -netdev bridge,br=qemubr0,id=n1 -device_
↳virtio-net,netdev=n1
```

-netdev socket,id=id[,fd=h][,listen=[host]:port][,connect=host:port] This host network backend can be used to connect the guest's network to another QEMU virtual machine using a TCP socket connection. If listen is specified, QEMU waits for incoming connections on port (host is optional). connect is used to connect to another QEMU instance using the listen option. fd=h specifies an already opened TCP socket.

Example:

```
# launch a first QEMU instance
qemu-system-x86_64 linux.img \
```



```
-device e1000,netdev=n1,mac=52:54:00:12:34:56 \  
-netdev socket,id=n1,listen=:1234  
# connect the network of this instance to the network of the first  
→instance  
qemu-system-x86_64 linux.img \  
-device e1000,netdev=n2,mac=52:54:00:12:34:57 \  
-netdev socket,id=n2,connect=127.0.0.1:1234
```

-netdev socket,id=id[,fd=h][,mcast=maddr:port[,localaddr=addr]] Configure a socket host network backend to share the guest's network traffic with another QEMU virtual machines using a UDP multicast socket, effectively making a bus for every QEMU with same multicast address maddr and port. NOTES:

1. Several QEMU can be running on different hosts and share same bus (assuming correct multicast setup for these hosts).
2. mcast support is compatible with User Mode Linux (argument ethN=mcast), see <http://user-mode-linux.sf.net>.
3. Use fd=h to specify an already opened UDP multicast socket.

Example:

```
# launch one QEMU instance  
qemu-system-x86_64 linux.img \  
-device e1000,netdev=n1,mac=52:54:00:12:34:56 \  
-netdev socket,id=n1,mcast=230.0.0.1:1234  
# launch another QEMU instance on same "bus"  
qemu-system-x86_64 linux.img \  
-device e1000,netdev=n2,mac=52:54:00:12:34:57 \  
-netdev socket,id=n2,mcast=230.0.0.1:1234  
# launch yet another QEMU instance on same "bus"  
qemu-system-x86_64 linux.img \  
-device e1000,netdev=n3,mac=52:54:00:12:34:58 \  
-netdev socket,id=n3,mcast=230.0.0.1:1234
```

Example (User Mode Linux compat.):

```
# launch QEMU instance (note mcast address selected is UML's default)  
qemu-system-x86_64 linux.img \  
-device e1000,netdev=n1,mac=52:54:00:12:34:56 \  
-netdev socket,id=n1,mcast=239.192.168.1:1102  
# launch UML  
/path/to/linux ubd0=/path/to/root_fs eth0=mcast
```

Example (send packets from host's 1.2.3.4):

```
qemu-system-x86_64 linux.img \  
-device e1000,netdev=n1,mac=52:54:00:12:34:56 \  
-netdev socket,id=n1,mcast=239.192.168.1:1102,  
→localaddr=1.2.3.4
```

-netdev l2tpv3,id=id,src=srcaddr,dst=dstaddr[,srcport=srcport][,dstport=dstport],txsession= Configure a L2TPv3 pseudowire host network backend. L2TPv3 (RFC3931) is a popular protocol to transport Ethernet (and other Layer 2) data frames between two systems. It is present in routers, firewalls and the Linux kernel (from version 3.3 onwards).

This transport allows a VM to communicate to another VM, router or firewall directly.

src=srcaddr source address (mandatory)

dst=dstaddr destination address (mandatory)

udp select udp encapsulation (default is ip).

srcport=srcport source udp port.

dstport=dstport destination udp port.

ipv6 force v6, otherwise defaults to v4.

rxcookie=rxcookie; txcookie=txcookie Cookies are a weak form of security in the l2tpv3 specification. Their function is mostly to prevent misconfiguration. By default they are 32 bit.

cookie64 Set cookie size to 64 bit instead of the default 32

counter=off Force a ‘cut-down’ L2TPv3 with no counter as in draft-mkonstan-l2tpext-keyed-ipv6-tunnel-00

pincounter=on Work around broken counter handling in peer. This may also help on networks which have packet reorder.

offset=offset Add an extra offset between header and data

For example, to attach a VM running on host 4.3.2.1 via L2TPv3 to the bridge br-lan on the remote Linux host 1.2.3.4:

```
# Setup tunnel on linux host using raw ip as encapsulation
# on 1.2.3.4
ip l2tp add tunnel remote 4.3.2.1 local 1.2.3.4 tunnel_id 1 peer_tunnel_
→id 1 \
    encap udp udp_sport 16384 udp_dport 16384
ip l2tp add session tunnel_id 1 name vmtunnel0 session_id \
    0xFFFFFFFF peer_session_id 0xFFFFFFFF
ifconfig vmtunnel0 mtu 1500
ifconfig vmtunnel0 up
brctl addif br-lan vmtunnel0

# on 4.3.2.1
# launch QEMU instance - if your network has reorder or is very lossy add_
→,pincounter

qemu-system-x86_64 linux.img -device e1000,netdev=n1 \
    -netdev l2tpv3,id=n1,src=4.2.3.1,dst=1.2.3.4,udp,srcport=16384,
→dstport=16384,rxsession=0xffffffff,txsession=0xffffffff,counter
```

-netdev vde,id=id[,sock=socketpath][,port=n][,group=groupname][,mode=octalmode]

Configure VDE backend to connect to PORT n of a vde switch running on host and listening for incoming connections on socketpath. Use GROUP groupname and MODE octalmode to change default ownership and permissions for communication port. This option is only available if QEMU has been compiled with vde support enabled.

Example:

```
# launch vde switch
vde_switch -F -sock /tmp/myswitch
# launch QEMU instance
qemu-system-x86_64 linux.img -nic vde,sock=/tmp/myswitch
```

-netdev vhost-user,chardev=id[,vhostforce=on|off][,queues=n] Establish a vhost-user netdev, backed by a chardev id. The chardev should be a unix domain socket backed one. The vhost-user uses a

specifically defined protocol to pass vhost ioctl replacement messages to an application on the other end of the socket. On non-MSIX guests, the feature can be forced with `vhostforce`. Use `'queues=n'` to specify the number of queues to be created for multiqueue vhost-user.

Example:

```
qemu -m 512 -object memory-backend-file,id=mem,size=512M,mem-path=/hugetlbfs,
↪share=on \
    -numa node,memdev=mem \
    -chardev socket,id=chr0,path=/path/to/socket \
    -netdev type=vhost-user,id=net0,chardev=chr0 \
    -device virtio-net-pci,netdev=net0
```

-netdev vhost-vdpa, vhostdev=/path/to/dev Establish a vhost-vdpa netdev.

vDPA device is a device that uses a datapath which complies with the virtio specifications with a vendor specific control path. vDPA devices can be both physically located on the hardware or emulated by software.

-netdev hubport, id=id, hubid=hubid[, netdev=nd] Create a hub port on the emulated hub with ID hubid.

The hubport netdev lets you connect a NIC to a QEMU emulated hub instead of a single netdev. Alternatively, you can also connect the hubport to another netdev with ID nd by using the `netdev=nd` option.

-net nic[, netdev=nd][, macaddr=mac][, model=type][, name=name][, addr=addr][, vectors=v]

Legacy option to configure or create an on-board (or machine default) Network Interface Card(NIC) and connect it either to the emulated hub with ID 0 (i.e. the default hub), or to the netdev nd. If model is omitted, then the default NIC model associated with the machine type is used. Note that the default NIC model may change in future QEMU releases, so it is highly recommended to always specify a model. Optionally, the MAC address can be changed to mac, the device address set to addr (PCI cards only), and a name can be assigned for use in monitor commands. Optionally, for PCI cards, you can specify the number v of MSI-X vectors that the card should have; this option currently only affects virtio cards; set v = 0 to disable MSI-X. If no `-net` option is specified, a single NIC is created. QEMU can emulate several different models of network card. Use `-net nic,model=help` for a list of available devices for your target.

-net user|tap|bridge|socket|l2tpv3|vde[, ...][, name=name] Configure a host network back-end (with the options corresponding to the same `-netdev` option) and connect it to the emulated hub 0 (the default hub). Use name to specify the name of the hub port.

1.2.7 Character device options

The general form of a character device option is:

-chardev backend, id=id[, mux=on|off][, options] Backend is one of: `null`, `socket`, `udp`, `msmouse`, `vc`, `ringbuf`, `file`, `pipe`, `console`, `serial`, `pty`, `stdio`, `braille`, `tty`, `parallel`, `parport`, `spicevmc`, `spiceport`. The specific backend will determine the applicable options.

Use `-chardev help` to print all available chardev backend types.

All devices must have an id, which can be any string up to 127 characters long. It is used to uniquely identify this device in other command line directives.

A character device may be used in multiplexing mode by multiple front-ends. Specify `mux=on` to enable this mode. A multiplexer is a “1:N” device, and here the “1” end is your specified chardev backend, and the “N” end is the various parts of QEMU that can talk to a chardev. If you create a chardev with `id=myid` and `mux=on`, QEMU will create a multiplexer with your specified ID, and you can then configure multiple front ends to use that chardev ID for their input/output. Up to four different front ends can be connected to a single multiplexed chardev. (Without multiplexing enabled, a chardev can only be used by a single front end.) For instance you could use this to allow a single stdio chardev to be used by two serial ports and the QEMU monitor:

```
-chardev stdio,mux=on,id=char0 \
-mon chardev=char0,mode=readline \
-serial chardev:char0 \
-serial chardev:char0
```

You can have more than one multiplexer in a system configuration; for instance you could have a TCP port multiplexed between UART 0 and UART 1, and stdio multiplexed between the QEMU monitor and a parallel port:

```
-chardev stdio,mux=on,id=char0 \
-mon chardev=char0,mode=readline \
-parallel chardev:char0 \
-chardev tcp,...,mux=on,id=char1 \
-serial chardev:char1 \
-serial chardev:char1
```

When you're using a multiplexed character device, some escape sequences are interpreted in the input. See the chapter about *Keys in the character backend multiplexer* in the System Emulation Users Guide for more details.

Note that some other command line options may implicitly create multiplexed character backends; for instance `-serial mon:stdio` creates a multiplexed stdio backend connected to the serial port and the QEMU monitor, and `-nographic` also multiplexes the console and the monitor to stdio.

There is currently no support for multiplexing in the other direction (where a single QEMU front end takes input and output from multiple chardevs).

Every backend supports the `logfile` option, which supplies the path to a file to record all data transmitted via the backend. The `logappend` option controls whether the log file will be truncated or appended to when opened.

The available backends are:

-chardev null,id=id A void device. This device will not emit any data, and will drop any data it receives. The null backend does not take any options.

-chardev socket,id=id[,TCP options or unix options][,server][,nowait][,telnet][,websocket]

Create a two-way stream socket, which can be either a TCP or a unix socket. A unix socket will be created if path is specified. Behaviour is undefined if TCP options are specified for a unix socket.

`server` specifies that the socket shall be a listening socket.

`nowait` specifies that QEMU should not block waiting for a client to connect to a listening socket.

`telnet` specifies that traffic on the socket should interpret telnet escape sequences.

`websocket` specifies that the socket uses WebSocket protocol for communication.

`reconnect` sets the timeout for reconnecting on non-server sockets when the remote end goes away. qemu will delay this many seconds and then attempt to reconnect. Zero disables reconnecting, and is the default.

`tls-creds` requests enablement of the TLS protocol for encryption, and specifies the id of the TLS credentials to use for the handshake. The credentials must be previously created with the `-object tls-creds` argument.

`tls-auth` provides the ID of the QAuthZ authorization object against which the client's x509 distinguished name will be validated. This object is only resolved at time of use, so can be deleted and recreated on the fly while the chardev server is active. If missing, it will default to denying access.

TCP and unix socket options are given below:

TCP options: `port=port[,host=host][,to=to][,ipv4][,ipv6][,nodelay]` `host` for a listening socket specifies the local address to be bound. For a connecting socket species the remote host to connect to. `host` is optional for listening sockets. If not specified it defaults to 0.0.0.0.

`port` for a listening socket specifies the local port to be bound. For a connecting socket specifies the port on the remote host to connect to. `port` can be given as either a port number or a service name. `port` is required.

`to` is only relevant to listening sockets. If it is specified, and `port` cannot be bound, QEMU will attempt to bind to subsequent ports up to and including `to` until it succeeds. `to` must be specified as a port number.

`ipv4` and `ipv6` specify that either IPv4 or IPv6 must be used. If neither is specified the socket may use either protocol.

`nodelay` disables the Nagle algorithm.

unix options: `path=path[,abstract=on|off][,tight=on|off]` `path` specifies the local path of the unix socket. `path` is required. `abstract` specifies the use of the abstract socket namespace, rather than the filesystem. Optional, defaults to false. `tight` sets the socket length of abstract sockets to their minimum, rather than the full `sun_path` length. Optional, defaults to true.

-chardev udp,id=id[,host=host],port=port[,localaddr=localaddr][,localport=localport][,ipv4]
Sends all traffic from the guest to a remote host over UDP.

`host` specifies the remote host to connect to. If not specified it defaults to `localhost`.

`port` specifies the port on the remote host to connect to. `port` is required.

`localaddr` specifies the local address to bind to. If not specified it defaults to 0.0.0.0.

`localport` specifies the local port to bind to. If not specified any available local port will be used.

`ipv4` and `ipv6` specify that either IPv4 or IPv6 must be used. If neither is specified the device may use either protocol.

-chardev msmouse,id=id Forward QEMU's emulated msmouse events to the guest. `msmouse` does not take any options.

-chardev vc,id=id[,width=width][,height=height][[,cols=cols][,rows=rows]]
Connect to a QEMU text console. `vc` may optionally be given a specific size.

`width` and `height` specify the width and height respectively of the console, in pixels.

`cols` and `rows` specify that the console be sized to fit a text console with the given dimensions.

-chardev ringbuf,id=id[,size=size] Create a ring buffer with fixed size `size`. `size` must be a power of two and defaults to 64K.

-chardev file,id=id,path=path Log all traffic received from the guest to a file.

`path` specifies the path of the file to be opened. This file will be created if it does not already exist, and overwritten if it does. `path` is required.

-chardev pipe,id=id,path=path Create a two-way connection to the guest. The behaviour differs slightly between Windows hosts and other hosts:

On Windows, a single duplex pipe will be created at `\\.pipe\\path`.

On other hosts, 2 pipes will be created called `path.in` and `path.out`. Data written to `path.in` will be received by the guest. Data written by the guest can be read from `path.out`. QEMU will not create these fifos, and requires them to be present.

`path` forms part of the pipe path as described above. `path` is required.

-chardev console, id=id Send traffic from the guest to QEMU's standard output. `console` does not take any options.

`console` is only available on Windows hosts.

-chardev serial, id=id, path=path Send traffic from the guest to a serial device on the host.

On Unix hosts `serial` will actually accept any tty device, not only serial lines.

`path` specifies the name of the serial device to open.

-chardev pty, id=id Create a new pseudo-terminal on the host and connect to it. `pty` does not take any options.

`pty` is not available on Windows hosts.

-chardev stdio, id=id[, signal=on|off] Connect to standard input and standard output of the QEMU process.

`signal` controls if signals are enabled on the terminal, that includes exiting QEMU with the key sequence Control-c. This option is enabled by default, use `signal=off` to disable it.

-chardev braille, id=id Connect to a local BrAPI server. `braille` does not take any options.

-chardev tty, id=id, path=path `tty` is only available on Linux, Sun, FreeBSD, NetBSD, OpenBSD and DragonFlyBSD hosts. It is an alias for `serial`.

`path` specifies the path to the tty. `path` is required.

-chardev parallel, id=id, path=path

-chardev parport, id=id, path=path `parallel` is only available on Linux, FreeBSD and DragonFlyBSD hosts.

Connect to a local parallel port.

`path` specifies the path to the parallel port device. `path` is required.

-chardev spicevmc, id=id, debug=debug, name=name `spicevmc` is only available when spice support is built in.

`debug` debug level for `spicevmc`

`name` name of spice channel to connect to

Connect to a spice virtual machine channel, such as `vdiport`.

-chardev spiceport, id=id, debug=debug, name=name `spiceport` is only available when spice support is built in.

`debug` debug level for `spicevmc`

`name` name of spice port to connect to

Connect to a spice port, allowing a Spice client to handle the traffic identified by a name (preferably a fqdn).

1.2.8 TPM device options

The general form of a TPM device option is:

-tpmdev backend, id=id[, options] The specific backend type will determine the applicable options. The `-tpmdev` option creates the TPM backend and requires a `-device` option that specifies the TPM frontend interface model.

Use `-tpmdev help` to print all available TPM backend types.

The available backends are:

-tpmdev passthrough, id=id, path=path, cancel-path=cancel-path (Linux-host only) Enable access to the host's TPM using the passthrough driver.

`path` specifies the path to the host's TPM device, i.e., on a Linux host this would be `/dev/tpm0`. `path` is optional and by default `/dev/tpm0` is used.

`cancel-path` specifies the path to the host TPM device's sysfs entry allowing for cancellation of an ongoing TPM command. `cancel-path` is optional and by default QEMU will search for the sysfs entry to use.

Some notes about using the host's TPM with the passthrough driver:

The TPM device accessed by the passthrough driver must not be used by any other application on the host.

Since the host's firmware (BIOS/UEFI) has already initialized the TPM, the VM's firmware (BIOS/UEFI) will not be able to initialize the TPM again and may therefore not show a TPM-specific menu that would otherwise allow the user to configure the TPM, e.g., allow the user to enable/disable or activate/deactivate the TPM. Further, if TPM ownership is released from within a VM then the host's TPM will get disabled and deactivated. To enable and activate the TPM again afterwards, the host has to be rebooted and the user is required to enter the firmware's menu to enable and activate the TPM. If the TPM is left disabled and/or deactivated most TPM commands will fail.

To create a passthrough TPM use the following two options:

```
-tpmdev passthrough, id=tpm0 -device tpm-tis, tpmdev=tpm0
```

Note that the `-tpmdev id` is `tpm0` and is referenced by `tpmdev=tpm0` in the device option.

-tpmdev emulator, id=id, chardev=dev (Linux-host only) Enable access to a TPM emulator using Unix domain socket based chardev backend.

`chardev` specifies the unique ID of a character device backend that provides connection to the software TPM server.

To create a TPM emulator backend device with chardev socket backend:

```
-chardev socket, id=chrtpm, path=/tmp/swtpm-sock -tpmdev emulator, id=tpm0,  
↪ chardev=chrtpm -device tpm-tis, tpmdev=tpm0
```

1.2.9 Linux/Multiboot boot specific

When using these options, you can use a given Linux or Multiboot kernel without installing it in the disk image. It can be useful for easier testing of various kernels.

-kernel bzImage Use `bzImage` as kernel image. The kernel can be either a Linux kernel or in multiboot format.

-append cmdline Use `cmdline` as kernel command line

-initrd file Use `file` as initial ram disk.

-initrd "file1 arg=foo, file2" This syntax is only available with multiboot.

Use `file1` and `file2` as modules and pass `arg=foo` as parameter to the first module.

-dtb file Use `file` as a device tree binary (dtb) image and pass it to the kernel on boot.

1.2.10 Debug/Expert options

-fw_cfg [name=]name, file=file Add named `fw_cfg` entry with contents from `file`.

-fw_cfg [name=]name, string=str Add named fw_cfg entry with contents from string str.

The terminating NUL character of the contents of str will not be included as part of the fw_cfg item data. To insert contents with embedded NUL characters, you have to use the file parameter.

The fw_cfg entries are passed by QEMU through to the guest.

Example:

```
-fw_cfg name=opt/com.mycompany/blob,file=./my_blob.bin
```

creates an fw_cfg entry named opt/com.mycompany/blob with contents from ./my_blob.bin.

-serial dev Redirect the virtual serial port to host character device dev. The default device is vc in graphical mode and stdio in non graphical mode.

This option can be used several times to simulate up to 4 serial ports.

Use `-serial none` to disable all serial ports.

Available character devices are:

vc[:WxH] Virtual console. Optionally, a width and height can be given in pixel with

```
vc:800x600
```

It is also possible to specify width or height in characters:

```
vc:80Cx24C
```

pty [Linux only] Pseudo TTY (a new PTY is automatically allocated)

none No device is allocated.

null void device

chardev:id Use a named character device defined with the `-chardev` option.

/dev/XXX [Linux only] Use host tty, e.g. `/dev/ttyS0`. The host serial port parameters are set according to the emulated ones.

/dev/parportN [Linux only, parallel port only] Use host parallel port N. Currently SPP and EPP parallel port features can be used.

file:filename Write output to filename. No character can be read.

stdio [Unix only] standard input/output

pipe:filename name pipe filename

COMn [Windows only] Use host serial port n

udp:[remote_host]:remote_port[@[src_ip]:src_port] This implements UDP Net Console. When remote_host or src_ip are not specified they default to 0.0.0.0. When not using a specified src_port a random port is automatically chosen.

If you just want a simple readonly console you can use `netcat` or `nc`, by starting QEMU with: `-serial udp::4555` and `nc` as: `nc -u -l -p 4555`. Any time QEMU writes something to that port it will appear in the netconsole session.

If you plan to send characters back via netconsole or you want to stop and start QEMU a lot of times, you should have QEMU use the same source port each time by using something like `-serial udp::4555@:4556` to QEMU. Another approach is to use a patched version of netcat which can listen to a TCP port and send and receive characters via udp. If you have a patched version of netcat which

activates telnet remote echo and single char transfer, then you can use the following options to set up a netcat redirector to allow telnet on port 5555 to access the QEMU port.

QEMU Options: `-serial udp::4555@:4556`

netcat options: `-u -P 4555 -L 0.0.0.0:4556 -t -p 5555 -I -T`

telnet options: `localhost 5555`

tcp: `[host]:port[,server][,nowait][,nodelay][,reconnect=seconds]` The TCP Net Console has two modes of operation. It can send the serial I/O to a location or wait for a connection from a location. By default the TCP Net Console is sent to host at the port. If you use the server option QEMU will wait for a client socket application to connect to the port before continuing, unless the `nowait` option was specified. The `nodelay` option disables the Nagle buffering algorithm. The `reconnect` option only applies if `noserver` is set, if the connection goes down it will attempt to reconnect at the given interval. If host is omitted, 0.0.0.0 is assumed. Only one TCP connection at a time is accepted. You can use `telnet` to connect to the corresponding character device.

Example to send tcp console to 192.168.0.2 port 4444 `-serial tcp:192.168.0.2:4444`

Example to listen and wait on port 4444 for connection `-serial tcp::4444,server`

Example to not wait and listen on ip 192.168.0.100 port 4444 `-serial tcp:192.168.0.100:4444,server,nowait`

telnet: `host:port[,server][,nowait][,nodelay]` The telnet protocol is used instead of raw tcp sockets. The options work the same as if you had specified `-serial tcp`. The difference is that the port acts like a telnet server or client using telnet option negotiation. This will also allow you to send the MAGIC_SYSRQ sequence if you use a telnet that supports sending the break sequence. Typically in unix telnet you do it with Control-] and then type “send break” followed by pressing the enter key.

websocket: `host:port,server[,nowait][,nodelay]` The WebSocket protocol is used instead of raw tcp socket. The port acts as a WebSocket server. Client mode is not supported.

unix: `path[,server][,nowait][,reconnect=seconds]` A unix domain socket is used instead of a tcp socket. The option works the same as if you had specified `-serial tcp` except the unix domain socket path is used for connections.

mon: `dev_string` This is a special option to allow the monitor to be multiplexed onto another serial port. The monitor is accessed with key sequence of Control-a and then pressing c. `dev_string` should be any one of the serial devices specified above. An example to multiplex the monitor onto a telnet server listening on port 4444 would be:

`-serial mon:telnet::4444,server,nowait`

When the monitor is multiplexed to stdio in this way, Ctrl+C will not terminate QEMU any more but will be passed to the guest instead.

braille Braille device. This will use BrlAPI to display the braille output on a real or fake device.

msmouse Three button serial mouse. Configure the guest to use Microsoft protocol.

-parallel dev Redirect the virtual parallel port to host device dev (same devices as the serial port). On Linux hosts, `/dev/parportN` can be used to use hardware devices connected on the corresponding host parallel port.

This option can be used several times to simulate up to 3 parallel ports.

Use `-parallel none` to disable all parallel ports.

-monitor dev Redirect the monitor to host device dev (same devices as the serial port). The default device is `vc` in graphical mode and `stdio` in non graphical mode. Use `-monitor none` to disable the default monitor.

- qmp dev** Like `-monitor` but opens in ‘control’ mode.
- qmp-pretty dev** Like `-qmp` but uses pretty JSON formatting.
- mon [chardev=]name[,mode=readline|control][,pretty=[on|off]]** Setup monitor on chardev name. `pretty` turns on JSON pretty printing easing human reading and debugging.
- debugcon dev** Redirect the debug console to host device dev (same devices as the serial port). The debug console is an I/O port which is typically port 0xe9; writing to that I/O port sends output to this device. The default device is `vc` in graphical mode and `stdio` in non graphical mode.
- pidfile file** Store the QEMU process PID in file. It is useful if you launch QEMU from a script.
- singlestep** Run the emulation in single step mode.
- preconfig** Pause QEMU for interactive configuration before the machine is created, which allows querying and configuring properties that will affect machine initialization. Use QMP command ‘x-exit-preconfig’ to exit the preconfig state and move to the next state (i.e. run guest if `-S` isn’t used or pause the second time if `-S` is used). This option is experimental.
- S** Do not start CPU at startup (you must type ‘c’ in the monitor).
- overcommit mem-lock=on|off**
- overcommit cpu-pm=on|off** Run qemu with hints about host resource overcommit. The default is to assume that host overcommits all resources.

Locking qemu and guest memory can be enabled via `mem-lock=on` (disabled by default). This works when host memory is not overcommitted and reduces the worst-case latency for guest.

Guest ability to manage power state of host cpus (increasing latency for other processes on the same host cpu, but decreasing latency for guest) can be enabled via `cpu-pm=on` (disabled by default). This works best when host CPU is not overcommitted. When used, host estimates of CPU cycle and power utilization will be incorrect, not taking into account guest idle time.
- gdb dev** Accept a gdb connection on device dev (see the [GDB usage](#) chapter in the System Emulation Users Guide). Note that this option does not pause QEMU execution – if you want QEMU to not start the guest until you connect with gdb and issue a `continue` command, you will need to also pass the `-S` option to QEMU.

The most usual configuration is to listen on a local TCP socket:

```
-gdb tcp::3117
```

but you can specify other backends; UDP, pseudo TTY, or even stdio are all reasonable use cases. For example, a stdio connection allows you to start QEMU from within gdb and establish the connection via a pipe:

```
(gdb) target remote | exec qemu-system-x86_64 -gdb stdio ...
```

- s** Shorthand for `-gdb tcp::1234`, i.e. open a gdbserver on TCP port 1234 (see the [GDB usage](#) chapter in the System Emulation Users Guide).
- d item1[, ...]** Enable logging of specified items. Use ‘-d help’ for a list of log items.
- D logfile** Output log in logfile instead of to stderr
- dfilter range1[, ...]** Filter debug output to that relevant to a range of target addresses. The filter spec can be either start+size, start-size or start..end where start end and size are the addresses and sizes required. For example:

```
-dfilter 0x8000..0x8fff,0xffffffffc000080000+0x200,0xffffffffc000060000-0x1000
```

Will dump output for any code in the 0x1000 sized block starting at 0x8000 and the 0x200 sized block starting at 0xffffffffc000080000 and another 0x1000 sized block starting at 0xffffffffc00005f000.

-seed number Force the guest to use a deterministic pseudo-random number generator, seeded with number. This does not affect crypto routines within the host.

-L path Set the directory for the BIOS, VGA BIOS and keymaps.

To list all the data directories, use `-L help`.

-bios file Set the filename for the BIOS.

-enable-kvm Enable KVM full virtualization support. This option is only available if KVM support is enabled when compiling.

-xen-domid id Specify xen guest domain id (XEN only).

-xen-attach Attach to existing xen domain. libxl will use this when starting QEMU (XEN only). Restrict set of available xen operations to specified domain id (XEN only).

-no-reboot Exit instead of rebooting.

-no-shutdown Don't exit QEMU on guest shutdown, but instead only stop the emulation. This allows for instance switching to monitor to commit changes to the disk image.

-action event=action The action parameter serves to modify QEMU's default behavior when certain guest events occur. It provides a generic method for specifying the same behaviors that are modified by the `-no-reboot` and `-no-shutdown` parameters.

Examples:

```
-action panic=none      -action reboot=shutdown,shutdown=pause    -watchdog
i6300esb -action watchdog=pause
```

-loadvm file Start right away with a saved state (loadvm in monitor)

-daemonize Daemonize the QEMU process after initialization. QEMU will not detach from standard IO until it is ready to receive connections on any of its devices. This option is a useful way for external programs to launch QEMU without having to cope with initialization race conditions.

-option-rom file Load the contents of file as an option ROM. This option is useful to load things like Ether-Boot.

-rtc [base=utc|localtime|datetime] [, clock=host|rt|vm] [, driftfix=none|slew]

Specify base as `utc` or `localtime` to let the RTC start at the current UTC or local time, respectively. `localtime` is required for correct date in MS-DOS or Windows. To start at a specific point in time, provide `datetime` in the format `2006-06-17T16:01:21` or `2006-06-17`. The default base is UTC.

By default the RTC is driven by the host system time. This allows using of the RTC as accurate reference clock inside the guest, specifically if the host time is smoothly following an accurate external reference clock, e.g. via NTP. If you want to isolate the guest time from the host, you can set `clock` to `rt` instead, which provides a host monotonic clock if host support it. To even prevent the RTC from progressing during suspension, you can set `clock` to `vm` (virtual clock). '`clock=vm`' is recommended especially in `icount` mode in order to preserve determinism; however, note that in `icount` mode the speed of the virtual clock is variable and can in general differ from the host clock.

Enable `driftfix` (i386 targets only) if you experience time drift problems, specifically with Windows' ACPI HAL. This option will try to figure out how many timer interrupts were not processed by the Windows guest and will re-inject them.

-icount [shift=N|auto] [, align=on|off] [, sleep=on|off] [, rr=record|replay, rrfile=filename [, rrsize=...]]

Enable virtual instruction counter. The virtual cpu will execute one instruction every 2^N ns of virtual time. If `auto` is specified then the virtual cpu speed will be automatically adjusted to keep virtual time within a few seconds of real time.

Note that while this option can give deterministic behavior, it does not provide cycle accurate emulation. Modern CPUs contain superscalar out of order cores with complex cache hierarchies. The number of instructions executed often has little or no correlation with actual performance.

When the virtual cpu is sleeping, the virtual time will advance at default speed unless `sleep=on` is specified. With `sleep=on`, the virtual time will jump to the next timer deadline instantly whenever the virtual cpu goes to sleep mode and will not advance if no timer is enabled. This behavior gives deterministic execution times from the guest point of view. The default if `icount` is enabled is `sleep=off`. `sleep=on` cannot be used together with either `shift=auto` or `align=on`.

`align=on` will activate the delay algorithm which will try to synchronise the host clock and the virtual clock. The goal is to have a guest running at the real frequency imposed by the `shift` option. Whenever the guest clock is behind the host clock and if `align=on` is specified then we print a message to the user to inform about the delay. Currently this option does not work when `shift` is `auto`. Note: The sync algorithm will work for those shift values for which the guest clock runs ahead of the host clock. Typically this happens when the shift value is high (how high depends on the host machine). The default if `icount` is enabled is `align=off`.

When the `rr` option is specified deterministic record/replay is enabled. The `rrfile=` option must also be provided to specify the path to the replay log. In record mode data is written to this file, and in replay mode it is read back. If the `rrsnapshot` option is given then it specifies a VM snapshot name. In record mode, a new VM snapshot with the given name is created at the start of execution recording. In replay mode this option specifies the snapshot name used to load the initial VM state.

-watchdog model Create a virtual hardware watchdog device. Once enabled (by a guest action), the watchdog must be periodically polled by an agent inside the guest or else the guest will be restarted. Choose a model for which your guest has drivers.

The model is the model of hardware watchdog to emulate. Use `-watchdog help` to list available hardware models. Only one watchdog can be enabled for a guest.

The following models may be available:

ib700 iBASE 700 is a very simple ISA watchdog with a single timer.

i6300esb Intel 6300ESB I/O controller hub is a much more featureful PCI-based dual-timer watchdog.

diag288 A virtual watchdog for s390x backed by the diagnose 288 hypercall (currently KVM only).

-watchdog-action action The action controls what QEMU will do when the watchdog timer expires. The default is `reset` (forcefully reset the guest). Other possible actions are: `shutdown` (attempt to gracefully shutdown the guest), `poweroff` (forcefully poweroff the guest), `inject-nmi` (inject a NMI into the guest), `pause` (pause the guest), `debug` (print a debug message and continue), or `none` (do nothing).

Note that the `shutdown` action requires that the guest responds to ACPI signals, which it may not be able to do in the sort of situations where the watchdog would have expired, and thus `-watchdog-action shutdown` is not recommended for production use.

Examples:

```
-watchdog i6300esb -watchdog-action pause;-watchdog ib700
```

-echr numeric_ascii_value Change the escape character used for switching to the monitor when using monitor and serial sharing. The default is `0x01` when using the `-nographic` option. `0x01` is equal to pressing Control-a. You can select a different character from the ascii control keys where 1 through 26 map to Control-a through Control-z. For instance you could use the either of the following to change the escape character to Control-t.

```
-echr 0x14;-echr 20
```

-incoming tcp:[host]:port[,to=maxport][,ipv4][,ipv6]

-incoming rdma:host:port[,ipv4][,ipv6] Prepare for incoming migration, listen on a given tcp port.

- incoming unix:socketpath** Prepare for incoming migration, listen on a given unix socket.
- incoming fd:fd** Accept incoming migration from a given filedescriptor.
- incoming exec:cmdline** Accept incoming migration as an output from specified external command.
- incoming defer** Wait for the URI to be specified via `migrate_incoming`. The monitor can be used to change settings (such as migration parameters) prior to issuing the `migrate_incoming` to allow the migration to begin.
- only-migratable** Only allow migratable devices. Devices will not be allowed to enter an unmigratable state.
- nodefaults** Don't create default devices. Normally, QEMU sets the default devices like serial port, parallel port, virtual console, monitor device, VGA adapter, floppy and CD-ROM drive and others. The `-nodefaults` option will disable all those default devices.
- chroot dir** Immediately before starting guest execution, chroot to the specified directory. Especially useful in combination with `-runas`.
- runas user** Immediately before starting guest execution, drop root privileges, switching to the specified user.
- prom-env variable=value** Set OpenBIOS nvram variable to given value (PPC, SPARC only).

```
qemu-system-sparc -prom-env 'auto-boot?=false' \  
-prom-env 'boot-device=sd(0,2,0):d' -prom-env 'boot-args=linux single'
```

```
qemu-system-ppc -prom-env 'auto-boot?=false' \  
-prom-env 'boot-device=hd:2,\yaboot' \  
-prom-env 'boot-args=conf=hd:2,\yaboot.conf'
```

- semihosting** Enable semihosting mode (ARM, M68K, Xtensa, MIPS, Nios II, RISC-V only).

Note that this allows guest direct access to the host filesystem, so should only be used with a trusted guest OS.

See the `-semihosting-config` option documentation for further information about the facilities this enables.

- semihosting-config [enable=on|off] [,target=native|gdb|auto] [,chardev=id] [,arg=str[,...]]**
Enable and configure semihosting (ARM, M68K, Xtensa, MIPS, Nios II, RISC-V only).

Note that this allows guest direct access to the host filesystem, so should only be used with a trusted guest OS.

On Arm this implements the standard semihosting API, version 2.0.

On M68K this implements the “ColdFire GDB” interface used by libgloss.

Xtensa semihosting provides basic file IO calls, such as `open/read/write/seek/select`. Tensilica baremetal libc for ISS and linux platform “sim” use this interface.

On RISC-V this implements the standard semihosting API, version 0.2.

target=native|gdb|auto Defines where the semihosting calls will be addressed, to QEMU (`native`) or to GDB (`gdb`). The default is `auto`, which means `gdb` during debug sessions and `native` otherwise.

chardev=str1 Send the output to a chardev backend output for native or auto output when not in `gdb`

arg=str1, arg=str2, ... Allows the user to pass input arguments, and can be used multiple times to build up a list. The old-style `-kernel/-append` method of passing a command line is still supported for backward compatibility. If both the `--semihosting-config arg` and the `-kernel/-append` are specified, the former is passed to semihosting as it always takes precedence.

- old-param** Old param mode (ARM only).

- sandbox arg[,obsolete=string] [,elevateprivileges=string] [,spawn=string] [,resourcecontrol=string]**
Enable Seccomp mode 2 system call filter. ‘on’ will enable syscall filtering and ‘off’ will disable it. The default is ‘off’.

- obsolete=string** Enable Obsolete system calls
- elevateprivileges=string** Disable set*uid/lgid system calls
- spawn=string** Disable *fork and execve
- resourcecontrol=string** Disable process affinity and scheduler priority
- readconfig file** Read device configuration from file. This approach is useful when you want to spawn QEMU process with many command line options but you don't want to exceed the command line character limit.
- writeconfig file** Write device configuration to file. The file can be either filename to save command line and device configuration into file or dash -) character to print the output to stdout. This can be later used as input file for **-readconfig** option.
- no-user-config** The **-no-user-config** option makes QEMU not load any of the user-provided config files on sysconfdir.
- trace [[enable=]pattern] [,events=FILE] [,file=FILE]** Specify tracing options.
- [enable=]PATTERN
- Immediately enable events matching *PATTERN* (either event name or a globbing pattern). This option is only available if QEMU has been compiled with the *simple*, *log* or *ftrace* tracing backend. To specify multiple events or patterns, specify the **-trace** option multiple times.
- Use **-trace help** to print a list of names of trace points.
- events=FILE
- Immediately enable events listed in *FILE*. The file must contain one event name (as listed in the *trace-events-all* file) per line; globbing patterns are accepted too. This option is only available if QEMU has been compiled with the *simple*, *log* or *ftrace* tracing backend.
- file=FILE
- Log output traces to *FILE*. This option is only available if QEMU has been compiled with the *simple* tracing backend.
- plugin file=file[,arg=string]** Load a plugin.
- file=file** Load the given plugin from a shared library file.
- arg=string** Argument string passed to the plugin. (Can be given multiple times.)
- enable-fips** Enable FIPS 140-2 compliance mode.
- msg [timestamp[=on|off]] [,guest-name[=on|off]]** Control error message format.
- timestamp=on|off** Prefix messages with a timestamp. Default is off.
- guest-name=on|off** Prefix messages with guest name but only if **-name guest** option is set otherwise the option is ignored. Default is off.
- dump-vmstate file** Dump json-encoded vmstate information for current machine type to file in file
- enable-sync-profile** Enable synchronization profiling.

1.2.11 Generic object creation

- object typename[,prop1=value1,...]** Create a new object of type *typename* setting properties in the order they are specified. Note that the 'id' property must be set. These objects are placed in the '/objects' path.

-object memory-backend-file, id=id, size=size, mem-path=dir, share=on|off, discard-data=on|off

Creates a memory file backend object, which can be used to back the guest RAM with huge pages.

The `id` parameter is a unique ID that will be used to reference this memory region when configuring the `-numa` argument.

The `size` option provides the size of the memory region, and accepts common suffixes, eg 500M.

The `mem-path` provides the path to either a shared memory or huge page filesystem mount.

The `share` boolean option determines whether the memory region is marked as private to QEMU, or shared. The latter allows a co-operating external process to access the QEMU memory region.

The `share` is also required for pvr dma devices due to limitations in the RDMA API provided by Linux.

Setting `share=on` might affect the ability to configure NUMA bindings for the memory backend under some circumstances, see `Documentation/vm/numa_memory_policy.txt` on the Linux kernel source tree for additional details.

Setting the `discard-data` boolean option to `on` indicates that file contents can be destroyed when QEMU exits, to avoid unnecessarily flushing data to the backing file. Note that `discard-data` is only an optimization, and QEMU might not discard file contents if it aborts unexpectedly or is terminated using `SIGKILL`.

The `merge` boolean option enables memory merge, also known as `MADV_MERGEABLE`, so that Kernel Samepage Merging will consider the pages for memory deduplication.

Setting the `dump` boolean option to `off` excludes the memory from core dumps. This feature is also known as `MADV_DONTDUMP`.

The `prealloc` boolean option enables memory preallocation.

The `host-nodes` option binds the memory range to a list of NUMA host nodes.

The `policy` option sets the NUMA policy to one of the following values:

default default host policy

preferred prefer the given host node list for allocation

bind restrict memory allocation to the given host node list

interleave interleave memory allocations across the given host node list

The `align` option specifies the base address alignment when QEMU `mmap(2)` `mem-path`, and accepts common suffixes, eg 2M. Some backend store specified by `mem-path` requires an alignment different than the default one used by QEMU, eg the device DAX `/dev/dax0.0` requires 2M alignment rather than 4K. In such cases, users can specify the required alignment via this option.

The `pmem` option specifies whether the backing file specified by `mem-path` is in host persistent memory that can be accessed using the SNIA NVM programming model (e.g. Intel NVDIMM). If `pmem` is set to `on`, QEMU will take necessary operations to guarantee the persistence of its own writes to `mem-path` (e.g. in vNVDIMM label emulation and live migration). Also, we will map the backend-file with `MAP_SYNC` flag, which ensures the file metadata is in sync for `mem-path` in case of host crash or a power failure. `MAP_SYNC` requires support from both the host kernel (since Linux kernel 4.15) and the filesystem of `mem-path` mounted with DAX option.

The `readonly` option specifies whether the backing file is opened read-only or read-write (default).

-object memory-backend-ram, id=id, merge=on|off, dump=on|off, share=on|off, prealloc=on|off

Creates a memory backend object, which can be used to back the guest RAM. Memory backend objects offer more control than the `-m` option that is traditionally used to define guest RAM. Please refer to `memory-backend-file` for a description of the options.

-object memory-backend-memfd, id=id, merge=on|off, dump=on|off, share=on|off, prealloc=on|off
Creates an anonymous memory file backend object, which allows QEMU to share the memory with an external process (e.g. when using vhost-user). The memory is allocated with memfd and optional sealing. (Linux only)

The `seal` option creates a sealed-file, that will block further resizing the memory ('on' by default).

The `hugetlb` option specify the file to be created resides in the `hugetlbfs` filesystem (since Linux 4.14). Used in conjunction with the `hugetlb` option, the `hugetlbsize` option specify the `hugetlb` page size on systems that support multiple `hugetlb` page sizes (it must be a power of 2 value supported by the system).

In some versions of Linux, the `hugetlb` option is incompatible with the `seal` option (requires at least Linux 4.16).

Please refer to `memory-backend-file` for a description of the other options.

The `share` boolean option is on by default with `memfd`.

-object rng-builtin, id=id Creates a random number generator backend which obtains entropy from QEMU builtin functions. The `id` parameter is a unique ID that will be used to reference this entropy backend from the `virtio-rng` device. By default, the `virtio-rng` device uses this RNG backend.

-object rng-random, id=id, filename=/dev/random Creates a random number generator backend which obtains entropy from a device on the host. The `id` parameter is a unique ID that will be used to reference this entropy backend from the `virtio-rng` device. The `filename` parameter specifies which file to obtain entropy from and if omitted defaults to `/dev/urandom`.

-object rng-egd, id=id, chardev=chardev Creates a random number generator backend which obtains entropy from an external daemon running on the host. The `id` parameter is a unique ID that will be used to reference this entropy backend from the `virtio-rng` device. The `chardev` parameter is the unique ID of a character device backend that provides the connection to the RNG daemon.

-object tls-creds-anon, id=id, endpoint=endpoint, dir=/path/to/cred/dir, verify-peer=on|off
Creates a TLS anonymous credentials object, which can be used to provide TLS support on network backends. The `id` parameter is a unique ID which network backends will use to access the credentials. The `endpoint` is either `server` or `client` depending on whether the QEMU network backend that uses the credentials will be acting as a client or as a server. If `verify-peer` is enabled (the default) then once the handshake is completed, the peer credentials will be verified, though this is a no-op for anonymous credentials.

The `dir` parameter tells QEMU where to find the credential files. For server endpoints, this directory may contain a file `dh-params.pem` providing diffie-hellman parameters to use for the TLS server. If the file is missing, QEMU will generate a set of DH parameters at startup. This is a computationally expensive operation that consumes random pool entropy, so it is recommended that a persistent set of parameters be generated upfront and saved.

-object tls-creds-psk, id=id, endpoint=endpoint, dir=/path/to/keys/dir[, username=username]
Creates a TLS Pre-Shared Keys (PSK) credentials object, which can be used to provide TLS support on network backends. The `id` parameter is a unique ID which network backends will use to access the credentials. The `endpoint` is either `server` or `client` depending on whether the QEMU network backend that uses the credentials will be acting as a client or as a server. For clients only, `username` is the username which will be sent to the server. If omitted it defaults to "qemu".

The `dir` parameter tells QEMU where to find the keys file. It is called "dir/keys.psk" and contains "username:key" pairs. This file can most easily be created using the GnuTLS `psktool` program.

For server endpoints, `dir` may also contain a file `dh-params.pem` providing diffie-hellman parameters to use for the TLS server. If the file is missing, QEMU will generate a set of DH parameters at startup. This is a computationally expensive operation that consumes random pool entropy, so it is recommended that a persistent set of parameters be generated up front and saved.

-object tls-creds-x509,id=id,endpoint=endpoint,dir=/path/to/cred/dir,priority=priority

Creates a TLS anonymous credentials object, which can be used to provide TLS support on network backends. The `id` parameter is a unique ID which network backends will use to access the credentials. The `endpoint` is either `server` or `client` depending on whether the QEMU network backend that uses the credentials will be acting as a client or as a server. If `verify-peer` is enabled (the default) then once the handshake is completed, the peer credentials will be verified. With x509 certificates, this implies that the clients must be provided with valid client certificates too.

The `dir` parameter tells QEMU where to find the credential files. For server endpoints, this directory may contain a file `dh-params.pem` providing diffie-hellman parameters to use for the TLS server. If the file is missing, QEMU will generate a set of DH parameters at startup. This is a computationally expensive operation that consumes random pool entropy, so it is recommended that a persistent set of parameters be generated upfront and saved.

For x509 certificate credentials the directory will contain further files providing the x509 certificates. The certificates must be stored in PEM format, in filenames `ca-cert.pem`, `ca-crl.pem` (optional), `server-cert.pem` (only servers), `server-key.pem` (only servers), `client-cert.pem` (only clients), and `client-key.pem` (only clients).

For the `server-key.pem` and `client-key.pem` files which contain sensitive private keys, it is possible to use an encrypted version by providing the `passwordid` parameter. This provides the ID of a previously created `secret` object containing the password for decryption.

The `priority` parameter allows to override the global default priority used by gnutls. This can be useful if the system administrator needs to use a weaker set of crypto priorities for QEMU without potentially forcing the weakness onto all applications. Or conversely if one wants a stronger default for QEMU than for all other applications, they can do this through this parameter. Its format is a gnutls priority string as described at https://gnutls.org/manual/html_node/Priority-Strings.html.

-object tls-cipher-suites,id=id,priority=priority Creates a TLS cipher suites object, which can be used to control the TLS cipher/protocol algorithms that applications are permitted to use.

The `id` parameter is a unique ID which frontends will use to access the ordered list of permitted TLS cipher suites from the host.

The `priority` parameter allows to override the global default priority used by gnutls. This can be useful if the system administrator needs to use a weaker set of crypto priorities for QEMU without potentially forcing the weakness onto all applications. Or conversely if one wants a stronger default for QEMU than for all other applications, they can do this through this parameter. Its format is a gnutls priority string as described at https://gnutls.org/manual/html_node/Priority-Strings.html.

An example of use of this object is to control UEFI HTTPS Boot. The `tls-cipher-suites` object exposes the ordered list of permitted TLS cipher suites from the host side to the guest firmware, via `fw_cfg`. The list is represented as an array of `IANA_TLS_CIPHER` objects. The firmware uses the `IANA_TLS_CIPHER` array for configuring guest-side TLS.

In the following example, the priority at which the host-side policy is retrieved is given by the `priority` property. Given that QEMU uses GNUTLS, `priority=@SYSTEM` may be used to refer to `/etc/crypto-policies/back-ends/gnutls.config`.

```
# qemu-system-x86_64 \  
-object tls-cipher-suites,id=mysuite0,priority=@SYSTEM \  
-fw_cfg name=etc/edk2/https/ciphers,gen_id=mysuite0
```

-object filter-buffer,id=id,netdev=netdev,interval=t[,queue=all|rx|tx][,status=on|off]

Interval `t` can't be 0, this filter batches the packet delivery: all packets arriving in a given interval on `netdev` `netdev` are delayed until the end of the interval. Interval is in microseconds. `status` is optional that indicate whether the netfilter is on (enabled) or off (disabled), the default status for netfilter will be 'on'.

queue `all|rx|tx` is an option that can be applied to any netfilter.

all: the filter is attached both to the receive and the transmit queue of the netdev (default).

rx: the filter is attached to the receive queue of the netdev, where it will receive packets sent to the netdev.

tx: the filter is attached to the transmit queue of the netdev, where it will receive packets sent by the netdev.

position head|tail**id=<id>** is an option to specify where the filter should be inserted in the filter list. It can be applied to any netfilter.

head: the filter is inserted at the head of the filter list, before any existing filters.

tail: the filter is inserted at the tail of the filter list, behind any existing filters (default).

id=<id>: the filter is inserted before or behind the filter specified by **<id>**, see the insert option below.

insert behind|before is an option to specify where to insert the new filter relative to the one specified with **position=id=<id>**. It can be applied to any netfilter.

before: insert before the specified filter.

behind: insert behind the specified filter (default).

-object filter-mirror,id=id,netdev=netdevid,outdev=chardevid,queue=all|rx|tx[,vnet_hdr]
filter-mirror on netdev netdevid,mirror net packet to chardev chardevid, if it has the vnet_hdr_support flag, filter-mirror will mirror packet with vnet_hdr_len.

-object filter-redirector,id=id,netdev=netdevid,indev=chardevid,outdev=chardevid,queue=
filter-redirector on netdev netdevid,redirect filter's net packet to chardev chardevid,and redirect indev's packet to filter.if it has the vnet_hdr_support flag, filter-redirector will redirect packet with vnet_hdr_len. Create a filter-redirector we need to differ outdev id from indev id, id can not be the same. we can just use indev or outdev, but at least one of indev or outdev need to be specified.

-object filter-rewriter,id=id,netdev=netdevid,queue=all|rx|tx,[vnet_hdr_support][,position=
Filter-rewriter is a part of COLO project.It will rewrite tcp packet to secondary from primary to keep secondary tcp connection,and rewrite tcp packet to primary from secondary make tcp packet can be handled by client.if it has the vnet_hdr_support flag, we can parse packet with vnet header.

usage: colo secondary: -object filter-redirector,id=f1,netdev=hn0,queue=tx,indev=red0 -object filter-redirector,id=f2,netdev=hn0,queue=rx,outdev=red1 -object filter-rewriter,id=rew0,netdev=hn0,queue=all

-object filter-dump,id=id,netdev=dev[,file=filename][,maxlen=len][,position=head|tail]
Dump the network traffic on netdev dev to the file specified by filename. At most len bytes (64k by default) per packet are stored. The file format is libpcap, so it can be analyzed with tools such as tcpdump or Wireshark.

-object colo-compare,id=id,primary_in=chardevid,secondary_in=chardevid,outdev=chardevi
Colo-compare gets packet from primary_in chardevid and secondary_in, then compare whether the payload of primary packet and secondary packet are the same. If same, it will output primary packet to out_dev, else it will notify COLO-framework to do checkpoint and send primary packet to out_dev. In order to improve efficiency, we need to put the task of comparison in another iothread. If it has the vnet_hdr_support flag, colo compare will send/rcv packet with vnet_hdr_len. The `compare_timeout=@var{ms}` determines the maximum time of the colo-compare hold the packet. The `expired_scan_cycle=@var{ms}` is to set the period of scanning expired primary node network packets. The `max_queue_size=@var{size}` is to set the max compare queue size depend on user environment. If user want to use Xen COLO, need to add the notify_dev to notify Xen colo-frame to do checkpoint.

COLO-compare must be used with the help of filter-mirror, filter-redirector and filter-rewriter.

KVM COLO

primary:

(continues on next page)

(continued from previous page)

```

-netdev tap,id=hn0,vhost=off,script=/etc/qemu-ifup,downscript=/etc/qemu-ifdown
-device e1000,id=e0,netdev=hn0,mac=52:a4:00:12:78:66
-chardev socket,id=mirror0,host=3.3.3.3,port=9003,server,nowait
-chardev socket,id=compare1,host=3.3.3.3,port=9004,server,nowait
-chardev socket,id=compare0,host=3.3.3.3,port=9001,server,nowait
-chardev socket,id=compare0-0,host=3.3.3.3,port=9001
-chardev socket,id=compare_out,host=3.3.3.3,port=9005,server,nowait
-chardev socket,id=compare_out0,host=3.3.3.3,port=9005
-object iothread,id=iothread1
-object filter-mirror,id=m0,netdev=hn0,queue=tx,outdev=mirror0
-object filter-redirector,netdev=hn0,id=redire0,queue=rx,indev=compare_out
-object filter-redirector,netdev=hn0,id=redire1,queue=rx,outdev=compare0
-object colo-compare,id=comp0,primary_in=compare0-0,secondary_in=compare1,
↳outdev=compare_out0,iothread=iothread1

secondary:
-netdev tap,id=hn0,vhost=off,script=/etc/qemu-ifup,down script=/etc/qemu-
↳ifdown
-device e1000,netdev=hn0,mac=52:a4:00:12:78:66
-chardev socket,id=red0,host=3.3.3.3,port=9003
-chardev socket,id=red1,host=3.3.3.3,port=9004
-object filter-redirector,id=f1,netdev=hn0,queue=tx,indev=red0
-object filter-redirector,id=f2,netdev=hn0,queue=rx,outdev=red1

Xen COLO

primary:
-netdev tap,id=hn0,vhost=off,script=/etc/qemu-ifup,downscript=/etc/qemu-ifdown
-device e1000,id=e0,netdev=hn0,mac=52:a4:00:12:78:66
-chardev socket,id=mirror0,host=3.3.3.3,port=9003,server,nowait
-chardev socket,id=compare1,host=3.3.3.3,port=9004,server,nowait
-chardev socket,id=compare0,host=3.3.3.3,port=9001,server,nowait
-chardev socket,id=compare0-0,host=3.3.3.3,port=9001
-chardev socket,id=compare_out,host=3.3.3.3,port=9005,server,nowait
-chardev socket,id=compare_out0,host=3.3.3.3,port=9005
-chardev socket,id=notify_way,host=3.3.3.3,port=9009,server,nowait
-object filter-mirror,id=m0,netdev=hn0,queue=tx,outdev=mirror0
-object filter-redirector,netdev=hn0,id=redire0,queue=rx,indev=compare_out
-object filter-redirector,netdev=hn0,id=redire1,queue=rx,outdev=compare0
-object iothread,id=iothread1
-object colo-compare,id=comp0,primary_in=compare0-0,secondary_in=compare1,
↳outdev=compare_out0,notify_dev=nofity_way,iothread=iothread1

secondary:
-netdev tap,id=hn0,vhost=off,script=/etc/qemu-ifup,down script=/etc/qemu-
↳ifdown
-device e1000,netdev=hn0,mac=52:a4:00:12:78:66
-chardev socket,id=red0,host=3.3.3.3,port=9003
-chardev socket,id=red1,host=3.3.3.3,port=9004
-object filter-redirector,id=f1,netdev=hn0,queue=tx,indev=red0
-object filter-redirector,id=f2,netdev=hn0,queue=rx,outdev=red1

```

If you want to know the detail of above command line, you can read the colo-compare git log.

-object cryptodev-backend-builtin,id=id[,queues=queues] Creates a cryptodev backend which executes crypto operation from the QEMU cipher APIS. The id parameter is a unique ID

that will be used to reference this cryptodev backend from the `virtio-crypto` device. The `queues` parameter is optional, which specify the queue number of cryptodev backend, the default of `queues` is 1.

```
# qemu-system-x86_64 \
[...] \
    -object cryptodev-backend-builtin,id=cryptodev0 \
    -device virtio-crypto-pci,id=crypto0,cryptodev=cryptodev0 \
[...]
```

-object cryptodev-vhost-user,id=id,chardev=chardev[,queues=queues] Creates a vhost-user cryptodev backend, backed by a chardev chardev. The `id` parameter is a unique ID that will be used to reference this cryptodev backend from the `virtio-crypto` device. The chardev should be a unix domain socket backed one. The vhost-user uses a specifically defined protocol to pass vhost ioctl replacement messages to an application on the other end of the socket. The `queues` parameter is optional, which specify the queue number of cryptodev backend for multiqueue vhost-user, the default of `queues` is 1.

```
# qemu-system-x86_64 \
[...] \
    -chardev socket,id=chardev0,path=/path/to/socket \
    -object cryptodev-vhost-user,id=cryptodev0,chardev=chardev0 \
    -device virtio-crypto-pci,id=crypto0,cryptodev=cryptodev0 \
[...]
```

-object secret,id=id,data=string,format=raw|base64[,keyid=secretid,iv=string]

-object secret,id=id,file=filename,format=raw|base64[,keyid=secretid,iv=string]

Defines a secret to store a password, encryption key, or some other sensitive data. The sensitive data can either be passed directly via the `data` parameter, or indirectly via the `file` parameter. Using the `data` parameter is insecure unless the sensitive data is encrypted.

The sensitive data can be provided in raw format (the default), or base64. When encoded as JSON, the raw format only supports valid UTF-8 characters, so base64 is recommended for sending binary data. QEMU will convert from which ever format is provided to the format it needs internally. eg, an RBD password can be provided in raw format, even though it will be base64 encoded when passed onto the RBD sever.

For added protection, it is possible to encrypt the data associated with a secret using the AES-256-CBC cipher. Use of encryption is indicated by providing the `keyid` and `iv` parameters. The `keyid` parameter provides the ID of a previously defined secret that contains the AES-256 decryption key. This key should be 32-bytes long and be base64 encoded. The `iv` parameter provides the random initialization vector used for encryption of this particular secret and should be a base64 encrypted string of the 16-byte IV.

The simplest (insecure) usage is to provide the secret inline

```
# qemu-system-x86_64 -object secret,id=sec0,data=letmein,format=raw
```

The simplest secure usage is to provide the secret via a file

```
# printf "letmein" > mypasswd.txt # QEMU_SYSTEM_MACRO -object se-
cret,id=sec0,file=myspasswd.txt,format=raw
```

For greater security, AES-256-CBC should be used. To illustrate usage, consider the `openssl` command line tool which can encrypt the data. Note that when encrypting, the plaintext must be padded to the cipher block size (32 bytes) using the standard PKCS#5/6 compatible padding algorithm.

First a master key needs to be created in base64 encoding:

```
# openssl rand -base64 32 > key.b64
# KEY=$(base64 -d key.b64 | hexdump -v -e '/1 "%02X"')
```

Each secret to be encrypted needs to have a random initialization vector generated. These do not need to be kept secret

```
# openssl rand -base64 16 > iv.b64
# IV=$(base64 -d iv.b64 | hexdump -v -e '/1 "%02X"')
```

The secret to be defined can now be encrypted, in this case we're telling openssl to base64 encode the result, but it could be left as raw bytes if desired.

```
# SECRET=$(printf "letmein" |
    openssl enc -aes-256-cbc -a -K $KEY -iv $IV)
```

When launching QEMU, create a master secret pointing to `key.b64` and specify that to be used to decrypt the user password. Pass the contents of `iv.b64` to the second secret

```
# qemu-system-x86_64 \
    -object secret,id=secmaster0,format=base64,file=key.b64 \
    -object secret,id=sec0,keyid=secmaster0,format=base64,\
        data=$SECRET,iv=$(cat iv.b64)
-object sev-guest,id=id,cbitpos=cbitpos,reduced-phys-bits=val,[sev-device=string,policy=string]
```

Create a Secure Encrypted Virtualization (SEV) guest object, which can be used to provide the guest memory encryption support on AMD processors.

When memory encryption is enabled, one of the physical address bit (aka the C-bit) is utilized to mark if a memory page is protected. The `cbitpos` is used to provide the C-bit position. The C-bit position is Host family dependent hence user must provide this value. On EPYC, the value should be 47.

When memory encryption is enabled, we loose certain bits in physical address space. The `reduced-phys-bits` is used to provide the number of bits we loose in physical address space. Similar to C-bit, the value is Host family dependent. On EPYC, the value should be 5.

The `sev-device` provides the device file to use for communicating with the SEV firmware running inside AMD Secure Processor. The default device is `/dev/sev`. If hardware supports memory encryption then `/dev/sev` devices are created by CCP driver.

The `policy` provides the guest policy to be enforced by the SEV firmware and restrict what configuration and operational commands can be performed on this guest by the hypervisor. The policy should be provided by the guest owner and is bound to the guest and cannot be changed throughout the lifetime of the guest. The default is 0.

If guest `policy` allows sharing the key with another SEV guest then `handle` can be use to provide handle of the guest from which to share the key.

The `dh-cert-file` and `session-file` provides the guest owner's Public Diffie-Hillman key defined in SEV spec. The PDH and session parameters are used for establishing a cryptographic session with the guest owner to negotiate keys used for attestation. The file must be encoded in base64.

e.g to launch a SEV guest

```
# qemu_system-x86_64 \
    ..... \
    -object sev-guest,id=sev0,cbitpos=47,reduced-phys-bits=5 \
    -machine ...,memory-encryption=sev0 \
    .....
```

```
-object authz-simple,id=id,identity=string
```

 Create an authorization object that will control access to network services.

The `identity` parameter identifies the user and its format depends on the network service that authorization object is associated with. For authorizing based on TLS x509 certificates, the identity must be the x509 distinguished name. Note that care must be taken to escape any commas in the distinguished name.

An example authorization object to validate a x509 distinguished name would look like:

```
# qemu-system-x86_64 \
... \
  -object 'authz-simple,id=auth0,identity=CN=laptop.example.com,,
↪O=Example Org,,L=London,,ST=London,,C=GB' \
...

```

Note the use of quotes due to the x509 distinguished name containing whitespace, and escaping of `,`.

-object authz-listfile,id=id,filename=path,refresh=on|off Create an authorization object that will control access to network services.

The `filename` parameter is the fully qualified path to a file containing the access control list rules in JSON format.

An example set of rules that match against SASL usernames might look like:

```
{
  "rules": [
    { "match": "fred", "policy": "allow", "format": "exact" },
    { "match": "bob", "policy": "allow", "format": "exact" },
    { "match": "danb", "policy": "deny", "format": "glob" },
    { "match": "dan*", "policy": "allow", "format": "exact" },
  ],
  "policy": "deny"
}
```

When checking access the object will iterate over all the rules and the first rule to match will have its `policy` value returned as the result. If no rules match, then the default `policy` value is returned.

The rules can either be an exact string match, or they can use the simple UNIX glob pattern matching to allow wildcards to be used.

If `refresh` is set to true the file will be monitored and automatically reloaded whenever its content changes.

As with the `authz-simple` object, the format of the identity strings being matched depends on the network service, but is usually a TLS x509 distinguished name, or a SASL username.

An example authorization object to validate a SASL username would look like:

```
# qemu-system-x86_64 \
... \
  -object authz-simple,id=auth0,filename=/etc/qemu/vnc-sasl.acl,
↪refresh=on \
...

```

-object authz-pam,id=id,service=string Create an authorization object that will control access to network services.

The `service` parameter provides the name of a PAM service to use for authorization. It requires that a file `/etc/pam.d/service` exist to provide the configuration for the account subsystem.

An example authorization object to validate a TLS x509 distinguished name would look like:

```
# qemu-system-x86_64 \
... \

```

```
-object authz-pam,id=auth0,service=qemu-vnc \  
...
```

There would then be a corresponding config file for PAM at `/etc/pam.d/qemu-vnc` that contains:

```
account requisite pam_listfile.so item=user sense=allow \  
file=/etc/qemu/vnc.allow
```

Finally the `/etc/qemu/vnc.allow` file would contain the list of x509 distinguished names that are permitted access

```
CN=laptop.example.com,O=Example Home,L=London,ST=London,C=GB
```

-object iothread,id=id,poll-max-ns=poll-max-ns,poll-grow=poll-grow,poll-shrink=poll-shrink

Creates a dedicated event loop thread that devices can be assigned to. This is known as an IOThread. By default device emulation happens in vCPU threads or the main event loop thread. This can become a scalability bottleneck. IOThreads allow device emulation and I/O to run on other host CPUs.

The `id` parameter is a unique ID that will be used to reference this IOThread from `-device ...`, `iothread=id`. Multiple devices can be assigned to an IOThread. Note that not all devices support an `iothread` parameter.

The `query-iothreads` QMP command lists IOThreads and reports their thread IDs so that the user can configure host CPU pinning/affinity.

IOThreads use an adaptive polling algorithm to reduce event loop latency. Instead of entering a blocking system call to monitor file descriptors and then pay the cost of being woken up when an event occurs, the polling algorithm spins waiting for events for a short time. The algorithm's default parameters are suitable for many cases but can be adjusted based on knowledge of the workload and/or host device latency.

The `poll-max-ns` parameter is the maximum number of nanoseconds to busy wait for events. Polling can be disabled by setting this value to 0.

The `poll-grow` parameter is the multiplier used to increase the polling time when the algorithm detects it is missing events due to not polling long enough.

The `poll-shrink` parameter is the divisor used to decrease the polling time when the algorithm detects it is spending too long polling without encountering events.

The polling parameters can be modified at run-time using the `qom-set` command (where `iothread1` is the IOThread's id):

```
(qemu) qom-set /objects/iothread1 poll-max-ns 100000
```

1.2.12 Device URL Syntax

In addition to using normal file images for the emulated storage devices, QEMU can also use networked resources such as iSCSI devices. These are specified using a special URL syntax.

iSCSI iSCSI support allows QEMU to access iSCSI resources directly and use as images for the guest storage. Both disk and cdrom images are supported.

Syntax for specifying iSCSI LUNs is “iscsi://<target-ip>[:<port>]/<target-iqn>/<lun>”

By default qemu will use the iSCSI initiator-name ‘iqn.2008-11.org.linux-kvm[:<name>]’ but this can also be set from the command line or a configuration file.

Since version Qemu 2.4 it is possible to specify a iSCSI request timeout to detect stalled requests and force a reestablishment of the session. The timeout is specified in seconds. The default is 0 which means no timeout. Libiscsi 1.15.0 or greater is required for this feature.

Example (without authentication):

```
qemu-system-x86_64 -iscsi initiator-name=iqn.2001-04.com.example:my-
→initiator \
    -cdrom iscsi://192.0.2.1/iqn.2001-04.com.example/2 \
    -drive file=iscsi://192.0.2.1/iqn.2001-04.com.example/1
```

Example (CHAP username/password via URL):

```
qemu-system-x86_64 -drive file=iscsi://user%password@192.0.2.1/iqn.2001-
→04.com.example/1
```

Example (CHAP username/password via environment variables):

```
LIBISCSI_CHAP_USERNAME="user" \
LIBISCSI_CHAP_PASSWORD="password" \
qemu-system-x86_64 -drive file=iscsi://192.0.2.1/iqn.2001-04.com.example/1
```

NBD QEMU supports NBD (Network Block Devices) both using TCP protocol as well as Unix Domain Sockets. With TCP, the default port is 10809.

Syntax for specifying a NBD device using TCP, in preferred URI form: “nbd://<server-ip>[:<port>]/[<export>]”

Syntax for specifying a NBD device using Unix Domain Sockets; remember that ‘?’ is a shell glob character and may need quoting: “nbd+unix:///[:<export>]?socket=<domain-socket>”

Older syntax that is also recognized: “nbd:<server-ip>:<port>[:exportname=<export>]”

Syntax for specifying a NBD device using Unix Domain Sockets “nbd:unix:<domain-socket>[:exportname=<export>]”

Example for TCP

```
qemu-system-x86_64 --drive file=nbd:192.0.2.1:30000
```

Example for Unix Domain Sockets

```
qemu-system-x86_64 --drive file=nbd:unix:/tmp/nbd-socket
```

SSH QEMU supports SSH (Secure Shell) access to remote disks.

Examples:

```
qemu-system-x86_64 -drive file=ssh://user@host/path/to/disk.img
qemu-system-x86_64 -drive file.driver=ssh,file.user=user,file.host=host,
→file.port=22,file.path=/path/to/disk.img
```

Currently authentication must be done using ssh-agent. Other authentication methods may be supported in future.

Sheepdog Sheepdog is a distributed storage system for QEMU. QEMU supports using either local sheepdog devices or remote networked devices.

Syntax for specifying a sheepdog device

```
sheepdog[+tcp|+unix]://[host:port]/vdiname[?socket=path][#snapid|tag]
```

Example

```
qemu-system-x86_64 --drive file=sheepdog://192.0.2.1:30000/
→MyVirtualMachine
```

See also <https://sheepdog.github.io/sheepdog/>.

GlusterFS GlusterFS is a user space distributed file system. QEMU supports the use of GlusterFS volumes for hosting VM disk images using TCP, Unix Domain Sockets and RDMA transport protocols.

Syntax for specifying a VM disk image on GlusterFS volume is

```
URI:
gluster[+type]://[host[:port]]/volume/path[?socket=...][,debug=N][,logfile=...]

JSON:
'json:{"driver":"qcow2","file":{"driver":"gluster","volume":"testvol","path":"a.
→img","debug":N,"logfile":"...",
                                "server":[{"type":"tcp","host":"...", "port":"...
→"},
                                {"type":"unix","socket":"..."}]}}'
```

Example

```
URI:
qemu-system-x86_64 --drive file=gluster://192.0.2.1/testvol/a.img,
                    file.debug=9,file.logfile=/var/log/qemu-
→gluster.log

JSON:
qemu-system-x86_64 'json:{"driver":"qcow2",
                        "file":{"driver":"gluster",
                                "volume":"testvol","path":"a.img",
                                "debug":9,"logfile":"/var/log/qemu-
→gluster.log",
                                "server":[{"type":"tcp","host":"1.2.3.
→4","port":24007},
                                {"type":"unix","socket":"/
→var/run/glusterd.socket"}]}}'
qemu-system-x86_64 -drive driver=qcow2,file.driver=gluster,file.
→volume=testvol,file.path=/path/a.img,
                    file.debug=9,file.logfile=/var/log/
→qemu-gluster.log,
                    file.server.0.type=tcp,file.server.
→0.host=1.2.3.4,file.server.0.port=24007,
                    file.server.1.type=unix,file.server.
→1.socket=/var/run/glusterd.socket
```

See also <http://www.gluster.org>.

HTTP/HTTPS/FTP/FTPS QEMU supports read-only access to files accessed over http(s) and ftp(s).

Syntax using a single filename:

```
<protocol>://[<username>[:<password>]@]<host>/<path>
```

where:

protocol 'http', 'https', 'ftp', or 'ftps'.

username Optional username for authentication to the remote server.

password Optional password for authentication to the remote server.

host Address of the remote server.

path Path on the remote server, including any query string.

The following options are also supported:

url The full URL when passing options to the driver explicitly.

readahead The amount of data to read ahead with each range request to the remote server. This value may optionally have the suffix 'T', 'G', 'M', 'K', 'k' or 'b'. If it does not have a suffix, it will be assumed to be in bytes. The value must be a multiple of 512 bytes. It defaults to 256k.

sslverify Whether to verify the remote server's certificate when connecting over SSL. It can have the value 'on' or 'off'. It defaults to 'on'.

cookie Send this cookie (it can also be a list of cookies separated by ';') with each outgoing request. Only supported when using protocols such as HTTP which support cookies, otherwise ignored.

timeout Set the timeout in seconds of the CURL connection. This timeout is the time that CURL waits for a response from the remote server to get the size of the image to be downloaded. If not set, the default timeout of 5 seconds is used.

Note that when passing options to qemu explicitly, `driver` is the value of `<protocol>`.

Example: boot from a remote Fedora 20 live ISO image

```
qemu_system-x86_64 --drive media=cdrom,file=https://archives.
↳ fedoraproject.org/pub/archive/fedora/linux/releases/20/Live/x86_64/
↳ Fedora-Live-Desktop-x86_64-20-1.iso,readonly
```

```
qemu_system-x86_64 --drive media=cdrom,file.driver=http,file.url=http:/
↳ /archives.fedoraproject.org/pub/fedora/linux/releases/20/Live/x86_64/
↳ Fedora-Live-Desktop-x86_64-20-1.iso,readonly
```

Example: boot from a remote Fedora 20 cloud image using a local overlay for writes, copy-on-read, and a readahead of 64k

```
qemu-img create -f qcow2 -o backing_file='json:{"file.driver":"http",,
↳ "file.url":"http://archives.fedoraproject.org/pub/archive/fedora/linux/
releases/20/Images/x86_64/Fedora-x86_64-20-20131211.1-sda.qcow2",, "file.
↳ readahead":"64k"}' /tmp/Fedora-x86_64-20-20131211.1-sda.qcow2
```

```
qemu_system-x86_64 -drive file=/tmp/Fedora-x86_64-20-20131211.1-sda.qcow2,
↳ copy-on-read=on
```

Example: boot from an image stored on a VMware vSphere server with a self-signed certificate using a local overlay for writes, a readahead of 64k and a timeout of 10 seconds.

```
qemu-img create -f qcow2 -o backing_file='json:{"file.driver":"https",
↳ , "file.url":"https://user:password@vsphere.example.com/folder/
test/test-flat.vmdk?dcPath=Datacenter&dsName=datastore1",, "file.
↳ sslverify":"off",, "file.readahead":"64k",, "file.timeout":10}' /tmp/
↳ test.qcow2
```

```
qemu_system-x86_64 -drive file=/tmp/test.qcow2
```

1.3 Keys in the graphical frontends

During the graphical emulation, you can use special key combinations to change modes. The default key mappings are shown below, but if you use `-alt-grab` then the modifier is Ctrl-Alt-Shift (instead of Ctrl-Alt) and if you use `-ctrl-grab` then the modifier is the right Ctrl key (instead of Ctrl-Alt):

Ctrl-Alt-f Toggle full screen

Ctrl-Alt-+ Enlarge the screen

Ctrl-Alt- Shrink the screen

Ctrl-Alt-u Restore the screen's un-scaled dimensions

Ctrl-Alt-n Switch to virtual console 'n'. Standard console mappings are:

- 1 Target system display
- 2 Monitor
- 3 Serial port

Ctrl-Alt Toggle mouse and keyboard grab.

In the virtual consoles, you can use Ctrl-Up, Ctrl-Down, Ctrl-PageUp and Ctrl-PageDown to move in the back log.

1.4 Keys in the character backend multiplexer

During emulation, if you are using a character backend multiplexer (which is the default if you are using `-nographic`) then several commands are available via an escape sequence. These key sequences all start with an escape character, which is Ctrl-a by default, but can be changed with `-echr`. The list below assumes you're using the default.

Ctrl-a h Print this help

Ctrl-a x Exit emulator

Ctrl-a s Save disk data back to file (if `-snapshot`)

Ctrl-a t Toggle console timestamps

Ctrl-a b Send break (magic `sysrq` in Linux)

Ctrl-a c Rotate between the frontends connected to the multiplexer (usually this switches between the monitor and the console)

Ctrl-a Ctrl-a Send the escape character to the frontend

1.5 QEMU Monitor

The QEMU monitor is used to give complex commands to the QEMU emulator. You can use it to:

- Remove or insert removable media images (such as CD-ROM or floppies).
- Freeze/unfreeze the Virtual Machine (VM) and save or restore its state from a disk file.
- Inspect the VM state without an external debugger.

1.5.1 Commands

The following commands are available:

help or ? [cmd] Show the help for all commands or just for command *cmd*.

commit Commit changes to the disk images (if `-snapshot` is used) or backing files. If the backing file is smaller than the snapshot, then the backing file will be resized to be the same size as the snapshot. If the snapshot is smaller than the backing file, the backing file will not be truncated. If you want the backing file to match the size of the smaller snapshot, you can safely truncate it yourself once the commit operation successfully completes.

quit or q Quit the emulator.

exit_preconfig This command makes QEMU exit the preconfig state and proceed with VM initialization using configuration data provided on the command line and via the QMP monitor during the preconfig state. The command is only available during the preconfig state (i.e. when the `-preconfig` command line option was in use).

block_resize Resize a block image while a guest is running. Usually requires guest action to see the updated size. Resize to a lower size is supported, but should be used with extreme caution. Note that this command only resizes image files, it can not resize block devices like LVM volumes.

block_stream Copy data from a backing file into a block device.

block_job_set_speed Set maximum speed for a background block operation.

block_job_cancel Stop an active background block operation (streaming, mirroring).

block_job_complete Manually trigger completion of an active background block operation. For mirroring, this will switch the device to the destination path.

block_job_pause Pause an active block streaming operation.

block_job_resume Resume a paused block streaming operation.

eject [-f] device Eject a removable medium (use -f to force it).

drive_del device Remove host block device. The result is that guest generated IO is no longer submitted against the host device underlying the disk. Once a drive has been deleted, the QEMU Block layer returns -EIO which results in IO errors in the guest for applications that are reading/writing to the device. These errors are always reported to the guest, regardless of the drive's error actions (drive options `error`, `werror`).

change device setting Change the configuration of a device.

change diskdevice filename [format [read-only-mode]] Change the medium for a removable disk device to point to *filename*. eg:

```
(qemu) change ide1-cd0 /path/to/some.iso
```

format is optional.

read-only-mode may be used to change the read-only status of the device. It accepts the following values:

retain Retains the current status; this is the default.

read-only Makes the device read-only.

read-write Makes the device writable.

change vnc password [password]

Change the password associated with the VNC server. If the new password is not supplied, the monitor will prompt for it to be entered. VNC passwords are only significant up to 8 letters. eg:

```
(qemu) change vnc password
Password: ****
```

screendump filename Save screen into PPM image *filename*.

logfile filename Output logs to *filename*.

trace-event changes status of a trace event

trace-file on|off|flush Open, close, or flush the trace file. If no argument is given, the status of the trace file is displayed.

log itemI[,...] Activate logging of the specified items.

savevm *tag* Create a snapshot of the whole virtual machine. If *tag* is provided, it is used as human readable identifier. If there is already a snapshot with the same tag, it is replaced. More info at [VM snapshots](#).

Since 4.0, savevm stopped allowing the snapshot id to be set, accepting only *tag* as parameter.

loadvm *tag* Set the whole virtual machine to the snapshot identified by the tag *tag*.

Since 4.0, loadvm stopped accepting snapshot id as parameter.

delvm *tag* Delete the snapshot identified by *tag*.

Since 4.0, delvm stopped deleting snapshots by snapshot id, accepting only *tag* as parameter.

singlestep [*off*] Run the emulation in single step mode. If called with option off, the emulation returns to normal mode.

stop Stop emulation.

cont or c Resume emulation.

system_wakeup Wakeup guest from suspend.

gdbserver [*port*] Start gdbserver session (default *port*=1234)

x/*fmt addr* Virtual memory dump starting at *addr*.

xp /*fmt addr* Physical memory dump starting at *addr*.

fmt is a format which tells the command how to format the data. Its syntax is: /{*count*}{*format*}{*size*}

count is the number of items to be dumped.

format can be x (hex), d (signed decimal), u (unsigned decimal), o (octal), c (char) or i (asm instruction).

size can be b (8 bits), h (16 bits), w (32 bits) or g (64 bits). On x86, h or w can be specified with the i format to respectively select 16 or 32 bit code instruction size.

Examples:

Dump 10 instructions at the current instruction pointer:

```
(qemu) x/10i $eip
0x90107063: ret
0x90107064: sti
0x90107065: lea    0x0(%esi,1),%esi
0x90107069: lea    0x0(%edi,1),%edi
0x90107070: ret
0x90107071: jmp    0x90107080
0x90107073: nop
0x90107074: nop
0x90107075: nop
0x90107076: nop
```

Dump 80 16 bit values at the start of the video memory:

```
(qemu) xp/80hx 0xb8000
0x000b8000: 0x0b50 0x0b6c 0x0b65 0x0b78 0x0b38 0x0b36 0x0b2f 0x0b42
0x000b8010: 0x0b6f 0x0b63 0x0b68 0x0b73 0x0b20 0x0b56 0x0b47 0x0b41
0x000b8020: 0x0b42 0x0b69 0x0b6f 0x0b73 0x0b20 0x0b63 0x0b75 0x0b72
0x000b8030: 0x0b72 0x0b65 0x0b6e 0x0b74 0x0b2d 0x0b63 0x0b76 0x0b73
0x000b8040: 0x0b20 0x0b30 0x0b35 0x0b20 0x0b4e 0x0b6f 0x0b76 0x0b20
0x000b8050: 0x0b32 0x0b30 0x0b30 0x0b33 0x0720 0x0720 0x0720 0x0720
0x000b8060: 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720
0x000b8070: 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720
```

(continues on next page)

(continued from previous page)

```
0x000b8080: 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720
0x000b8090: 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720 0x0720
```

gpa2hva *addr* Print the host virtual address at which the guest's physical address *addr* is mapped.

gpa2hpa *addr* Print the host physical address at which the guest's physical address *addr* is mapped.

gva2gpa *addr* Print the guest physical address at which the guest's virtual address *addr* is mapped based on the mapping for the current CPU.

print or p/*fmt* *expr* Print expression value. Only the *format* part of *fmt* is used.

i/*fmt* *addr* [*index*] Read I/O port.

o/*fmt* *addr* *val* Write to I/O port.

sendkey *keys* Send *keys* to the guest. *keys* could be the name of the key or the raw value in hexadecimal format. Use **-** to press several keys simultaneously. Example:

```
sendkey ctrl-alt-f1
```

This command is useful to send keys that your graphical user interface intercepts at low level, such as **ctrl-alt-f1** in X Window.

sync-profile [*on|off|reset*] Enable, disable or reset synchronization profiling. With no arguments, prints whether profiling is on or off.

system_reset Reset the system.

system_powerdown Power down the system (if supported).

sum *addr* *size* Compute the checksum of a memory region.

device_add *config* Add device.

device_del *id* Remove device *id*. *id* may be a short ID or a QOM object path.

cpu *index* Set the default CPU.

mouse_move *dx* *dy* [*dz*] Move the active mouse to the specified coordinates *dx* *dy* with optional scroll axis *dz*.

mouse_button *val* Change the active mouse button state *val* (1=L, 2=M, 4=R).

mouse_set *index* Set which mouse device receives events at given *index*, index can be obtained with:

```
info mice
```

wavcapture *filename* *audiodev* [*frequency* [*bits* [*channels*]]] Capture audio into *filename* from *audiodev*, using sample rate *frequency* bits per sample *bits* and number of channels *channels*.

Defaults:

- Sample rate = 44100 Hz - CD quality
- Bits = 16
- Number of channels = 2 - Stereo

stopcapture *index* Stop capture with a given *index*, index can be obtained with:

```
info capture
```

memsave *addr* *size* *file* save to disk virtual memory dump starting at *addr* of size *size*.

pmemsave *addr* *size* *file* save to disk physical memory dump starting at *addr* of size *size*.

boot_set bootdevicelist Define new values for the boot device list. Those values will override the values specified on the command line through the `-boot` option.

The values that can be specified here depend on the machine type, but are the same that can be specified in the `-boot` command line option.

nmi cpu Inject an NMI on the default CPU (x86/s390) or all CPUs (ppc64).

ringbuf_write device data Write *data* to ring buffer character device *device*. *data* must be a UTF-8 string.

ringbuf_read device Read and print up to *size* bytes from ring buffer character device *device*. Certain non-printable characters are printed `\uXXXX`, where `XXXX` is the character code in hexadecimal. Character `\` is printed `\\`. Bug: can screw up when the buffer contains invalid UTF-8 sequences, NUL characters, after the ring buffer lost data, and when reading stops because the size limit is reached.

announce_self Trigger a round of GARP/RARP broadcasts; this is useful for explicitly updating the network infrastructure after a reconfiguration or some forms of migration. The timings of the round are set by the migration announce parameters. An optional comma separated *interfaces* list restricts the announce to the named set of interfaces. An optional *id* can be used to start a separate announce timer and to change the parameters of it later.

migrate [-d] [-b] [-i] uri Migrate to *uri* (using `-d` to not wait for completion).

`-b` for migration with full copy of disk

`-i` for migration with incremental copy of disk (base image is shared)

migrate_cancel Cancel the current VM migration.

migrate_continue state Continue migration from the paused state *state*

migrate_incoming uri Continue an incoming migration using the *uri* (that has the same syntax as the `-incoming` option).

migrate_recover uri Continue a paused incoming postcopy migration using the *uri*.

migrate_pause Pause an ongoing migration. Currently it only supports postcopy.

migrate_set_cache_size value Set cache size to *value* (in bytes) for xbzrl migrations.

migrate_set_speed value Set maximum speed to *value* (in bytes) for migrations.

migrate_set_downtime second Set maximum tolerated downtime (in seconds) for migration.

migrate_set_capability capability state Enable/Disable the usage of a capability *capability* for migration.

migrate_set_parameter parameter value Set the parameter *parameter* for migration.

migrate_start_postcopy Switch in-progress migration to postcopy mode. Ignored after the end of migration (or once already in postcopy).

x_colo_lost_heartbeat Tell COLO that heartbeat is lost, a failover or takeover is needed.

client_migrate_info protocol hostname port tls-port cert-subject Set migration information for remote display. This makes the server ask the client to automatically reconnect using the new parameters once migration finished successfully. Only implemented for SPICE.

dump-guest-memory [-p] filename begin length

dump-guest-memory [-z|-l|-s|-w] filename Dump guest memory to *protocol*. The file can be processed with crash or gdb. Without `-z|-l|-s|-w`, the dump format is ELF.

`-p` do paging to get guest's memory mapping.

`-z` dump in kdump-compressed format, with zlib compression.

`-l` dump in kdump-compressed format, with lzo compression.

- s** dump in kdump-compressed format, with snappy compression.
- w** dump in Windows crashdump format (can be used instead of ELF-dump converting), for Windows x64 guests with vmcoreinfo driver only

filename dump file name.

begin the starting physical address. It's optional, and should be specified together with *length*.

length the memory size, in bytes. It's optional, and should be specified together with *begin*.

dump-skeys filename Save guest storage keys to a file.

migration_mode mode Enables or disables migration mode.

snapshot_blkdev Snapshot device, using snapshot file as target if provided

snapshot_blkdev_internal Take an internal snapshot on device if it support

snapshot_delete_blkdev_internal Delete an internal snapshot on device if it support

drive_mirror Start mirroring a block device's writes to a new destination, using the specified target.

drive_backup Start a point-in-time copy of a block device to a specified target.

drive_add Add drive to PCI storage controller.

pcie_aer_inject_error Inject PCIe AER error

netdev_add Add host network device.

netdev_del Remove host network device.

object_add Create QOM object.

object_del Destroy QOM object.

hostfwd_add Redirect TCP or UDP connections from host to guest (requires -net user).

hostfwd_remove Remove host-to-guest TCP or UDP redirection.

balloon value Request VM to change its memory allocation to *value* (in MB).

set_link name [on|off] Switch link *name* on (i.e. up) or off (i.e. down).

watchdog_action Change watchdog action.

acl_show aclname List all the matching rules in the access control list, and the default policy. There are currently two named access control lists, *vnc.x509dname* and *vnc.username* matching on the x509 client certificate distinguished name, and SASL username respectively.

acl_policy aclname allow|deny Set the default access control list policy, used in the event that none of the explicit rules match. The default policy at startup is always *deny*.

acl_add aclname match allow|deny [index] Add a match rule to the access control list, allowing or denying access. The match will normally be an exact username or x509 distinguished name, but can optionally include wildcard globs. eg **@EXAMPLE.COM* to allow all users in the *EXAMPLE.COM* kerberos realm. The match will normally be appended to the end of the ACL, but can be inserted earlier in the list if the optional *index* parameter is supplied.

acl_remove aclname match Remove the specified match rule from the access control list.

acl_reset aclname Remove all matches from the access control list, and set the default policy back to *deny*.

nbd_server_start host:port Start an NBD server on the given host and/or port. If the *-a* option is included, all of the virtual machine's block devices that have an inserted media on them are automatically exported; in this case, the *-w* option makes the devices writable too.

nbd_server_add device [*name*] Export a block device through QEMU's NBD server, which must be started beforehand with `nbd_server_start`. The `-w` option makes the exported device writable too. The export name is controlled by *name*, defaulting to *device*.

nbd_server_remove [`-f`] *name* Stop exporting a block device through QEMU's NBD server, which was previously started with `nbd_server_add`. The `-f` option forces the server to drop the export immediately even if clients are connected; otherwise the command fails unless there are no clients.

nbd_server_stop Stop the QEMU embedded NBD server.

mce cpu bank status mcgstatus addr misc Inject an MCE on the given CPU (x86 only).

getfd fdname If a file descriptor is passed alongside this command using the SCM_RIGHTS mechanism on unix sockets, it is stored using the name *fdname* for later use by other monitor commands.

closefd fdname Close the file descriptor previously assigned to *fdname* using the `getfd` command. This is only needed if the file descriptor was never used by another monitor command.

block_passwd device password Set the encrypted device *device* password to *password*

This command is now obsolete and will always return an error since 2.10

block_set_io_throttle device bps bps_rd bps_wr iops iops_rd iops_wr Change I/O throttle limits for a block drive to *bps bps_rd bps_wr iops iops_rd iops_wr*. *device* can be a block device name, a qdev ID or a QOM path.

set_password [*vnc* | *spice*] **password** [*action-if-connected*] Change spice/vnc password. Use zero to make the password stay valid forever. *action-if-connected* specifies what should happen in case a connection is established: *fail* makes the password change fail. *disconnect* changes the password and disconnects the client. *keep* changes the password and keeps the connection up. *keep* is the default.

expire_password [*vnc* | *spice*] **expire-time** Specify when a password for spice/vnc becomes invalid. *expire-time* accepts:

now Invalidate password instantly.

never Password stays valid forever.

+nsec Password stays valid for *nsec* seconds starting now.

nsec Password is invalidated at the given time. *nsec* are the seconds passed since 1970, i.e. unix epoch.

chardev-add args `chardev-add` accepts the same parameters as the `-chardev` command line switch.

chardev-change args `chardev-change` accepts existing chardev *id* and then the same arguments as the `-chardev` command line switch (except for "id").

chardev-remove id Removes the chardev *id*.

chardev-send-break id Send a break on the chardev *id*.

qemu-io device command Executes a `qemu-io` command on the given block device.

qom-list [path] Print QOM properties of object at location *path*

qom-get path property Print QOM property *property* of object at location *path*

qom-set path property value Set QOM property *property* of object at location *path* to value *value*

replay_break icount Set replay breakpoint at instruction count *icount*. Execution stops when the specified instruction is reached. There can be at most one breakpoint. When breakpoint is set, any prior one is removed. The breakpoint may be set only in replay mode and only "in the future", i.e. at instruction counts greater than the current one. The current instruction count can be observed with `info replay`.

replay_delete_break Remove replay breakpoint which was previously set with `replay_break`. The command is ignored when there are no replay breakpoints.

replay_seek *icount* Automatically proceed to the instruction count *icount*, when replaying the execution. The command automatically loads nearest snapshot and replays the execution to find the desired instruction. When there is no preceding snapshot or the execution is not replayed, then the command fails. *icount* for the reference may be observed with `info replay` command.

info subcommand Show various information about the system state.

info version Show the version of QEMU.

info network Show the network state.

info chardev Show the character devices.

info block Show info of one block device or all block devices.

info blockstats Show block device statistics.

info block-jobs Show progress of ongoing block device operations.

info registers Show the cpu registers.

info lapic Show local APIC state

info ioapic Show io APIC state

info cpus Show infos for each CPU.

info history Show the command line history.

info irq Show the interrupts statistics (if available).

info pic Show PIC state.

info rdma Show RDMA state.

info pci Show PCI information.

info tlb Show virtual to physical memory mappings.

info mem Show the active virtual memory mappings.

info mtree Show memory tree.

info jit Show dynamic compiler info.

info opcount Show dynamic compiler opcode counters

info sync-profile [-m|-n] [*max*] Show synchronization profiling info, up to *max* entries (default: 10), sorted by total wait time.

-m sort by mean wait time

-n do not coalesce objects with the same call site

When different objects that share the same call site are coalesced, the “Object” field shows—enclosed in brackets—the number of objects being coalesced.

info kvm Show KVM information.

info numa Show NUMA information.

info usb Show guest USB devices.

info usbhost Show host USB devices.

info profile Show profiling information.

info capture Show capture information.

info snapshots Show the currently saved VM snapshots.

info status Show the current VM status (running|paused).

info mice Show which guest mouse is receiving events.

info vnc Show the vnc server status.

info spice Show the spice server status.

info name Show the current VM name.

info uuid Show the current VM UUID.

info cpustats Show CPU statistics.

info usernet Show user network stack connection states.

info migrate Show migration status.

info migrate_capabilities Show current migration capabilities.

info migrate_parameters Show current migration parameters.

info migrate_cache_size Show current migration xbzrle cache size.

info balloon Show balloon information.

info qtree Show device tree.

info qdm Show qdev device model list.

info qom-tree Show QOM composition tree.

info roms Show roms.

info trace-events Show available trace-events & their state.

info tpm Show the TPM device.

info memdev Show memory backends

info memory-devices Show memory devices.

info iothreads Show iothread's identifiers.

info rocker name Show rocker switch.

info rocker-ports name-ports Show rocker ports.

info rocker-of-dpa-flows name [tbl_id] Show rocker OF-DPA flow tables.

info rocker-of-dpa-groups name [type] Show rocker OF-DPA groups.

info skeys address Display the value of a storage key (s390 only)

info cmma address Display the values of the CMMA storage attributes for a range of pages (s390 only)

info dump Display the latest dump status.

info ramblock Dump all the ramblocks of the system.

info hotpluggable-cpus Show information about hotpluggable CPUs

info vm-generation-id Show Virtual Machine Generation ID

info memory_size_summary Display the amount of initially allocated and present hotpluggable (if enabled) memory in bytes.

info sev Show SEV information.

info replay Display the record/replay information: mode and the current icount.

1.5.2 Integer expressions

The monitor understands integers expressions for every integer argument. You can use register names to get the value of specific CPU registers by prefixing them with \$.

1.6 Disk Images

QEMU supports many disk image formats, including growable disk images (their size increase as non empty sectors are written), compressed and encrypted disk images.

1.6.1 Quick start for disk image creation

You can create a disk image with the command:

```
qemu-img create myimage.img mysize
```

where myimage.img is the disk image filename and mysize is its size in kilobytes. You can add an M suffix to give the size in megabytes and a G suffix for gigabytes.

See the qemu-img invocation documentation for more information.

1.6.2 Snapshot mode

If you use the option `-snapshot`, all disk images are considered as read only. When sectors are written, they are written in a temporary file created in `/tmp`. You can however force the write back to the raw disk images by using the `commit` monitor command (or `C-a s` in the serial console).

1.6.3 VM snapshots

VM snapshots are snapshots of the complete virtual machine including CPU state, RAM, device state and the content of all the writable disks. In order to use VM snapshots, you must have at least one non removable and writable block device using the `qcow2` disk image format. Normally this device is the first virtual hard drive.

Use the monitor command `savevm` to create a new VM snapshot or replace an existing one. A human readable name can be assigned to each snapshot in addition to its numerical ID.

Use `loadvm` to restore a VM snapshot and `delvm` to remove a VM snapshot. `info snapshots` lists the available snapshots with their associated information:

```
(qemu) info snapshots
Snapshot devices: hda
Snapshot list (from hda):
ID      TAG      VM SIZE      DATE      VM CLOCK
1        start    41M 2006-08-06 12:38:02 00:00:14.954
2        40M 2006-08-06 12:43:29 00:00:18.633
3        msys     40M 2006-08-06 12:44:04 00:00:23.514
```

A VM snapshot is made of a VM state info (its size is shown in `info snapshots`) and a snapshot of every writable disk image. The VM state info is stored in the first `qcow2` non removable and writable block device. The disk image snapshots are stored in every disk image. The size of a snapshot in a disk image is difficult to evaluate and is not shown by `info snapshots` because the associated disk sectors are shared among all the snapshots to save disk space (otherwise each snapshot would need a full copy of all the disk images).

When using the (unrelated) `-snapshot` option (*Snapshot mode*), you can always make VM snapshots, but they are deleted as soon as you exit QEMU.

VM snapshots currently have the following known limitations:

- They cannot cope with removable devices if they are removed or inserted after a snapshot is done.
- A few device drivers still have incomplete snapshot support so their state is not saved or restored properly (in particular USB).

1.6.4 Disk image file formats

QEMU supports many image file formats that can be used with VMs as well as with any of the tools (like `qemu-img`). This includes the preferred formats `raw` and `qcow2` as well as formats that are supported for compatibility with older QEMU versions or other hypervisors.

Depending on the image format, different options can be passed to `qemu-img create` and `qemu-img convert` using the `-o` option. This section describes each format and the options that are supported for it.

raw

Raw disk image format. This format has the advantage of being simple and easily exportable to all other emulators. If your file system supports *holes* (for example in `ext2` or `ext3` on Linux or NTFS on Windows), then only the written sectors will reserve space. Use `qemu-img info` to know the real size used by the image or `ls -ls` on Unix/Linux.

Supported options:

preallocation

Preallocation mode (allowed values: `off`, `falloc`, `full`). `falloc` mode preallocates space for image by calling `posix_fallocate()`. `full` mode preallocates space for image by writing data to underlying storage. This data may or may not be zero, depending on the storage location.

qcow2

QEMU image format, the most versatile format. Use it to have smaller images (useful if your filesystem does not supports holes, for example on Windows), zlib based compression and support of multiple VM snapshots.

Supported options:

compat

Determines the `qcow2` version to use. `compat=0.10` uses the traditional image format that can be read by any QEMU since 0.10. `compat=1.1` enables image format extensions that only QEMU 1.1 and newer understand (this is the default). Amongst others, this includes zero clusters, which allow efficient copy-on-read for sparse images.

backing_file

File name of a base image (see `create` subcommand)

backing_fmt

Image format of the base image

encryption

This option is deprecated and equivalent to `encrypt.format=aes`

encrypt.format

If this is set to `luks`, it requests that the `qcow2` payload (not `qcow2` header) be encrypted using the LUKS format. The passphrase to use to unlock the LUKS key slot is given by the `encrypt.key-secret` parameter. LUKS encryption parameters can be tuned with the other `encrypt.*` parameters.

If this is set to `aes`, the image is encrypted with 128-bit AES-CBC. The encryption key is given by the `encrypt.key-secret` parameter. This encryption format is considered to be flawed by modern cryptography standards, suffering from a number of design problems:

- The AES-CBC cipher is used with predictable initialization vectors based on the sector number. This makes it vulnerable to chosen plaintext attacks which can reveal the existence of encrypted data.
- The user passphrase is directly used as the encryption key. A poorly chosen or short passphrase will compromise the security of the encryption.
- In the event of the passphrase being compromised there is no way to change the passphrase to protect data in any qcow images. The files must be cloned, using a different encryption passphrase in the new file. The original file must then be securely erased using a program like shred, though even this is ineffective with many modern storage technologies.

The use of this is no longer supported in system emulators. Support only remains in the command line utilities, for the purposes of data liberation and interoperability with old versions of QEMU. The `luks` format should be used instead.

`encrypt.key-secret`

Provides the ID of a `secret` object that contains the passphrase (`encrypt.format=luks`) or encryption key (`encrypt.format=aes`).

`encrypt.cipher-alg`

Name of the cipher algorithm and key length. Currently defaults to `aes-256`. Only used when `encrypt.format=luks`.

`encrypt.cipher-mode`

Name of the encryption mode to use. Currently defaults to `xts`. Only used when `encrypt.format=luks`.

`encrypt.ivgen-alg`

Name of the initialization vector generator algorithm. Currently defaults to `plain64`. Only used when `encrypt.format=luks`.

`encrypt.ivgen-hash-alg`

Name of the hash algorithm to use with the initialization vector generator (if required). Defaults to `sha256`. Only used when `encrypt.format=luks`.

`encrypt.hash-alg`

Name of the hash algorithm to use for PBKDF algorithm Defaults to `sha256`. Only used when `encrypt.format=luks`.

`encrypt.iter-time`

Amount of time, in milliseconds, to use for PBKDF algorithm per key slot. Defaults to 2000. Only used when `encrypt.format=luks`.

`cluster_size`

Changes the qcow2 cluster size (must be between 512 and 2M). Smaller cluster sizes can improve the image file size whereas larger cluster sizes generally provide better performance.

`preallocation`

Preallocation mode (allowed values: `off`, `metadata`, `falloc`, `full`). An image with preallocated metadata is initially larger but can improve performance when the image needs to grow. `falloc` and `full` preallocations are like the same options of `raw` format, but sets up metadata also.

`lazy_refcounts`

If this option is set to `on`, reference count updates are postponed with the goal of avoiding metadata I/O and improving performance. This is particularly interesting with `cache=writethrough` which doesn't batch metadata updates. The tradeoff is that after a host crash, the reference count tables must be rebuilt, i.e. on the next open an (automatic) `qemu-img check -r all` is required, which may take some time.

This option can only be enabled if `compat=1.1` is specified.

nocow

If this option is set to `on`, it will turn off COW of the file. It's only valid on `btrfs`, no effect on other file systems.

`Btrfs` has low performance when hosting a VM image file, even more when the guest on the VM also using `btrfs` as file system. Turning off COW is a way to mitigate this bad performance. Generally there are two ways to turn off COW on `btrfs`:

- Disable it by mounting with `nodatacow`, then all newly created files will be `NOCOW`.
- For an empty file, add the `NOCOW` file attribute. That's what this option does.

Note: this option is only valid to new or empty files. If there is an existing file which is COW and has data blocks already, it couldn't be changed to `NOCOW` by setting `nocow=on`. One can issue `lsattr filename` to check if the `NOCOW` flag is set or not (Capital 'C' is `NOCOW` flag).

qed

Old QEMU image format with support for backing files and compact image files (when your filesystem or transport medium does not support holes).

When converting QED images to `qcow2`, you might want to consider using the `lazy_refcounts=on` option to get a more QED-like behaviour.

Supported options:

backing_file

File name of a base image (see `create` subcommand).

backing_fmt

Image file format of backing file (optional). Useful if the format cannot be autodetected because it has no header, like some `vhd/vpc` files.

cluster_size

Changes the cluster size (must be power-of-2 between 4K and 64K). Smaller cluster sizes can improve the image file size whereas larger cluster sizes generally provide better performance.

table_size

Changes the number of clusters per L1/L2 table (must be power-of-2 between 1 and 16). There is normally no need to change this value but this option can be used for performance benchmarking.

qcow

Old QEMU image format with support for backing files, compact image files, encryption and compression.

Supported options:

backing_file

File name of a base image (see `create` subcommand)

encryption

This option is deprecated and equivalent to `encrypt.format=aes`

encrypt.format

If this is set to `aes`, the image is encrypted with 128-bit AES-CBC. The encryption key is given by the `encrypt.key-secret` parameter. This encryption format is considered to be flawed by modern cryptography standards, suffering from a number of design problems enumerated previously against the `qcow2` image format.

The use of this is no longer supported in system emulators. Support only remains in the command line utilities, for the purposes of data liberation and interoperability with old versions of QEMU.

Users requiring native encryption should use the `qcow2` format instead with `encrypt.format=luks`.

encrypt.key-secret

Provides the ID of a `secret` object that contains the encryption key (`encrypt.format=aes`).

luks

LUKS v1 encryption format, compatible with Linux `dm-crypt/cryptsetup`

Supported options:

key-secret

Provides the ID of a `secret` object that contains the passphrase.

cipher-alg

Name of the cipher algorithm and key length. Currently defaults to `aes-256`.

cipher-mode

Name of the encryption mode to use. Currently defaults to `xts`.

ivgen-alg

Name of the initialization vector generator algorithm. Currently defaults to `plain64`.

ivgen-hash-alg

Name of the hash algorithm to use with the initialization vector generator (if required). Defaults to `sha256`.

hash-alg

Name of the hash algorithm to use for PBKDF algorithm Defaults to `sha256`.

iter-time

Amount of time, in milliseconds, to use for PBKDF algorithm per key slot. Defaults to 2000.

vdi

VirtualBox 1.1 compatible image format.

Supported options:

static

If this option is set to `on`, the image is created with metadata preallocation.

vmdk

VMware 3 and 4 compatible image format.

Supported options:

backing_file

File name of a base image (see `create` subcommand).

compat6

Create a VMDK version 6 image (instead of version 4)

hwversion

Specify vmdk virtual hardware version. `Compat6` flag cannot be enabled if `hwversion` is specified.

subformat

Specifies which VMDK subformat to use. Valid options are `monolithicSparse` (default), `monolithicFlat`, `twoGbMaxExtentSparse`, `twoGbMaxExtentFlat` and `streamOptimized`.

vpc

VirtualPC compatible image format (VHD).

Supported options:

subformat

Specifies which VHD subformat to use. Valid options are `dynamic` (default) and `fixed`.

VHDX

Hyper-V compatible image format (VHDX).

Supported options:

subformat

Specifies which VHDX subformat to use. Valid options are `dynamic` (default) and `fixed`.

block_state_zero

Force use of payload blocks of type 'ZERO'. Can be set to `on` (default) or `off`. When set to `off`, new blocks will be created as `PAYLOAD_BLOCK_NOT_PRESENT`, which means parsers are free to return arbitrary data for those blocks. Do not set to `off` when using `qemu-img convert` with `subformat=dynamic`.

block_size

Block size; min 1 MB, max 256 MB. 0 means auto-calculate based on image size.

log_size

Log size; min 1 MB.

1.6.5 Read-only formats

More disk image file formats are supported in a read-only mode.

bochs

Bochs images of `growing` type.

cloop

Linux Compressed Loop image, useful only to reuse directly compressed CD-ROM images present for example in the Knoppix CD-ROMs.

dmg

Apple disk image.

parallels

Parallels disk image format.

1.6.6 Using host drives

In addition to disk image files, QEMU can directly access host devices. We describe here the usage for QEMU version `>= 0.8.3`.

Linux

On Linux, you can directly use the host device filename instead of a disk image filename provided you have enough privileges to access it. For example, use `/dev/cdrom` to access to the CDROM.

CD You can specify a CDROM device even if no CDROM is loaded. QEMU has specific code to detect CDROM insertion or removal. CDROM ejection by the guest OS is supported. Currently only data CDs are supported.

Floppy You can specify a floppy device even if no floppy is loaded. Floppy removal is currently not detected accurately (if you change floppy without doing floppy access while the floppy is not loaded, the guest OS will think that the same floppy is loaded). Use of the host's floppy device is deprecated, and support for it will be removed in a future release.

Hard disks Hard disks can be used. Normally you must specify the whole disk (`/dev/hdb` instead of `/dev/hdb1`) so that the guest OS can see it as a partitioned disk. **WARNING:** unless you know what you do, it is better to only make READ-ONLY accesses to the hard disk otherwise you may corrupt your host data (use the `-snapshot` command line option or modify the device permissions accordingly).

Windows

CD The preferred syntax is the drive letter (e.g. `d:`). The alternate syntax `\\.\\d:` is supported. `/dev/cdrom` is supported as an alias to the first CDROM drive.

Currently there is no specific code to handle removable media, so it is better to use the `change` or `eject` monitor commands to change or eject media.

Hard disks Hard disks can be used with the syntax: `\\.\\PhysicalDriveN` where *N* is the drive number (0 is the first hard disk).

WARNING: unless you know what you do, it is better to only make READ-ONLY accesses to the hard disk otherwise you may corrupt your host data (use the `-snapshot` command line so that the modifications are written in a temporary file).

Mac OS X

`/dev/cdrom` is an alias to the first CDROM.

Currently there is no specific code to handle removable media, so it is better to use the `change` or `eject` monitor commands to change or eject media.

1.6.7 Virtual FAT disk images

QEMU can automatically create a virtual FAT disk image from a directory tree. In order to use it, just type:

```
qemu-system-x86_64 linux.img -hdb fat:/my_directory
```

Then you access access to all the files in the `/my_directory` directory without having to copy them in a disk image or to export them via SAMBA or NFS. The default access is *read-only*.

Floppies can be emulated with the `:floppy:` option:

```
qemu-system-x86_64 linux.img -fda fat:floppy:/my_directory
```

A read/write support is available for testing (beta stage) with the `:rw:` option:

```
qemu-system-x86_64 linux.img -fda fat:floppy:rw:/my_directory
```

What you should *never* do:

- use non-ASCII filenames
- use “-snapshot” together with “:rw:”
- expect it to work when loadvm’ing
- write to the FAT directory on the host system while accessing it with the guest system

1.6.8 NBD access

QEMU can access directly to block device exported using the Network Block Device protocol.

```
qemu-system-x86_64 linux.img -hdb nbd://my_nbd_server.mydomain.org:1024/
```

If the NBD server is located on the same host, you can use an unix socket instead of an inet socket:

```
qemu-system-x86_64 linux.img -hdb nbd+unix:///socket=/tmp/my_socket
```

In this case, the block device must be exported using `qemu-nbd`:

```
qemu-nbd --socket=/tmp/my_socket my_disk.qcow2
```

The use of `qemu-nbd` allows sharing of a disk between several guests:

```
qemu-nbd --socket=/tmp/my_socket --share=2 my_disk.qcow2
```

and then you can use it with two guests:

```
qemu-system-x86_64 linux1.img -hdb nbd+unix:///socket=/tmp/my_socket
qemu-system-x86_64 linux2.img -hdb nbd+unix:///socket=/tmp/my_socket
```

If the `nbd-server` uses named exports (supported since NBD 2.9.18, or with QEMU's own embedded NBD server), you must specify an export name in the URI:

```
qemu-system-x86_64 -cdrom nbd://localhost/debian-500-ppc-netinst
qemu-system-x86_64 -cdrom nbd://localhost/openSUSE-11.1-ppc-netinst
```

The URI syntax for NBD is supported since QEMU 1.3. An alternative syntax is also available. Here are some example of the older syntax:

```
qemu-system-x86_64 linux.img -hdb nbd:my_nbd_server.mydomain.org:1024
qemu-system-x86_64 linux2.img -hdb nbd:unix:/tmp/my_socket
qemu-system-x86_64 -cdrom nbd:localhost:10809:exportname=debian-500-ppc-
↪netinst
```

1.6.9 Sheepdog disk images

Sheepdog is a distributed storage system for QEMU. It provides highly available block level storage volumes that can be attached to QEMU-based virtual machines.

You can create a Sheepdog disk image with the command:

```
qemu-img create sheepdog:///IMAGE SIZE
```

where *IMAGE* is the Sheepdog image name and *SIZE* is its size.

To import the existing *FILENAME* to Sheepdog, you can use a convert command.

```
qemu-img convert FILENAME sheepdog:///IMAGE
```

You can boot from the Sheepdog disk image with the command:

```
qemu-system-x86_64 sheepdog:///IMAGE
```

You can also create a snapshot of the Sheepdog image like `qcow2`.

```
qemu-img snapshot -c TAG sheepdog:///IMAGE
```

where *TAG* is a tag name of the newly created snapshot.

To boot from the Sheepdog snapshot, specify the tag name of the snapshot.

```
qemu-system-x86_64 sheepdog:///IMAGE#TAG
```

You can create a cloned image from the existing snapshot.

```
qemu-img create -b sheepdog:///BASE#TAG sheepdog:///IMAGE
```

where *BASE* is an image name of the source snapshot and *TAG* is its tag name.

You can use an unix socket instead of an inet socket:

```
qemu-system-x86_64 sheepdog+unix:///IMAGE?socket=PATH
```

If the Sheepdog daemon doesn't run on the local host, you need to specify one of the Sheepdog servers to connect to.

```
qemu-img create sheepdog://HOSTNAME:PORT/IMAGE SIZE
qemu-system-x86_64 sheepdog://HOSTNAME:PORT/IMAGE
```

1.6.10 iSCSI LUNs

iSCSI is a popular protocol used to access SCSI devices across a computer network.

There are two different ways iSCSI devices can be used by QEMU.

The first method is to mount the iSCSI LUN on the host, and make it appear as any other ordinary SCSI device on the host and then to access this device as a /dev/sd device from QEMU. How to do this differs between host OSes.

The second method involves using the iSCSI initiator that is built into QEMU. This provides a mechanism that works the same way regardless of which host OS you are running QEMU on. This section will describe this second method of using iSCSI together with QEMU.

In QEMU, iSCSI devices are described using special iSCSI URLs. URL syntax:

```
iscsi://[<username>[%<password>]@]<host>[:<port>]/<target-iqn-name>/<lun>
```

Username and password are optional and only used if your target is set up using CHAP authentication for access control. Alternatively the username and password can also be set via environment variables to have these not show up in the process list:

```
export LIBISCSI_CHAP_USERNAME=<username>
export LIBISCSI_CHAP_PASSWORD=<password>
iscsi://<host>/<target-iqn-name>/<lun>
```

Various session related parameters can be set via special options, either in a configuration file provided via '-readconfig' or directly on the command line.

If the initiator-name is not specified qemu will use a default name of 'iqn.2008-11.org.linux-kvm[:<uuid>]' where <uuid> is the UUID of the virtual machine. If the UUID is not specified qemu will use 'iqn.2008-11.org.linux-kvm[:<name>]' where <name> is the name of the virtual machine.

Setting a specific initiator name to use when logging in to the target:

```
-iscsi initiator-name=iqn.qemu.test:my-initiator
```

Controlling which type of header digest to negotiate with the target:

```
-iscsi header-digest=CRC32C|CRC32C-NONE|NONE-CRC32C|NONE
```

These can also be set via a configuration file:

```
[iscsi]
user = "CHAP username"
password = "CHAP password"
initiator-name = "iqn.qemu.test:my-initiator"
# header digest is one of CRC32C/CRC32C-NONE/NONE-CRC32C/NONE
header-digest = "CRC32C"
```

Setting the target name allows different options for different targets:

```
[iscsi "iqn.target.name"]
user = "CHAP username"
password = "CHAP password"
initiator-name = "iqn.qemu.test:my-initiator"
# header digest is one of CRC32C/CRC32C-NONE/NONE-CRC32C/NONE
header-digest = "CRC32C"
```

How to use a configuration file to set iSCSI configuration options:

```
cat >iscsi.conf <<EOF
[iscsi]
user = "me"
password = "my password"
initiator-name = "iqn.qemu.test:my-initiator"
header-digest = "CRC32C"
EOF
```

```
qemu-system-x86_64 -drive file=iscsi://127.0.0.1/iqn.qemu.test/1 \
  -readconfig iscsi.conf
```

How to set up a simple iSCSI target on loopback and access it via QEMU: this example shows how to set up an iSCSI target with one CDROM and one DISK using the Linux STGT software target. This target is available on Red Hat based systems as the package ‘scsi-target-utils’.

```
tgtadm --iscsi portal=127.0.0.1:3260
tgtadm --lld iscsi --op new --mode target --tid 1 -T iqn.qemu.test
tgtadm --lld iscsi --mode logicalunit --op new --tid 1 --lun 1 \
  -b /IMAGES/disk.img --device-type=disk
tgtadm --lld iscsi --mode logicalunit --op new --tid 1 --lun 2 \
  -b /IMAGES/cd.iso --device-type=cd
tgtadm --lld iscsi --op bind --mode target --tid 1 -I ALL

qemu-system-x86_64 -iscsi initiator-name=iqn.qemu.test:my-initiator \
  -boot d -drive file=iscsi://127.0.0.1/iqn.qemu.test/1 \
  -cdrom iscsi://127.0.0.1/iqn.qemu.test/2
```

1.6.11 GlusterFS disk images

GlusterFS is a user space distributed file system.

You can boot from the GlusterFS disk image with the command:

URI:

```
qemu-system-x86_64 -drive file=gluster[+TYPE]://[HOST][:PORT]/VOLUME/PATH
                    [?socket=...] [,file.debug=9] [,file.logfile=...]
```

JSON:

```
qemu-system-x86_64 'json:{"driver":"qcow2",
                        "file":{"driver":"gluster",
                                "volume":"testvol","path":"a.img","debug":9,
                                ↪"logfile":"...",
                                "server":[{"type":"tcp","host":"...",
                                ↪"port":"..."}],
                                {"type":"unix","socket":"..."}]}}'
```

gluster is the protocol.

TYPE specifies the transport type used to connect to gluster management daemon (glusterd). Valid transport types are tcp and unix. In the URI form, if a transport type isn't specified, then tcp type is assumed.

HOST specifies the server where the volume file specification for the given volume resides. This can be either a hostname or an ipv4 address. If transport type is unix, then *HOST* field should not be specified. Instead *socket* field needs to be populated with the path to unix domain socket.

PORT is the port number on which glusterd is listening. This is optional and if not specified, it defaults to port 24007. If the transport type is unix, then *PORT* should not be specified.

VOLUME is the name of the gluster volume which contains the disk image.

PATH is the path to the actual disk image that resides on gluster volume.

debug is the logging level of the gluster protocol driver. Debug levels are 0-9, with 9 being the most verbose, and 0 representing no debugging output. The default level is 4. The current logging levels defined in the gluster source are 0 - None, 1 - Emergency, 2 - Alert, 3 - Critical, 4 - Error, 5 - Warning, 6 - Notice, 7 - Info, 8 - Debug, 9 - Trace

logfile is a commandline option to mention log file path which helps in logging to the specified file and also help in persisting the gfapi logs. The default is stderr.

You can create a GlusterFS disk image with the command:

```
qemu-img create gluster://HOST/VOLUME/PATH SIZE
```

Examples

```
qemu-system-x86_64 -drive file=gluster://1.2.3.4/testvol/a.img
qemu-system-x86_64 -drive file=gluster+tcp://1.2.3.4/testvol/a.img
qemu-system-x86_64 -drive file=gluster+tcp://1.2.3.4:24007/testvol/dir/a.img
qemu-system-x86_64 -drive file=gluster+tcp://[1:2:3:4:5:6:7:8]/testvol/dir/a.
↪img
qemu-system-x86_64 -drive file=gluster+tcp://[1:2:3:4:5:6:7:8]:24007/testvol/
↪dir/a.img
qemu-system-x86_64 -drive file=gluster+tcp://server.domain.com:24007/testvol/
↪dir/a.img
qemu-system-x86_64 -drive file=gluster+unix:///testvol/dir/a.img?socket=/tmp/
↪glusterd.socket
qemu-system-x86_64 -drive file=gluster+rdma://1.2.3.4:24007/testvol/a.img
qemu-system-x86_64 -drive file=gluster://1.2.3.4/testvol/a.img,file.debug=9,
↪file.logfile=/var/log/qemu-gluster.log
qemu-system-x86_64 'json:{"driver":"qcow2",
                        "file":{"driver":"gluster",
                                "volume":"testvol","path":"a.img",
                                "debug":9,"logfile":"/var/log/qemu-gluster.
↪log",
                                "server":[{"type":"tcp","host":"1.2.3.4",
↪"port":24007},
```

```
                                {"type": "unix", "socket": "/var/run/
↪glusterd.socket"}]]}]'
qemu-system-x86_64 -drive driver=qcow2,file.driver=gluster,file.
↪volume=testvol,file.path=/path/a.img,
                                file.debug=9,file.logfile=/var/log/qemu-
↪gluster.log,
                                file.server.0.type=tcp,file.server.0.
↪host=1.2.3.4,file.server.0.port=24007,
                                file.server.1.type=unix,file.server.1.
↪socket=/var/run/glusterd.socket
```

1.6.12 Secure Shell (ssh) disk images

You can access disk images located on a remote ssh server by using the ssh protocol:

```
qemu-system-x86_64 -drive file=ssh://[USER@]SERVER[:PORT]/PATH[?host_key_
↪check=HOST_KEY_CHECK]
```

Alternative syntax using properties:

```
qemu-system-x86_64 -drive file.driver=ssh[,file.user=USER],file.host=SERVER[,
↪file.port=PORT],file.path=PATH[,file.host_key_check=HOST_KEY_CHECK]
```

ssh is the protocol.

USER is the remote user. If not specified, then the local username is tried.

SERVER specifies the remote ssh server. Any ssh server can be used, but it must implement the sftp-server protocol. Most Unix/Linux systems should work without requiring any extra configuration.

PORT is the port number on which sshd is listening. By default the standard ssh port (22) is used.

PATH is the path to the disk image.

The optional *HOST_KEY_CHECK* parameter controls how the remote host's key is checked. The default is *yes* which means to use the local *.ssh/known_hosts* file. Setting this to *no* turns off known-hosts checking. Or you can check that the host key matches a specific fingerprint: *host_key_check=md5:78:45:8e:14:57:4f:d5:45:83:0a:0e:f3:49:82:c9:c8* (sha1: can also be used as a prefix, but note that OpenSSH tools only use MD5 to print fingerprints).

Currently authentication must be done using ssh-agent. Other authentication methods may be supported in future.

Note: Many ssh servers do not support an *fsync*-style operation. The ssh driver cannot guarantee that disk flush requests are obeyed, and this causes a risk of disk corruption if the remote server or network goes down during writes. The driver will print a warning when *fsync* is not supported:

```
warning: ssh server ssh.example.com:22 does not support fsync
```

With sufficiently new versions of libssh and OpenSSH, *fsync* is supported.

1.6.13 NVMe disk images

NVM Express (NVMe) storage controllers can be accessed directly by a userspace driver in QEMU. This bypasses the host kernel file system and block layers while retaining QEMU block layer functionalities, such as block jobs, I/O throttling, image formats, etc. Disk I/O performance is typically higher than with `-drive file=/dev/sda` using either thread pool or linux-aio.

The controller will be exclusively used by the QEMU process once started. To be able to share storage between multiple VMs and other applications on the host, please use the file based protocols.

Before starting QEMU, bind the host NVMe controller to the host vfio-pci driver. For example:

```
# modprobe vfio-pci
# lspci -n -s 0000:06:0d.0
06:0d.0 0401: 1102:0002 (rev 08)
# echo 0000:06:0d.0 > /sys/bus/pci/devices/0000:06:0d.0/driver/unbind
# echo 1102 0002 > /sys/bus/pci/drivers/vfio-pci/new_id

# qemu-system-x86_64 -drive file=nvme://HOST:BUS:SLOT.FUNC/NAMESPACE
```

Alternative syntax using properties:

```
qemu-system-x86_64 -drive file.driver=nvme,file.device=HOST:BUS:SLOT.FUNC,
↳file.namespace=NAMESPACE
```

HOST:BUS:SLOT.FUNC is the NVMe controller's PCI device address on the host.

NAMESPACE is the NVMe namespace number, starting from 1.

1.6.14 Disk image file locking

By default, QEMU tries to protect image files from unexpected concurrent access, as long as it's supported by the block protocol driver and host operating system. If multiple QEMU processes (including QEMU emulators and utilities) try to open the same image with conflicting accessing modes, all but the first one will get an error.

This feature is currently supported by the file protocol on Linux with the Open File Descriptor (OFD) locking API, and can be configured to fall back to POSIX locking if the POSIX host doesn't support Linux OFD locking.

To explicitly enable image locking, specify "locking=on" in the file protocol driver options. If OFD locking is not possible, a warning will be printed and the POSIX locking API will be used. In this case there is a risk that the lock will get silently lost when doing hot plugging and block jobs, due to the shortcomings of the POSIX locking API.

QEMU transparently handles lock handover during shared storage migration. For shared virtual disk images between multiple VMs, the "share-rw" device option should be used.

By default, the guest has exclusive write access to its disk image. If the guest can safely share the disk image with other writers the `-device ..., share-rw=on` parameter can be used. This is only safe if the guest is running software, such as a cluster file system, that coordinates disk accesses to avoid corruption.

Note that `share-rw=on` only declares the guest's ability to share the disk. Some QEMU features, such as image file formats, require exclusive write access to the disk image and this is unaffected by the `share-rw=on` option.

Alternatively, locking can be fully disabled by "locking=off" block device option. In the command line, the option is usually in the form of "file.locking=off" as the protocol driver is normally placed as a "file" child under a format driver. For example:

```
-blockdev driver=qcow2,file.filename=/path/to/image,file.locking=off,file.driver=file
```

To check if image locking is active, check the output of the "lslocks" command on host and see if there are locks held by the QEMU process on the image file. More than one byte could be locked by the QEMU instance, each byte of which reflects a particular permission that is acquired or protected by the running block driver.

1.6.15 Filter drivers

QEMU supports several filter drivers, which don't store any data, but perform some additional tasks, hooking io requests.

preallocate

The preallocate filter driver is intended to be inserted between format and protocol nodes and preallocates some additional space (expanding the protocol file) when writing past the file's end. This can be useful for file-systems with slow allocation.

Supported options:

prealloc-align

On preallocation, align the file length to this value (in bytes), default 1M.

prealloc-size

How much to preallocate (in bytes), default 128M.

1.7 Network emulation

QEMU can simulate several network cards (e.g. PCI or ISA cards on the PC target) and can connect them to a network backend on the host or an emulated hub. The various host network backends can either be used to connect the NIC of the guest to a real network (e.g. by using a TAP devices or the non-privileged user mode network stack), or to other guest instances running in another QEMU process (e.g. by using the socket host network backend).

1.7.1 Using TAP network interfaces

This is the standard way to connect QEMU to a real network. QEMU adds a virtual network device on your host (called `tapN`), and you can then configure it as if it was a real ethernet card.

Linux host

As an example, you can download the `linux-test-xxx.tar.gz` archive and copy the script `qemu-ifup` in `/etc` and configure properly `sudo` so that the command `ifconfig` contained in `qemu-ifup` can be executed as root. You must verify that your host kernel supports the TAP network interfaces: the device `/dev/net/tun` must be present.

See *Invocation* to have examples of command lines using the TAP network interfaces.

Windows host

There is a virtual ethernet driver for Windows 2000/XP systems, called TAP-Win32. But it is not included in standard QEMU for Windows, so you will need to get it separately. It is part of OpenVPN package, so download OpenVPN from : <https://openvpn.net/>.

1.7.2 Using the user mode network stack

By using the option `-net user` (default configuration if no `-net` option is specified), QEMU uses a completely user mode network stack (you don't need root privilege to use the virtual network). The virtual network configuration is the following:

```
guest (10.0.2.15) <-----> Firewall/DHCP server <-----> Internet
                        |
                        | (10.0.2.2)
                        |
                        ----> DNS server (10.0.2.3)
                        |
                        ----> SMB server (10.0.2.4)
```


The QEMU VM behaves as if it was behind a firewall which blocks all incoming connections. You can use a DHCP client to automatically configure the network in the QEMU VM. The DHCP server assign addresses to the hosts starting from 10.0.2.15.

In order to check that the user mode network is working, you can ping the address 10.0.2.2 and verify that you got an address in the range 10.0.2.x from the QEMU virtual DHCP server.

Note that ICMP traffic in general does not work with user mode networking. `ping`, aka. ICMP echo, to the local router (10.0.2.2) shall work, however. If you're using QEMU on Linux ≥ 3.0 , it can use unprivileged ICMP ping sockets to allow `ping` to the Internet. The host admin has to set the `ping_group_range` in order to grant access to those sockets. To allow ping for GID 100 (usually users group):

```
echo 100 100 > /proc/sys/net/ipv4/ping_group_range
```

When using the built-in TFTP server, the router is also the TFTP server.

When using the `'-netdev user,hostfwd=...'` option, TCP or UDP connections can be redirected from the host to the guest. It allows for example to redirect X11, telnet or SSH connections.

1.7.3 Hubs

QEMU can simulate several hubs. A hub can be thought of as a virtual connection between several network devices. These devices can be for example QEMU virtual ethernet cards or virtual Host ethernet devices (TAP devices). You can connect guest NICs or host network backends to such a hub using the `-netdev hubport` or `-nic hubport` options. The legacy `-net` option also connects the given device to the emulated hub with ID 0 (i.e. the default hub) unless you specify a netdev with `-net nic,netdev=xxx` here.

1.7.4 Connecting emulated networks between QEMU instances

Using the `-netdev socket` (or `-nic socket` or `-net socket`) option, it is possible to create emulated networks that span several QEMU instances. See the description of the `-netdev socket` option in *Invocation* to have a basic example.

1.8 QEMU virtio-net standby (net_failover)

This document explains the setup and usage of virtio-net standby feature which is used to create a `net_failover` pair of devices.

The general idea is that we have a pair of devices, a (vfi-)pci and a virtio-net device. Before migration the vfi device is unplugged and data flows through the virtio-net device, on the target side another vfi-pci device is plugged in to take over the data-path. In the guest the `net_failover` kernel module will pair net devices with the same MAC address.

The two devices are called primary and standby device. The fast hardware based networking device is called the primary device and the virtio-net device is the standby device.

1.8.1 Restrictions

Currently only PCIe devices are allowed as primary devices, this restriction can be lifted in the future with enhanced QEMU support. Also, only networking devices are allowed as primary device. The user needs to ensure that primary and standby devices are not plugged into the same PCIe slot.

1.8.2 Usecase

Virtio-net standby allows easy migration while using a passed-through fast networking device by falling back to a virtio-net device for the duration of the migration. It is like a simple version of a bond, the difference is that it requires no configuration in the guest. When a guest is live-migrated to another host QEMU will unplug the primary device via the PCIe based hotplug handler and traffic will go through the virtio-net device. On the target system the primary device will be automatically plugged back and the `net_failover` module registers it again as the primary device.

1.8.3 Usage

The primary device can be hotplugged or be part of the startup configuration

-device virtio-net-pci,netdev=hostnet1,id=net1,mac=52:54:00:6f:55:cc, bus=root2,failover=on

With the parameter `failover=on` the `VIRTIO_NET_F_STANDBY` feature will be enabled.

-device vfio-pci,host=5e:00:02,id=hostdev0,bus=root1,failover_pair_id=net1

`failover_pair_id` references the id of the virtio-net standby device. This is only for pairing the devices within QEMU. The guest kernel module `net_failover` will match devices with identical MAC addresses.

1.8.4 Hotplug

Both primary and standby device can be hotplugged via the QEMU monitor. Note that if the virtio-net device is plugged first a warning will be issued that it couldn't find the primary device.

1.8.5 Migration

A new migration state `wait-unplug` was added for this feature. If failover primary devices are present in the configuration, migration will go into this state. It will wait until the device unplug is completed in the guest and then move into active state. On the target system the primary devices will be automatically hotplugged when the feature bit was negotiated for the virtio-net standby device.

1.9 USB emulation

QEMU can emulate a PCI UHCI, OHCI, EHCI or XHCI USB controller. You can plug virtual USB devices or real host USB devices (only works with certain host operating systems). QEMU will automatically create and connect virtual USB hubs as necessary to connect multiple USB devices.

1.9.1 Connecting USB devices

USB devices can be connected with the `-device usb-...` command line option or the `device_add` monitor command. Available devices are:

usb-mouse Virtual Mouse. This will override the PS/2 mouse emulation when activated.

usb-tablet Pointer device that uses absolute coordinates (like a touchscreen). This means QEMU is able to report the mouse position without having to grab the mouse. Also overrides the PS/2 mouse emulation when activated.

usb-storage,drive=drive_id Mass storage device backed by `drive_id` (see the *Disk Images* chapter in the System Emulation Users Guide)

usb-uas USB attached SCSI device, see [usb-storage.txt](#) for details

usb-bot Bulk-only transport storage device, see [usb-storage.txt](#) for details here, too

usb-mtp, rootdir=dir Media transfer protocol device, using dir as root of the file tree that is presented to the guest.

usb-host, hostbus=bus, hostaddr=addr Pass through the host device identified by bus and addr

usb-host, vendorid=vendor, productid=product Pass through the host device identified by vendor and product ID

usb-wacom-tablet Virtual Wacom PenPartner tablet. This device is similar to the `tablet` above but it can be used with the tslib library because in addition to touch coordinates it reports touch pressure.

usb-kbd Standard USB keyboard. Will override the PS/2 keyboard (if present).

usb-serial, chardev=id Serial converter. This emulates an FTDI FT232BM chip connected to host character device id.

usb-braille, chardev=id Braille device. This will use BrlAPI to display the braille output on a real or fake device referenced by id.

usb-net [, netdev=id] Network adapter that supports CDC ethernet and RNDIS protocols. id specifies a netdev defined with `-netdev ..., id=id`. For instance, user-mode networking can be used with

```
qemu-system-x86_64 [...] -netdev user,id=net0 -device usb-net,netdev=net0
```

usb-ccid Smartcard reader device

usb-audio USB audio device

u2f-{emulated, passthru} Universal Second Factor device

1.9.2 Using host USB devices on a Linux host

WARNING: this is an experimental feature. QEMU will slow down when using it. USB devices requiring real time streaming (i.e. USB Video Cameras) are not supported yet.

1. If you use an early Linux 2.4 kernel, verify that no Linux driver is actually using the USB device. A simple way to do that is simply to disable the corresponding kernel module by renaming it from `mydriver.o` to `mydriver.o.disabled`.
2. Verify that `/proc/bus/usb` is working (most Linux distributions should enable it by default). You should see something like that:

```
ls /proc/bus/usb
001  devices  drivers
```

3. Since only root can access to the USB devices directly, you can either launch QEMU as root or change the permissions of the USB devices you want to use. For testing, the following suffices:

```
chown -R myuid /proc/bus/usb
```

4. Launch QEMU and do in the monitor:

```
info usbhost
Device 1.2, speed 480 Mb/s
Class 00: USB device 1234:5678, USB DISK
```

You should see the list of the devices you can use (Never try to use hubs, it won't work).

5. Add the device in QEMU by using:

```
device_add usb-host,vendorid=0x1234,productid=0x5678
```

Normally the guest OS should report that a new USB device is plugged. You can use the option `-device usb-host, ...` to do the same.

6. Now you can try to use the host USB device in QEMU.

When relaunching QEMU, you may have to unplug and plug again the USB device to make it work again (this is a bug).

1.10 Inter-VM Shared Memory device

On Linux hosts, a shared memory device is available. The basic syntax is:

```
qemu_system-x86_64 -device ivshmem-plain,memdev=hostmem
```

where `hostmem` names a host memory backend. For a POSIX shared memory backend, use something like

```
-object memory-backend-file,size=1M,share,mem-path=/dev/shm/ivshmem,id=hostmem
```

If desired, interrupts can be sent between guest VMs accessing the same shared memory region. Interrupt support requires using a shared memory server and using a chardev socket to connect to it. The code for the shared memory server is [qemu.git/contrib/ivshmem-server](https://github.com/qemu/qemu/blob/master/contrib/ivshmem-server). An example syntax when using the shared memory server is:

```
# First start the ivshmem server once and for all
ivshmem-server -p pidfile -S path -m shm-name -l shm-size -n vectors

# Then start your qemu instances with matching arguments
qemu_system-x86_64 -device ivshmem-doorbell,vectors=vectors,chardev=id
                  -chardev socket,path=path,id=id
```

When using the server, the guest will be assigned a VM ID (≥ 0) that allows guests using the same server to communicate via interrupts. Guests can read their VM ID from a device register (see `ivshmem-spec.txt`).

1.10.1 Migration with ivshmem

With device property `master=on`, the guest will copy the shared memory on migration to the destination host. With `master=off`, the guest will not be able to migrate with the device attached. In the latter case, the device should be detached and then reattached after migration using the PCI hotplug support.

At most one of the devices sharing the same memory can be master. The master must complete migration before you plug back the other devices.

1.10.2 ivshmem and hugepages

Instead of specifying the `<shm size>` using POSIX `shm`, you may specify a memory backend that has hugepage support:

```
qemu_system-x86_64 -object memory-backend-file,size=1G,mem-path=/dev/
    ↪hugepages/my-shmem-file,share,id=mb1
                  -device ivshmem-plain,memdev=mb1
```

`ivshmem-server` also supports hugepages mount points with the `-m` memory path argument.

1.11 Direct Linux Boot

This section explains how to launch a Linux kernel inside QEMU without having to make a full bootable image. It is very useful for fast Linux kernel testing.

The syntax is:

```
qemu-system-x86_64 -kernel bzImage -hda rootdisk.img -append "root=/dev/hda"
```

Use `-kernel` to provide the Linux kernel image and `-append` to give the kernel command line arguments. The `-initrd` option can be used to provide an INITRD image.

If you do not need graphical output, you can disable it and redirect the virtual serial port and the QEMU monitor to the console with the `-nographic` option. The typical command line is:

```
qemu-system-x86_64 -kernel bzImage -hda rootdisk.img -append "root=/dev/hda console=ttyS0" -nographic
```

Use Ctrl-a c to switch between the serial console and the monitor (see *Keys in the graphical frontends*).

1.12 VNC security

The VNC server capability provides access to the graphical console of the guest VM across the network. This has a number of security considerations depending on the deployment scenarios.

1.12.1 Without passwords

The simplest VNC server setup does not include any form of authentication. For this setup it is recommended to restrict it to listen on a UNIX domain socket only. For example

```
qemu-system-x86_64 [...OPTIONS...] -vnc unix:/home/joebloggs/.qemu-myvm-vnc
```

This ensures that only users on local box with read/write access to that path can access the VNC server. To securely access the VNC server from a remote machine, a combination of netcat+ssh can be used to provide a secure tunnel.

1.12.2 With passwords

The VNC protocol has limited support for password based authentication. Since the protocol limits passwords to 8 characters it should not be considered to provide high security. The password can be fairly easily brute-forced by a client making repeat connections. For this reason, a VNC server using password authentication should be restricted to only listen on the loopback interface or UNIX domain sockets. Password authentication is not supported when operating in FIPS 140-2 compliance mode as it requires the use of the DES cipher. Password authentication is requested with the `password` option, and then once QEMU is running the password is set with the monitor. Until the monitor is used to set the password all clients will be rejected.

```
qemu-system-x86_64 [...OPTIONS...] -vnc :1,password -monitor stdio
(qemu) change vnc password
Password: ****
(qemu)
```

1.12.3 With x509 certificates

The QEMU VNC server also implements the VeNCrypt extension allowing use of TLS for encryption of the session, and x509 certificates for authentication. The use of x509 certificates is strongly recommended, because TLS on its own is susceptible to man-in-the-middle attacks. Basic x509 certificate support provides a secure session, but no authentication. This allows any client to connect, and provides an encrypted session.

```
qemu-system-x86_64 [...OPTIONS...] -object tls-creds-x509,id=tls0,dir=/etc/  
→pki/qemu,endpoint=server,verify-peer=off -vnc :1,tls-creds=tls0 -monitor_  
→stdio
```

In the above example `/etc/pki/qemu` should contain at least three files, `ca-cert.pem`, `server-cert.pem` and `server-key.pem`. Unprivileged users will want to use a private directory, for example `$HOME/.pki/qemu`. NB the `server-key.pem` file should be protected with file mode 0600 to only be readable by the user owning it.

1.12.4 With x509 certificates and client verification

Certificates can also provide a means to authenticate the client connecting. The server will request that the client provide a certificate, which it will then validate against the CA certificate. This is a good choice if deploying in an environment with a private internal certificate authority. It uses the same syntax as previously, but with `verify-peer` set to `on` instead.

```
qemu-system-x86_64 [...OPTIONS...] -object tls-creds-x509,id=tls0,dir=/etc/  
→pki/qemu,endpoint=server,verify-peer=on -vnc :1,tls-creds=tls0 -monitor_  
→stdio
```

1.12.5 With x509 certificates, client verification and passwords

Finally, the previous method can be combined with VNC password authentication to provide two layers of authentication for clients.

```
qemu-system-x86_64 [...OPTIONS...] -object tls-creds-x509,id=tls0,dir=/etc/  
→pki/qemu,endpoint=server,verify-peer=on -vnc :1,tls-creds=tls0,password -  
→monitor stdio  
(qemu) change vnc password  
Password: ****  
(qemu)
```

1.12.6 With SASL authentication

The SASL authentication method is a VNC extension, that provides an easily extendable, pluggable authentication method. This allows for integration with a wide range of authentication mechanisms, such as PAM, GSSAPI/Kerberos, LDAP, SQL databases, one-time keys and more. The strength of the authentication depends on the exact mechanism configured. If the chosen mechanism also provides a SSF layer, then it will encrypt the datastream as well.

Refer to the later docs on how to choose the exact SASL mechanism used for authentication, but assuming use of one supporting SSF, then QEMU can be launched with:

```
qemu-system-x86_64 [...OPTIONS...] -vnc :1,sasl -monitor stdio
```

1.12.7 With x509 certificates and SASL authentication

If the desired SASL authentication mechanism does not supported SSF layers, then it is strongly advised to run it in combination with TLS and x509 certificates. This provides securely encrypted data stream, avoiding risk of

compromising of the security credentials. This can be enabled, by combining the ‘sasl’ option with the aforementioned TLS + x509 options:

```
qemu-system-x86_64 [...OPTIONS...] -object tls-creds-x509,id=tls0,dir=/
↳etc/pki/qemu,endpoint=server,verify-peer=on -vnc :1,tls-creds=tls0,sasl -
↳monitor stdio
```

1.12.8 Configuring SASL mechanisms

The following documentation assumes use of the Cyrus SASL implementation on a Linux host, but the principles should apply to any other SASL implementation or host. When SASL is enabled, the mechanism configuration will be loaded from system default SASL service config `/etc/sasl2/qemu.conf`. If running QEMU as an unprivileged user, an environment variable `SASL_CONF_PATH` can be used to make it search alternate locations for the service config file.

If the TLS option is enabled for VNC, then it will provide session encryption, otherwise the SASL mechanism will have to provide encryption. In the latter case the list of possible plugins that can be used is drastically reduced. In fact only the GSSAPI SASL mechanism provides an acceptable level of security by modern standards. Previous versions of QEMU referred to the DIGEST-MD5 mechanism, however, it has multiple serious flaws described in detail in RFC 6331 and thus should never be used any more. The SCRAM-SHA-1 mechanism provides a simple username/password auth facility similar to DIGEST-MD5, but does not support session encryption, so can only be used in combination with TLS.

When not using TLS the recommended configuration is

```
mech_list: gssapi
keytab: /etc/qemu/krb5.tab
```

This says to use the ‘GSSAPI’ mechanism with the Kerberos v5 protocol, with the server principal stored in `/etc/qemu/krb5.tab`. For this to work the administrator of your KDC must generate a Kerberos principal for the server, with a name of ‘`qemu/somehost.example.com@EXAMPLE.COM`’ replacing ‘`somehost.example.com`’ with the fully qualified host name of the machine running QEMU, and ‘`EXAMPLE.COM`’ with the Kerberos Realm.

When using TLS, if username+password authentication is desired, then a reasonable configuration is

```
mech_list: scram-sha-1
sasldb_path: /etc/qemu/passwd.db
```

The `saslpasswd2` program can be used to populate the `passwd.db` file with accounts.

Other SASL configurations will be left as an exercise for the reader. Note that all mechanisms, except GSSAPI, should be combined with use of TLS to ensure a secure data channel.

1.13 TLS setup for network services

Almost all network services in QEMU have the ability to use TLS for session data encryption, along with x509 certificates for simple client authentication. What follows is a description of how to generate certificates suitable for usage with QEMU, and applies to the VNC server, character devices with the TCP backend, NBD server and client, and migration server and client.

At a high level, QEMU requires certificates and private keys to be provided in PEM format. Aside from the core fields, the certificates should include various extension data sets, including v3 basic constraints data, key purpose, key usage and subject alt name.

The GnuTLS package includes a command called `certtool` which can be used to easily generate certificates and keys in the required format with expected data present. Alternatively a certificate management service may be used.

At a minimum it is necessary to setup a certificate authority, and issue certificates to each server. If using x509 certificates for authentication, then each client will also need to be issued a certificate.

Assuming that the QEMU network services will only ever be exposed to clients on a private intranet, there is no need to use a commercial certificate authority to create certificates. A self-signed CA is sufficient, and in fact likely to be more secure since it removes the ability of malicious 3rd parties to trick the CA into mis-issuing certs for impersonating your services. The only likely exception where a commercial CA might be desirable is if enabling the VNC websockets server and exposing it directly to remote browser clients. In such a case it might be useful to use a commercial CA to avoid needing to install custom CA certs in the web browsers.

The recommendation is for the server to keep its certificates in either `/etc/pki/qemu` or for unprivileged users in `$HOME/.pki/qemu`.

1.13.1 Setup the Certificate Authority

This step only needs to be performed once per organization / organizational unit. First the CA needs a private key. This key must be kept VERY secret and secure. If this key is compromised the entire trust chain of the certificates issued with it is lost.

```
# certtool --generate-privkey > ca-key.pem
```

To generate a self-signed certificate requires one core piece of information, the name of the organization. A template file `ca.info` should be populated with the desired data to avoid having to deal with interactive prompts from `certtool`:

```
# cat > ca.info <<EOF
cn = Name of your organization
ca
cert_signing_key
EOF
# certtool --generate-self-signed \
           --load-privkey ca-key.pem \
           --template ca.info \
           --outfile ca-cert.pem
```

The `ca` keyword in the template sets the v3 basic constraints extension to indicate this certificate is for a CA, while `cert_signing_key` sets the key usage extension to indicate this will be used for signing other keys. The generated `ca-cert.pem` file should be copied to all servers and clients wishing to utilize TLS support in the VNC server. The `ca-key.pem` must not be disclosed/copied anywhere except the host responsible for issuing certificates.

1.13.2 Issuing server certificates

Each server (or host) needs to be issued with a key and certificate. When connecting the certificate is sent to the client which validates it against the CA certificate. The core pieces of information for a server certificate are the hostnames and/or IP addresses that will be used by clients when connecting. The hostname / IP address that the client specifies when connecting will be validated against the hostname(s) and IP address(es) recorded in the server certificate, and if no match is found the client will close the connection.

Thus it is recommended that the server certificate include both the fully qualified and unqualified hostnames. If the server will have permanently assigned IP address(es), and clients are likely to use them when connecting, they may also be included in the certificate. Both IPv4 and IPv6 addresses are supported. Historically certificates only included 1 hostname in the CN field, however, usage of this field for validation is now deprecated. Instead modern TLS clients will validate against the Subject Alt Name extension data, which allows for multiple entries. In the future usage of the CN field may be discontinued entirely, so providing SAN extension data is strongly recommended.

On the host holding the CA, create template files containing the information for each server, and use it to issue server certificates.


```
# cat > server-hostNNN.info <<EOF
organization = Name of your organization
cn = hostNNN.foo.example.com
dns_name = hostNNN
dns_name = hostNNN.foo.example.com
ip_address = 10.0.1.87
ip_address = 192.8.0.92
ip_address = 2620:0:cafe::87
ip_address = 2001:24::92
tls_www_server
encryption_key
signing_key
EOF
# certtool --generate-privkey > server-hostNNN-key.pem
# certtool --generate-certificate \
    --load-ca-certificate ca-cert.pem \
    --load-ca-privkey ca-key.pem \
    --load-privkey server-hostNNN-key.pem \
    --template server-hostNNN.info \
    --outfile server-hostNNN-cert.pem
```

The `dns_name` and `ip_address` fields in the template are setting the subject alt name extension data. The `tls_www_server` keyword is the key purpose extension to indicate this certificate is intended for usage in a web server. Although QEMU network services are not in fact HTTP servers (except for VNC websockets), setting this key purpose is still recommended. The `encryption_key` and `signing_key` keyword is the key usage extension to indicate this certificate is intended for usage in the data session.

The `server-hostNNN-key.pem` and `server-hostNNN-cert.pem` files should now be securely copied to the server for which they were generated, and renamed to `server-key.pem` and `server-cert.pem` when added to the `/etc/pki/qemu` directory on the target host. The `server-key.pem` file is security sensitive and should be kept protected with file mode 0600 to prevent disclosure.

1.13.3 Issuing client certificates

The QEMU x509 TLS credential setup defaults to enabling client verification using certificates, providing a simple authentication mechanism. If this default is used, each client also needs to be issued a certificate. The client certificate contains enough metadata to uniquely identify the client with the scope of the certificate authority. The client certificate would typically include fields for organization, state, city, building, etc.

Once again on the host holding the CA, create template files containing the information for each client, and use it to issue client certificates.

```
# cat > client-hostNNN.info <<EOF
country = GB
state = London
locality = City Of London
organization = Name of your organization
cn = hostNNN.foo.example.com
tls_www_client
encryption_key
signing_key
EOF
# certtool --generate-privkey > client-hostNNN-key.pem
# certtool --generate-certificate \
    --load-ca-certificate ca-cert.pem \
    --load-ca-privkey ca-key.pem \
```

(continues on next page)

(continued from previous page)

```
--load-privkey client-hostNNN-key.pem \
--template client-hostNNN.info \
--outfile client-hostNNN-cert.pem
```

The subject alt name extension data is not required for clients, so the `dns_name` and `ip_address` fields are not included. The `tls_www_client` keyword is the key purpose extension to indicate this certificate is intended for usage in a web client. Although QEMU network clients are not in fact HTTP clients, setting this key purpose is still recommended. The `encryption_key` and `signing_key` keyword is the key usage extension to indicate this certificate is intended for usage in the data session.

The `client-hostNNN-key.pem` and `client-hostNNN-cert.pem` files should now be securely copied to the client for which they were generated, and renamed to `client-key.pem` and `client-cert.pem` when added to the `/etc/pki/qemu` directory on the target host. The `client-key.pem` file is security sensitive and should be kept protected with file mode 0600 to prevent disclosure.

If a single host is going to be using TLS in both a client and server role, it is possible to create a single certificate to cover both roles. This would be quite common for the migration and NBD services, where a QEMU process will be started by accepting a TLS protected incoming migration, and later itself be migrated out to another host. To generate a single certificate, simply include the template data from both the client and server instructions in one.

```
# cat > both-hostNNN.info <<EOF
country = GB
state = London
locality = City Of London
organization = Name of your organization
cn = hostNNN.foo.example.com
dns_name = hostNNN
dns_name = hostNNN.foo.example.com
ip_address = 10.0.1.87
ip_address = 192.8.0.92
ip_address = 2620:0:cafe::87
ip_address = 2001:24::92
tls_www_server
tls_www_client
encryption_key
signing_key
EOF
# certtool --generate-privkey > both-hostNNN-key.pem
# certtool --generate-certificate \
--load-ca-certificate ca-cert.pem \
--load-ca-privkey ca-key.pem \
--load-privkey both-hostNNN-key.pem \
--template both-hostNNN.info \
--outfile both-hostNNN-cert.pem
```

When copying the PEM files to the target host, save them twice, once as `server-cert.pem` and `server-key.pem`, and again as `client-cert.pem` and `client-key.pem`.

1.13.4 TLS x509 credential configuration

QEMU has a standard mechanism for loading x509 credentials that will be used for network services and clients. It requires specifying the `tls-creds-x509` class name to the `--object` command line argument for the system emulators. Each set of credentials loaded should be given a unique string identifier via the `id` parameter. A single set of TLS credentials can be used for multiple network backends, so VNC, migration, NBD, character devices can all

share the same credentials. Note, however, that credentials for use in a client endpoint must be loaded separately from those used in a server endpoint.

When specifying the object, the `dir` parameters specifies which directory contains the credential files. This directory is expected to contain files with the names mentioned previously, `ca-cert.pem`, `server-key.pem`, `server-cert.pem`, `client-key.pem` and `client-cert.pem` as appropriate. It is also possible to include a set of pre-generated Diffie-Hellman (DH) parameters in a file `dh-params.pem`, which can be created using the `certtool --generate-dh-params` command. If omitted, QEMU will dynamically generate DH parameters when loading the credentials.

The `endpoint` parameter indicates whether the credentials will be used for a network client or server, and determines which PEM files are loaded.

The `verify` parameter determines whether x509 certificate validation should be performed. This defaults to enabled, meaning clients will always validate the server hostname against the certificate subject alt name fields and/or CN field. It also means that servers will request that clients provide a certificate and validate them. Verification should never be turned off for client endpoints, however, it may be turned off for server endpoints if an alternative mechanism is used to authenticate clients. For example, the VNC server can use SASL to authenticate clients instead.

To load server credentials with client certificate validation enabled

```
qemu-system-x86_64 -object tls-creds-x509,id=tls0,dir=/etc/pki/qemu,
↳endpoint=server
```

while to load client credentials use

```
qemu-system-x86_64 -object tls-creds-x509,id=tls0,dir=/etc/pki/qemu,
↳endpoint=client
```

Network services which support TLS will all have a `tls-creds` parameter which expects the ID of the TLS credentials object. For example with VNC:

```
qemu-system-x86_64 -vnc 0.0.0.0:0,tls-creds=tls0
```

1.13.5 TLS Pre-Shared Keys (PSK)

Instead of using certificates, you may also use TLS Pre-Shared Keys (TLS-PSK). This can be simpler to set up than certificates but is less scalable.

Use the GnuTLS `psktool` program to generate a `keys.psk` file containing one or more usernames and random keys:

```
mkdir -m 0700 /tmp/keys
psktool -u rich -p /tmp/keys/keys.psk
```

TLS-enabled servers such as `qemu-nbd` can use this directory like so:

```
qemu-nbd \
  -t -x / \
  --object tls-creds-psk,id=tls0,endpoint=server,dir=/tmp/keys \
  --tls-creds tls0 \
  image.qcow2
```

When connecting from a qemu-based client you must specify the directory containing `keys.psk` and an optional username (defaults to “qemu”):

```
qemu-img info \
  --object tls-creds-psk,id=tls0,dir=/tmp/keys,username=rich,endpoint=client \
```

(continues on next page)

(continued from previous page)

```
--image-opts \
file.driver=nbd,file.host=localhost,file.port=10809,file.tls-creds=tls0,file.
↪export=/
```

1.14 GDB usage

QEMU supports working with gdb via gdb's remote-connection facility (the "gdbstub"). This allows you to debug guest code in the same way that you might with a low-level debug facility like JTAG on real hardware. You can stop and start the virtual machine, examine state like registers and memory, and set breakpoints and watchpoints.

In order to use gdb, launch QEMU with the `-s` and `-S` options. The `-s` option will make QEMU listen for an incoming connection from gdb on TCP port 1234, and `-S` will make QEMU not start the guest until you tell it to from gdb. (If you want to specify which TCP port to use or to use something other than TCP for the gdbstub connection, use the `-gdb dev` option instead of `-s`.)

```
qemu-system-x86_64 -s -S -kernel bzImage -hda rootdisk.img -append "root=/dev/
↪hda"
```

QEMU will launch but will silently wait for gdb to connect.

Then launch gdb on the 'vmlinux' executable:

```
> gdb vmlinux
```

In gdb, connect to QEMU:

```
(gdb) target remote localhost:1234
```

Then you can use gdb normally. For example, type 'c' to launch the kernel:

```
(gdb) c
```

Here are some useful tips in order to use gdb on system code:

1. Use `info reg` to display all the CPU registers.
2. Use `x/10i $eip` to display the code at the PC position.
3. Use `set architecture i8086` to dump 16 bit code. Then use `x/10i $cs*16+$eip` to dump the code at the PC position.

Advanced debugging options:

The default single stepping behavior is step with the IRQs and timer service routines off. It is set this way because when gdb executes a single step it expects to advance beyond the current instruction. With the IRQs and timer service routines on, a single step might jump into the one of the interrupt or exception vectors instead of executing the current instruction. This means you may hit the same breakpoint a number of times before executing the instruction gdb wants to have executed. Because there are rare circumstances where you want to single step into an interrupt vector the behavior can be controlled from GDB. There are three commands you can query and set the single step behavior:

maintenance packet qqemu.sstepbits This will display the MASK bits used to control the single stepping IE:

```
(gdb) maintenance packet qqemu.sstepbits
sending: "qqemu.sstepbits"
received: "ENABLE=1,NOIRQ=2,NOTIMER=4"
```

maintenance packet qqemu.sstep This will display the current value of the mask used when single stepping IE:

```
(gdb) maintenance packet qqemu.sstep
sending: "qqemu.sstep"
received: "0x7"
```

maintenance packet Qqemu.sstep=HEX_VALUE This will change the single step mask, so if wanted to enable IRQs on the single step, but not timers, you would use:

```
(gdb) maintenance packet Qqemu.sstep=0x5
sending: "qemu.sstep=0x5"
received: "OK"
```

Another feature that QEMU gdbstub provides is to toggle the memory GDB works with, by default GDB will show the current process memory respecting the virtual address translation.

If you want to examine/change the physical memory you can set the gdbstub to work with the physical memory rather with the virtual one.

The memory mode can be checked by sending the following command:

maintenance packet qqemu.PhyMemMode This will return either 0 or 1, 1 indicates you are currently in the physical memory mode.

maintenance packet Qqemu.PhyMemMode:1 This will change the memory mode to physical memory.

maintenance packet Qqemu.PhyMemMode:0 This will change it back to normal memory mode.

1.15 Managed start up options

In system mode emulation, it's possible to create a VM in a paused state using the `-S` command line option. In this state the machine is completely initialized according to command line options and ready to execute VM code but VCPU threads are not executing any code. The VM state in this paused state depends on the way QEMU was started. It could be in:

- initial state (after reset/power on state)
- with direct kernel loading, the initial state could be amended to execute code loaded by QEMU in the VM's RAM and with incoming migration
- with incoming migration, initial state will be amended with the migrated machine state after migration completes

This paused state is typically used by users to query machine state and/or additionally configure the machine (by hotplugging devices) in runtime before allowing VM code to run.

However, at the `-S` pause point, it's impossible to configure options that affect initial VM creation (like: `-smp/-m/-numa ...`) or cold plug devices. The experimental `--preconfig` command line option allows pausing QEMU before the initial VM creation, in a "preconfig" state, where additional queries and configuration can be performed via QMP before moving on to the resulting configuration startup. In the preconfig state, QEMU only allows a limited set of commands over the QMP monitor, where the commands do not depend on an initialized machine, including but not limited to:

- `qmp_capabilities`
- `query-qmp-schema`
- `query-commands`
- `query-status`

- x-exit-preconfig

1.16 Virtual CPU hotplug

A complete example of vCPU hotplug (and hot-unplug) using QMP `device_add` and `device_del`.

1.16.1 vCPU hotplug

- (1) Launch QEMU as follows (note that the “maxcpus” is mandatory to allow vCPU hotplug):

```
$ qemu-system-x86_64 -display none -no-user-config -m 2048 \  
-nodefaults -monitor stdio -machine pc,accel=kvm,usb=off \  
-smp 1,maxcpus=2 -cpu IvyBridge-IBRS \  
-qmp unix:/tmp/qmp-sock,server,nowait
```

- (2) Run ‘qmp-shell’ (located in the source tree, under: “scripts/qmp/) to connect to the just-launched QEMU:

```
$> ./qmp-shell -p -v /tmp/qmp-sock  
[...]  
(QEMU)
```

- (3) Find out which CPU types could be plugged, and into which sockets:

```
(QEMU) query-hotpluggable-cpus  
{  
  "execute": "query-hotpluggable-cpus",  
  "arguments": {}  
}  
{  
  "return": [  
    {  
      "type": "IvyBridge-IBRS-x86_64-cpu",  
      "vcpus-count": 1,  
      "props": {  
        "socket-id": 1,  
        "core-id": 0,  
        "thread-id": 0  
      }  
    },  
    {  
      "qom-path": "/machine/unattached/device[0]",  
      "type": "IvyBridge-IBRS-x86_64-cpu",  
      "vcpus-count": 1,  
      "props": {  
        "socket-id": 0,  
        "core-id": 0,  
        "thread-id": 0  
      }  
    }  
  ]  
}  
(QEMU)
```

- (4) The `query-hotpluggable-cpus` command returns an object for CPUs that are present (containing a “qom-path” member) or which may be hot-plugged (no “qom-path” member). From its output in step (3),

we can see that IvyBridge-IBRS-x86_64-cpu is present in socket 0, while hot-plugging a CPU into socket 1 requires passing the listed properties to QMP `device_add`:

```
(QEMU) device_add id=cpu-2 driver=IvyBridge-IBRS-x86_64-cpu socket-id=1 core-id=0_
↪thread-id=0
{
  "execute": "device_add",
  "arguments": {
    "socket-id": 1,
    "driver": "IvyBridge-IBRS-x86_64-cpu",
    "id": "cpu-2",
    "core-id": 0,
    "thread-id": 0
  }
}
{
  "return": {}
}
(QEMU)
```

(5) Optionally, run QMP *query-cpus-fast* for some details about the vCPUs:

```
(QEMU) query-cpus-fast
{
  "execute": "query-cpus-fast",
  "arguments": {}
}
{
  "return": [
    {
      "qom-path": "/machine/unattached/device[0]",
      "target": "x86_64",
      "thread-id": 11534,
      "cpu-index": 0,
      "props": {
        "socket-id": 0,
        "core-id": 0,
        "thread-id": 0
      },
      "arch": "x86"
    },
    {
      "qom-path": "/machine/peripheral/cpu-2",
      "target": "x86_64",
      "thread-id": 12106,
      "cpu-index": 1,
      "props": {
        "socket-id": 1,
        "core-id": 0,
        "thread-id": 0
      },
      "arch": "x86"
    }
  ]
}
(QEMU)
```

1.16.2 vCPU hot-unplug

From the ‘qmp-shell’, invoke the QMP `device_del` command:

```
(QEMU) device_del id=cpu-2
{
  "execute": "device_del",
  "arguments": {
    "id": "cpu-2"
  }
}
{
  "return": {}
}
(QEMU)
```

Note: vCPU hot-unplug requires guest cooperation; so the `device_del` command above does not guarantee vCPU removal – it’s a “request to unplug”. At this point, the guest will get a System Control Interrupt (SCI) and calls the ACPI handler for the affected vCPU device. Then the guest kernel will bring the vCPU offline and tell QEMU to unplug it.

1.17 virtio pmem

This document explains the setup and usage of the virtio pmem device. The virtio pmem device is a paravirtualized persistent memory device on regular (i.e non-NVDIMM) storage.

1.17.1 Usecase

Virtio pmem allows to bypass the guest page cache and directly use host page cache. This reduces guest memory footprint as the host can make efficient memory reclaim decisions under memory pressure.

1.17.2 How does virtio-pmem compare to the nvdimm emulation?

NVDIMM emulation on regular (i.e. non-NVDIMM) host storage does not persist the guest writes as there are no defined semantics in the device specification. The virtio pmem device provides guest write persistence on non-NVDIMM host storage.

1.17.3 virtio pmem usage

A virtio pmem device backed by a memory-backend-file can be created on the QEMU command line as in the following example:

```
-object memory-backend-file,id=mem1,share,mem-path=./virtio_pmem.img,size=4G
-device virtio-pmem-pci,memdev=mem1,id=nv1
```

where:

- “object memory-backend-file,id=mem1,share,mem-path=<image>, size=<image size>” creates a backend file with the specified size.

- “device virtio-pmem-pci,id=nvdimm1,memdev=mem1” creates a virtio pmem pci device whose storage is provided by above memory backend device.

Multiple virtio pmem devices can be created if multiple pairs of “-object” and “-device” are provided.

1.17.4 Hotplug

Virtio pmem devices can be hotplugged via the QEMU monitor. First, the memory backing has to be added via ‘object_add’; afterwards, the virtio pmem device can be added via ‘device_add’.

For example, the following commands add another 4GB virtio pmem device to the guest:

```
(qemu) object_add memory-backend-file,id=mem2,share=on,mem-path=virtio_pmem2.img,
↪size=4G
(qemu) device_add virtio-pmem-pci,id=virtio_pmem2,memdev=mem2
```

1.17.5 Guest Data Persistence

Guest data persistence on non-NVDIMM requires guest userspace applications to perform fsync/msync. This is different from a real nvdimm backend where no additional fsync/msync is required. This is to persist guest writes in host backing file which otherwise remains in host page cache and there is risk of losing the data in case of power failure.

With virtio pmem device, MAP_SYNC mmap flag is not supported. This provides a hint to application to perform fsync for write persistence.

1.17.6 Limitations

- Real nvdimm device backend is not supported.
- virtio pmem hotunplug is not supported.
- ACPI NVDIMM features like regions/namespaces are not supported.
- ndctl command is not supported.

1.18 Persistent reservation managers

SCSI persistent reservations allow restricting access to block devices to specific initiators in a shared storage setup. When implementing clustering of virtual machines, it is a common requirement for virtual machines to send persistent reservation SCSI commands. However, the operating system restricts sending these commands to unprivileged programs because incorrect usage can disrupt regular operation of the storage fabric.

For this reason, QEMU’s SCSI passthrough devices, `scsi-block` and `scsi-generic` (both are only available on Linux) can delegate implementation of persistent reservations to a separate object, the “persistent reservation manager”. Only PERSISTENT RESERVE OUT and PERSISTENT RESERVE IN commands are passed to the persistent reservation manager object; other commands are processed by QEMU as usual.

1.18.1 Defining a persistent reservation manager

A persistent reservation manager is an instance of a subclass of the “pr-manager” QOM class.

Right now only one subclass is defined, `pr-manager-helper`, which forwards the commands to an external privileged helper program over Unix sockets. The helper program only allows sending persistent reservation commands to

devices for which QEMU has a file descriptor, so that QEMU will not be able to effect persistent reservations unless it has access to both the socket and the device.

`pr-manager-helper` has a single string property, `path`, which accepts the path to the helper program's Unix socket. For example, the following command line defines a `pr-manager-helper` object and attaches it to a SCSI passthrough device:

```
$ qemu-system-x86_64
  -device virtio-scsi \
  -object pr-manager-helper,id=helper0,path=/var/run/qemu-pr-helper.sock
  -drive if=none,id=hd,driver=raw,file.filename=/dev/sdb,file.pr-manager=helper0
  -device scsi-block,drive=hd
```

Alternatively, using `-blockdev`:

```
$ qemu-system-x86_64
  -device virtio-scsi \
  -object pr-manager-helper,id=helper0,path=/var/run/qemu-pr-helper.sock
  -blockdev node-name=hd,driver=raw,file.driver=host_device,file.filename=/dev/sdb,
  ↪file.pr-manager=helper0
  -device scsi-block,drive=hd
```

You will also need to ensure that the helper program **qemu-pr-helper** is running, and that it has been set up to use the same socket filename as your QEMU commandline specifies. See the `qemu-pr-helper` documentation or manpage for further details.

1.18.2 Multipath devices and persistent reservations

Proper support of persistent reservation for multipath devices requires communication with the multipath daemon, so that the reservation is registered and applied when a path is newly discovered or becomes online again. **qemu-pr-helper** can do this if the `libmpathpersist` library was available on the system at build time.

As of August 2017, a reservation key must be specified in `multipath.conf` for `multipathd` to check for persistent reservation for newly discovered paths or reinstated paths. The attribute can be added to the `defaults` section or the `multipaths` section; for example:

```
multipaths {
    multipath {
        wwid      XXXXXXXXXXXXXXXXX
        alias      yellow
        reservation_key  0x123abc
    }
}
```

Linking **qemu-pr-helper** to `libmpathpersist` does not impede its usage on regular SCSI devices.

1.19 QEMU System Emulator Targets

QEMU is a generic emulator and it emulates many machines. Most of the options are similar for all machines. Specific information about the various targets are mentioned in the following sections.

Contents:

1.19.1 x86 System emulator

Board-specific documentation

‘microvm’ virtual platform (microvm)

microvm is a machine type inspired by Firecracker and constructed after its machine model.

It’s a minimalist machine type without PCI nor ACPI support, designed for short-lived guests. microvm also establishes a baseline for benchmarking and optimizing both QEMU and guest operating systems, since it is optimized for both boot time and footprint.

Supported devices

The microvm machine type supports the following devices:

- ISA bus
- i8259 PIC (optional)
- i8254 PIT (optional)
- MC146818 RTC (optional)
- One ISA serial port (optional)
- LAPIC
- IOAPIC (with kernel-irqchip=split by default)
- kvmclock (if using KVM)
- fw_cfg
- Up to eight virtio-mmio devices (configured by the user)

Limitations

Currently, microvm does *not* support the following features:

- PCI-only devices.
- Hotplug of any kind.
- Live migration across QEMU versions.

Using the microvm machine type

Machine-specific options

It supports the following machine-specific options:

- microvm.x-option-roms=bool (Set off to disable loading option ROMs)
- microvm.pit=OnOffAuto (Enable i8254 PIT)
- microvm.isa-serial=bool (Set off to disable the instantiation an ISA serial port)
- microvm.pic=OnOffAuto (Enable i8259 PIC)

- `microvm.rtc=OnOffAuto` (Enable MC146818 RTC)
- `microvm.auto-kernel-cmdline=bool` (Set off to disable adding virtio-mmio devices to the kernel cmdline)

Boot options

By default, microvm uses `qboot` as its BIOS, to obtain better boot times, but it's also compatible with `SeaBIOS`.

As no current FW is able to boot from a block device using `virtio-mmio` as its transport, a microvm-based VM needs to be run using a host-side kernel and, optionally, an `initrd` image.

Running a microvm-based VM

By default, microvm aims for maximum compatibility, enabling both legacy and non-legacy devices. In this example, a VM is created without passing any additional machine-specific option, using the legacy ISA serial device as console:

```
$ qemu-system-x86_64 -M microvm \  
-enable-kvm -cpu host -m 512m -smp 2 \  
-kernel vmlinux -append "earlyprintk=ttyS0 console=ttyS0 root=/dev/vda" \  
-nodefaults -no-user-config -nographic \  
-serial stdio \  
-drive id=test,file=test.img,format=raw,if=none \  
-device virtio-blk-device,drive=test \  
-netdev tap,id=tap0,script=no,downscript=no \  
-device virtio-net-device,netdev=tap0
```

While the example above works, you might be interested in reducing the footprint further by disabling some legacy devices. If you're using KVM, you can disable the RTC, making the Guest rely on `kvmclock` exclusively. Additionally, if your host's CPUs have the `TSC_DEADLINE` feature, you can also disable both the i8259 PIC and the i8254 PIT (make sure you're also emulating a CPU with such feature in the guest).

This is an example of a VM with all optional legacy features disabled:

```
$ qemu-system-x86_64 \  
-M microvm,x-option-roms=off,pit=off,pic=off,isa-serial=off,rtc=off \  
-enable-kvm -cpu host -m 512m -smp 2 \  
-kernel vmlinux -append "console=hvc0 root=/dev/vda" \  
-nodefaults -no-user-config -nographic \  
-chardev stdio,id=virtiocon0 \  
-device virtio-serial-device \  
-device virtconsole,chardev=virtiocon0 \  
-drive id=test,file=test.img,format=raw,if=none \  
-device virtio-blk-device,drive=test \  
-netdev tap,id=tap0,script=no,downscript=no \  
-device virtio-net-device,netdev=tap0
```

Triggering a guest-initiated shut down

As the microvm machine type includes just a small set of system devices, some x86 mechanisms for rebooting or shutting down the system, like sending a key sequence to the keyboard or writing to an ACPI register, doesn't have any effect in the VM.

The recommended way to trigger a guest-initiated shut down is by generating a `triple-fault`, which will cause the VM to initiate a reboot. Additionally, if the `-no-reboot` argument is present in the command line, QEMU will detect this event and terminate its own execution gracefully.

Linux does support this mechanism, but by default will only be used after other options have been tried and failed, causing the reboot to be delayed by a small number of seconds. It's possible to instruct it to try the triple-fault mechanism first, by adding `reboot=t` to the kernel's command line.

i440fx PC (`pc-i440fx`, `pc`)

Peripherals

The QEMU PC System emulator simulates the following peripherals:

- i440FX host PCI bridge and PIIX3 PCI to ISA bridge
- Cirrus CLGD 5446 PCI VGA card or dummy VGA card with Bochs VESA extensions (hardware level, including all non standard modes).
- PS/2 mouse and keyboard
- 2 PCI IDE interfaces with hard disk and CD-ROM support
- Floppy disk
- PCI and ISA network adapters
- Serial ports
- IPMI BMC, either and internal or external one
- Creative SoundBlaster 16 sound card
- ENSONIQ AudioPCI ES1370 sound card
- Intel 82801AA AC97 Audio compatible sound card
- Intel HD Audio Controller and HDA codec
- Adlib (OPL2) - Yamaha YM3812 compatible chip
- Gravis Ultrasound GF1 sound card
- CS4231A compatible sound card
- PC speaker
- PCI UHCI, OHCI, EHCI or XHCI USB controller and a virtual USB-1.1 hub.

SMP is supported with up to 255 CPUs.

QEMU uses the PC BIOS from the Seabios project and the Plex86/Bochs LGPL VGA BIOS.

QEMU uses YM3812 emulation by Tatsuyuki Satoh.

QEMU uses GUS emulation (GUSEMU32 <http://www.deinmeister.de/gusemu/>) by Tibor "TS" Schütz.

Note that, by default, GUS shares IRQ(7) with parallel ports and so QEMU must be told to not have parallel ports to have working GUS.

```
qemu_system-x86_64 dos.img -device gus -parallel none
```

Alternatively:

```
qemu_system-x86_64 dos.img -device gus,irq=5
```

Or some other unclaimed IRQ.

CS4231A is the chip used in Windows Sound System and GUSMAX products

The PC speaker audio device can be configured using the `pcspk-audiodev` machine property, i.e.

```
qemu_system-x86_64 some.img -audiodev <backend>,id=<name> -machine pcspk-  
↪audiodev=<name>
```

Recommendations for KVM CPU model configuration on x86 hosts

The information that follows provides recommendations for configuring CPU models on x86 hosts. The goals are to maximise performance, while protecting guest OS against various CPU hardware flaws, and optionally enabling live migration between hosts with heterogeneous CPU models.

Two ways to configure CPU models with QEMU / KVM

(1) Host passthrough

This passes the host CPU model features, model, stepping, exactly to the guest. Note that KVM may filter out some host CPU model features if they cannot be supported with virtualization. Live migration is unsafe when this mode is used as libvirt / QEMU cannot guarantee a stable CPU is exposed to the guest across hosts. This is the recommended CPU to use, provided live migration is not required.

(2) Named model

QEMU comes with a number of predefined named CPU models, that typically refer to specific generations of hardware released by Intel and AMD. These allow the guest VMs to have a degree of isolation from the host CPU, allowing greater flexibility in live migrating between hosts with differing hardware. @end table

In both cases, it is possible to optionally add or remove individual CPU features, to alter what is presented to the guest by default.

Libvirt supports a third way to configure CPU models known as “Host model”. This uses the QEMU “Named model” feature, automatically picking a CPU model that is similar the host CPU, and then adding extra features to approximate the host model as closely as possible. This does not guarantee the CPU family, stepping, etc will precisely match the host CPU, as they would with “Host passthrough”, but gives much of the benefit of passthrough, while making live migration safe.

Preferred CPU models for Intel x86 hosts

The following CPU models are preferred for use on Intel hosts. Administrators / applications are recommended to use the CPU model that matches the generation of the host CPUs in use. In a deployment with a mixture of host CPU models between machines, if live migration compatibility is required, use the newest CPU model that is compatible across all desired hosts.

Cascadelake-Server, Cascadelake-Server-noTSX Intel Xeon Processor (Cascade Lake, 2019), with “stepping” levels 6 or 7 only. (The Cascade Lake Xeon processor with *stepping 5 is vulnerable to MDS variants.*)

Skylake-Server, Skylake-Server-IBRS, Skylake-Server-IBRS-noTSX Intel Xeon Processor (Skylake, 2016)

Skylake-Client, Skylake-Client-IBRS, Skylake-Client-noTSX-IBRS Intel Core Processor (Skylake, 2015)

Broadwell, Broadwell-IBRS, Broadwell-noTSX, Broadwell-noTSX-IBRS Intel Core Processor (Broadwell, 2014)

Haswell, Haswell-IBRS, Haswell-noTSX, Haswell-noTSX-IBRS Intel Core Processor (Haswell, 2013)

IvyBridge, IvyBridge-IBR Intel Xeon E3-12xx v2 (Ivy Bridge, 2012)

SandyBridge, SandyBridge-IBRS Intel Xeon E312xx (Sandy Bridge, 2011)

Westmere, Westmere-IBRS Westmere E56xx/L56xx/X56xx (Nehalem-C, 2010)

Nehalem, Nehalem-IBRS Intel Core i7 9xx (Nehalem Class Core i7, 2008)

Penryn Intel Core 2 Duo P9xxx (Penryn Class Core 2, 2007)

Conroe Intel Celeron_4x0 (Conroe/Merom Class Core 2, 2006)

Important CPU features for Intel x86 hosts

The following are important CPU features that should be used on Intel x86 hosts, when available in the host CPU. Some of them require explicit configuration to enable, as they are not included by default in some, or all, of the named CPU models listed above. In general all of these features are included if using “Host passthrough” or “Host model”.

pcid Recommended to mitigate the cost of the Meltdown (CVE-2017-5754) fix.

Included by default in Haswell, Broadwell & Skylake Intel CPU models.

Should be explicitly turned on for Westmere, SandyBridge, and IvyBridge Intel CPU models. Note that some desktop/mobile Westmere CPUs cannot support this feature.

spec-ctrl Required to enable the Spectre v2 (CVE-2017-5715) fix.

Included by default in Intel CPU models with -IBRS suffix.

Must be explicitly turned on for Intel CPU models without -IBRS suffix.

Requires the host CPU microcode to support this feature before it can be used for guest CPUs.

stibp Required to enable stronger Spectre v2 (CVE-2017-5715) fixes in some operating systems.

Must be explicitly turned on for all Intel CPU models.

Requires the host CPU microcode to support this feature before it can be used for guest CPUs.

ssbd Required to enable the CVE-2018-3639 fix.

Not included by default in any Intel CPU model.

Must be explicitly turned on for all Intel CPU models.

Requires the host CPU microcode to support this feature before it can be used for guest CPUs.

pdpe1gb Recommended to allow guest OS to use 1GB size pages.

Not included by default in any Intel CPU model.

Should be explicitly turned on for all Intel CPU models.

Note that not all CPU hardware will support this feature.

md-clear Required to confirm the MDS (CVE-2018-12126, CVE-2018-12127, CVE-2018-12130, CVE-2019-11091) fixes.

Not included by default in any Intel CPU model.

Must be explicitly turned on for all Intel CPU models.

Requires the host CPU microcode to support this feature before it can be used for guest CPUs.

mds-no Recommended to inform the guest OS that the host is *not* vulnerable to any of the MDS variants ([MFBDS] CVE-2018-12130, [MLPDS] CVE-2018-12127, [MSBDS] CVE-2018-12126).

This is an MSR (Model-Specific Register) feature rather than a CPUID feature, so it will not appear in the Linux `/proc/cpuinfo` in the host or guest. Instead, the host kernel uses it to populate the MDS vulnerability file in `sysfs`.

So it should only be enabled for VMs if the host reports `@code{Not affected}` in the `/sys/devices/system/cpu/vulnerabilities/mds` file.

taa-no Recommended to inform that the guest that the host is *not* vulnerable to CVE-2019-11135, TSX Asynchronous Abort (TAA).

This too is an MSR feature, so it does not show up in the Linux `/proc/cpuinfo` in the host or guest.

It should only be enabled for VMs if the host reports `Not affected` in the `/sys/devices/system/cpu/vulnerabilities/tsx_async_abort` file.

tsx-ctrl Recommended to inform the guest that it can disable the Intel TSX (Transactional Synchronization Extensions) feature; or, if the processor is vulnerable, use the Intel VERW instruction (a processor-level instruction that performs checks on memory access) as a mitigation for the TAA vulnerability. (For details, refer to Intel's [deep dive into MDS](#).)

Expose this to the guest OS if and only if: (a) the host has TSX enabled; *and* (b) the guest has `rtm` CPU flag enabled.

By disabling TSX, KVM-based guests can avoid paying the price of mitigating TSX-based attacks.

Note that `tsx-ctrl` too is an MSR feature, so it does not show up in the Linux `/proc/cpuinfo` in the host or guest.

To validate that Intel TSX is indeed disabled for the guest, there are two ways: (a) check for the *absence* of `rtm` in the guest's `/proc/cpuinfo`; or (b) the `/sys/devices/system/cpu/vulnerabilities/tsx_async_abort` file in the guest should report `Mitigation: TSX disabled`.

Preferred CPU models for AMD x86 hosts

The following CPU models are preferred for use on Intel hosts. Administrators / applications are recommended to use the CPU model that matches the generation of the host CPUs in use. In a deployment with a mixture of host CPU models between machines, if live migration compatibility is required, use the newest CPU model that is compatible across all desired hosts.

EPYC, EPYC-IBPB AMD EPYC Processor (2017)

Opteron_G5 AMD Opteron 63xx class CPU (2012)

Opteron_G4 AMD Opteron 62xx class CPU (2011)

Opteron_G3 AMD Opteron 23xx (Gen 3 Class Opteron, 2009)

Opteron_G2 AMD Opteron 22xx (Gen 2 Class Opteron, 2006)

Opteron_G1 AMD Opteron 240 (Gen 1 Class Opteron, 2004)

Important CPU features for AMD x86 hosts

The following are important CPU features that should be used on AMD x86 hosts, when available in the host CPU. Some of them require explicit configuration to enable, as they are not included by default in some, or all, of the named CPU models listed above. In general all of these features are included if using “Host passthrough” or “Host model”.

ibpb Required to enable the Spectre v2 (CVE-2017-5715) fix.

Included by default in AMD CPU models with -IBPB suffix.

Must be explicitly turned on for AMD CPU models without -IBPB suffix.

Requires the host CPU microcode to support this feature before it can be used for guest CPUs.

stibp Required to enable stronger Spectre v2 (CVE-2017-5715) fixes in some operating systems.

Must be explicitly turned on for all AMD CPU models.

Requires the host CPU microcode to support this feature before it can be used for guest CPUs.

virt-ssbd Required to enable the CVE-2018-3639 fix

Not included by default in any AMD CPU model.

Must be explicitly turned on for all AMD CPU models.

This should be provided to guests, even if amd-ssbd is also provided, for maximum guest compatibility.

Note for some QEMU / libvirt versions, this must be force enabled when when using “Host model”, because this is a virtual feature that doesn’t exist in the physical host CPUs.

amd-ssbd Required to enable the CVE-2018-3639 fix

Not included by default in any AMD CPU model.

Must be explicitly turned on for all AMD CPU models.

This provides higher performance than `virt-ssbd` so should be exposed to guests whenever available in the host. `virt-ssbd` should none the less also be exposed for maximum guest compatibility as some kernels only know about `virt-ssbd`.

amd-no-ssb Recommended to indicate the host is not vulnerable CVE-2018-3639

Not included by default in any AMD CPU model.

Future hardware generations of CPU will not be vulnerable to CVE-2018-3639, and thus the guest should be told not to enable its mitigations, by exposing `amd-no-ssb`. This is mutually exclusive with `virt-ssbd` and `amd-ssbd`.

pdpe1gb Recommended to allow guest OS to use 1GB size pages

Not included by default in any AMD CPU model.

Should be explicitly turned on for all AMD CPU models.

Note that not all CPU hardware will support this feature.

Default x86 CPU models

The default QEMU CPU models are designed such that they can run on all hosts. If an application does not wish to do perform any host compatibility checks before launching guests, the default is guaranteed to work.

The default CPU models will, however, leave the guest OS vulnerable to various CPU hardware flaws, so their use is strongly discouraged. Applications should follow the earlier guidance to setup a better CPU configuration, with host passthrough recommended if live migration is not needed.

qemu32, qemu64 QEMU Virtual CPU version 2.5+ (32 & 64 bit variants)

`qemu64` is used for `x86_64` guests and `qemu32` is used for `i686` guests, when no `-cpu` argument is given to QEMU, or no `<cpu>` is provided in libvirt XML.

Other non-recommended x86 CPUs

The following CPUs models are compatible with most AMD and Intel x86 hosts, but their usage is discouraged, as they expose a very limited featureset, which prevents guests having optimal performance.

kvm32, kvm64 Common KVM processor (32 & 64 bit variants).

Legacy models just for historical compatibility with ancient QEMU versions.

486, athlon, phenom, coreduo, core2duo, n270, pentium, pentium2, pentium3 Various very old x86 CPU models, mostly predating the introduction of hardware assisted virtualization, that should thus not be required for running virtual machines.

Syntax for configuring CPU models

The examples below illustrate the approach to configuring the various CPU models / features in QEMU and libvirt.

QEMU command line

Host passthrough:

```
qemu-system-x86_64 -cpu host
```

Host passthrough with feature customization:

```
qemu-system-x86_64 -cpu host,-vmx,...
```

Named CPU models:

```
qemu-system-x86_64 -cpu Westmere
```

Named CPU models with feature customization:

```
qemu-system-x86_64 -cpu Westmere,+pcid,...
```

Libvirt guest XML

Host passthrough:

```
<cpu mode='host-passthrough'/>
```

Host passthrough with feature customization:

```
<cpu mode='host-passthrough'>
  <feature name="vmx" policy="disable"/>
  ...
</cpu>
```

Host model:

```
<cpu mode='host-model'/>
```

Host model with feature customization:

```
<cpu mode='host-model'>
  <feature name="vmx" policy="disable"/>
  ...
</cpu>
```

Named model:

```
<cpu mode='custom'>
  <model name="Westmere"/>
</cpu>
```

Named model with feature customization:

```
<cpu mode='custom'>
  <model name="Westmere"/>
  <feature name="pcid" policy="require"/>
  ...
</cpu>
```

OS requirements

On x86_64 hosts, the default set of CPU features enabled by the KVM accelerator require the host to be running Linux v4.5 or newer. Red Hat Enterprise Linux 7 is also supported, since the required functionality was backported.

1.19.2 PowerPC System emulator

Use the executable `qemu-system-ppc` to simulate a complete 40P (PREP) or PowerMac PowerPC system.

QEMU emulates the following PowerMac peripherals:

- UniNorth or Grackle PCI Bridge
- PCI VGA compatible card with VESA Bochs Extensions
- 2 PMAC IDE interfaces with hard disk and CD-ROM support
- NE2000 PCI adapters
- Non Volatile RAM
- VIA-CUDA with ADB keyboard and mouse.

QEMU emulates the following 40P (PREP) peripherals:

- PCI Bridge
- PCI VGA compatible card with VESA Bochs Extensions
- 2 IDE interfaces with hard disk and CD-ROM support
- Floppy disk
- PCnet network adapters
- Serial port
- PREP Non Volatile RAM
- PC compatible keyboard and mouse.

Since version 0.9.1, QEMU uses OpenBIOS <https://www.openbios.org/> for the g3beige and mac99 PowerMac and the 40p machines. OpenBIOS is a free (GPL v2) portable firmware implementation. The goal is to implement a 100% IEEE 1275-1994 (referred to as Open Firmware) compliant firmware.

More information is available at http://perso.magic.fr/l_indien/qemu-ppc/.

1.19.3 Sparc32 System emulator

Use the executable `qemu-system-sparc` to simulate the following Sun4m architecture machines:

- SPARCstation 4
- SPARCstation 5
- SPARCstation 10
- SPARCstation 20
- SPARCserver 600MP
- SPARCstation LX
- SPARCstation Voyager
- SPARCclassic
- SPARCbook

The emulation is somewhat complete. SMP up to 16 CPUs is supported, but Linux limits the number of usable CPUs to 4.

QEMU emulates the following sun4m peripherals:

- IOMMU
- TCX or cgthree Frame buffer
- Lance (Am7990) Ethernet
- Non Volatile RAM M48T02/M48T08
- Slave I/O: timers, interrupt controllers, Zilog serial ports, keyboard and power/reset logic
- ESP SCSI controller with hard disk and CD-ROM support
- Floppy drive (not on SS-600MP)
- CS4231 sound device (only on SS-5, not working yet)

The number of peripherals is fixed in the architecture. Maximum memory size depends on the machine type, for SS-5 it is 256MB and for others 2047MB.

Since version 0.8.2, QEMU uses OpenBIOS <https://www.openbios.org/>. OpenBIOS is a free (GPL v2) portable firmware implementation. The goal is to implement a 100% IEEE 1275-1994 (referred to as Open Firmware) compliant firmware.

A sample Linux 2.6 series kernel and ram disk image are available on the QEMU web site. There are still issues with NetBSD and OpenBSD, but most kernel versions work. Please note that currently older Solaris kernels don't work probably due to interface issues between OpenBIOS and Solaris.

1.19.4 Sparc64 System emulator

Use the executable `qemu-system-sparc64` to simulate a Sun4u (UltraSPARC PC-like machine), Sun4v (T1 PC-like machine), or generic Niagara (T1) machine. The Sun4u emulator is mostly complete, being able to run Linux, NetBSD and OpenBSD in headless (`-nographic`) mode. The Sun4v emulator is still a work in progress.

The Niagara T1 emulator makes use of firmware and OS binaries supplied in the `S10image/` directory of the OpenSPARC T1 project http://download.oracle.com/technetwork/systems/opensparc/OpenSPARCT1_Arch.1.5.tar.bz2 and is able to boot the `disk.s10hw2` Solaris image.

```
qemu-system-sparc64 -M niagara -L /path-to/S10image/ \
                    -nographic -m 256 \
                    -drive if=pflash,readonly=on,file=/S10image/disk.s10hw2
```

QEMU emulates the following peripherals:

- UltraSparc Iii APB PCI Bridge
- PCI VGA compatible card with VESA Bochs Extensions
- PS/2 mouse and keyboard
- Non Volatile RAM M48T59
- PC-compatible serial ports
- 2 PCI IDE interfaces with hard disk and CD-ROM support
- Floppy disk

1.19.5 MIPS System emulator

Four executables cover simulation of 32 and 64-bit MIPS systems in both endian options, `qemu-system-mips`, `qemu-system-mipsel`, `qemu-system-mips64` and `qemu-system-mips64el`. Five different machine types are emulated:

- A generic ISA PC-like machine “mips”
- The MIPS Malta prototype board “malta”
- An ACER Pica “pica61”. This machine needs the 64-bit emulator.
- MIPS emulator pseudo board “mipssim”
- A MIPS Magnum R4000 machine “magnum”. This machine needs the 64-bit emulator.

The generic emulation is supported by Debian ‘Etch’ and is able to install Debian into a virtual disk image. The following devices are emulated:

- A range of MIPS CPUs, default is the 24Kf
- PC style serial port
- PC style IDE disk
- NE2000 network card

The Malta emulation supports the following devices:

- Core board with MIPS 24Kf CPU and Galileo system controller
- PIIX4 PCI/USB/SMBus controller
- The Multi-I/O chip’s serial device

- PCI network cards (PCnet32 and others)
- Malta FPGA serial device
- Cirrus (default) or any other PCI VGA graphics card

The Boston board emulation supports the following devices:

- Xilinx FPGA, which includes a PCIe root port and an UART
- Intel EG20T PCH connects the I/O peripherals, but only the SATA bus is emulated

The ACER Pica emulation supports:

- MIPS R4000 CPU
- PC-style IRQ and DMA controllers
- PC Keyboard
- IDE controller

The MIPS Magnum R4000 emulation supports:

- MIPS R4000 CPU
- PC-style IRQ controller
- PC Keyboard
- SCSI controller
- G364 framebuffer

The Fuloong 2E emulation supports:

- Loongson 2E CPU
- Bonito64 system controller as North Bridge
- VT82C686 chipset as South Bridge
- RTL8139D as a network card chipset

The Loongson-3 virtual platform emulation supports:

- Loongson 3A CPU
- LIOINTC as interrupt controller
- GPEX and virtio as peripheral devices
- Both KVM and TCG supported

The mipssim pseudo board emulation provides an environment similar to what the proprietary MIPS emulator uses for running Linux. It supports:

- A range of MIPS CPUs, default is the 24Kf
- PC style serial port
- MIPSnet network emulation

Supported CPU model configurations on MIPS hosts

QEMU supports variety of MIPS CPU models:

Supported CPU models for MIPS32 hosts

The following CPU models are supported for use on MIPS32 hosts. Administrators / applications are recommended to use the CPU model that matches the generation of the host CPUs in use. In a deployment with a mixture of host CPU models between machines, if live migration compatibility is required, use the newest CPU model that is compatible across all desired hosts.

mips32r6-generic MIPS32 Processor (Release 6, 2015)
P5600 MIPS32 Processor (P5600, 2014)
M14K, M14Kc MIPS32 Processor (M14K, 2009)
74Kf MIPS32 Processor (74K, 2007)
34Kf MIPS32 Processor (34K, 2006)
24Kc, 24KEc, 24Kf MIPS32 Processor (24K, 2003)
4Kc, 4Km, 4KEcR1, 4KEmR1, 4KEc, 4KEm MIPS32 Processor (4K, 1999)

Supported CPU models for MIPS64 hosts

The following CPU models are supported for use on MIPS64 hosts. Administrators / applications are recommended to use the CPU model that matches the generation of the host CPUs in use. In a deployment with a mixture of host CPU models between machines, if live migration compatibility is required, use the newest CPU model that is compatible across all desired hosts.

I6400 MIPS64 Processor (Release 6, 2014)
Loongson-2E MIPS64 Processor (Loongson 2, 2006)
Loongson-2F MIPS64 Processor (Loongson 2, 2008)
Loongson-3A1000 MIPS64 Processor (Loongson 3, 2010)
Loongson-3A4000 MIPS64 Processor (Loongson 3, 2018)
mips64dspr2 MIPS64 Processor (Release 2, 2006)
MIPS64R2-generic, 5KEc, 5KEf MIPS64 Processor (Release 2, 2002)
20Kc MIPS64 Processor (20K, 2000)
5Kc, 5Kf MIPS64 Processor (5K, 1999)
VR5432 MIPS64 Processor (VR, 1998)
R4000 MIPS64 Processor (MIPS III, 1991)

Supported CPU models for nanoMIPS hosts

The following CPU models are supported for use on nanoMIPS hosts. Administrators / applications are recommended to use the CPU model that matches the generation of the host CPUs in use. In a deployment with a mixture of host CPU models between machines, if live migration compatibility is required, use the newest CPU model that is compatible across all desired hosts.

I7200 MIPS I7200 (nanoMIPS, 2018)

Preferred CPU models for MIPS hosts

The following CPU models are preferred for use on different MIPS hosts:

MIPS III R4000

MIPS32R2 34Kf

MIPS64R6 I6400

nanoMIPS I7200

nanoMIPS System emulator

Executable `qemu-system-mipsel` also covers simulation of 32-bit nanoMIPS system in little endian mode:

- nanoMIPS I7200 CPU

Example of `qemu-system-mipsel` usage for nanoMIPS is shown below:

Download `<disk_image_file>` from <https://mipsdistros.mips.com/LinuxDistro/nanomips/buildroot/index.html>.

Download `<kernel_image_file>` from <https://mipsdistros.mips.com/LinuxDistro/nanomips/kernels/v4.15.18-432-gb2eb9a8b07a1-20180627102142/index.html>.

Start system emulation of Malta board with nanoMIPS I7200 CPU:

```
qemu-system-mipsel -cpu I7200 -kernel <kernel_image_file> \  
-M malta -serial stdio -m <memory_size> -hda <disk_image_file> \  
-append "mem=256m@0x0 rw console=ttyS0 vga=cirrus vesa=0x111 root=/dev/sda"
```

1.19.6 Arm System emulator

QEMU can emulate both 32-bit and 64-bit Arm CPUs. Use the `qemu-system-aarch64` executable to simulate a 64-bit Arm machine. You can use either `qemu-system-arm` or `qemu-system-aarch64` to simulate a 32-bit Arm machine: in general, command lines that work for `qemu-system-arm` will behave the same when used with `qemu-system-aarch64`.

QEMU has generally good support for Arm guests. It has support for nearly fifty different machines. The reason we support so many is that Arm hardware is much more widely varying than x86 hardware. Arm CPUs are generally built into “system-on-chip” (SoC) designs created by many different companies with different devices, and these SoCs are then built into machines which can vary still further even if they use the same SoC. Even with fifty boards QEMU does not cover more than a small fraction of the Arm hardware ecosystem.

The situation for 64-bit Arm is fairly similar, except that we don’t implement so many different machines.

As well as the more common “A-profile” CPUs (which have MMUs and will run Linux) QEMU also supports “M-profile” CPUs such as the Cortex-M0, Cortex-M4 and Cortex-M33 (which are microcontrollers used in very embedded boards). For most boards the CPU type is fixed (matching what the hardware has), so typically you don’t need to specify the CPU type by hand, except for special cases like the `virt` board.

Choosing a board model

For QEMU’s Arm system emulation, you must specify which board model you want to use with the `-M` or `--machine` option; there is no default.

Because Arm systems differ so much and in fundamental ways, typically operating system or firmware images intended to run on one machine will not run at all on any other. This is often surprising for new users who are used to the x86

world where every system looks like a standard PC. (Once the kernel has booted, most userspace software cares much less about the detail of the hardware.)

If you already have a system image or a kernel that works on hardware and you want to boot with QEMU, check whether QEMU lists that machine in its `--machine help` output. If it is listed, then you can probably use that board model. If it is not listed, then unfortunately your image will almost certainly not boot on QEMU. (You might be able to extract the filesystem and use that with a different kernel which boots on a system that QEMU does emulate.)

If you don't care about reproducing the idiosyncrasies of a particular bit of hardware, such as small amount of RAM, no PCI or other hard disk, etc., and just want to run Linux, the best option is to use the `virt` board. This is a platform which doesn't correspond to any real hardware and is designed for use in virtual machines. You'll need to compile Linux with a suitable configuration for running on the `virt` board. `virt` supports PCI, virtio, recent CPUs and large amounts of RAM. It also supports 64-bit CPUs.

Board-specific documentation

Unfortunately many of the Arm boards QEMU supports are currently undocumented; you can get a complete list by running `qemu-system-aarch64 --machine help`.

Arm Integrator/CP (`integratorcp`)

The Arm Integrator/CP board is emulated with the following devices:

- ARM926E, ARM1026E, ARM946E, ARM1136 or Cortex-A8 CPU
- Two PL011 UARTs
- SMC 91c111 Ethernet adapter
- PL110 LCD controller
- PL050 KMI with PS/2 keyboard and mouse.
- PL181 MultiMedia Card Interface with SD card.

Arm MPS2 boards (`mps2-an385`, `mps2-an386`, `mps2-an500`, `mps2-an505`, `mps2-an511`, `mps2-an521`)

These board models all use Arm M-profile CPUs.

The Arm MPS2 and MPS2+ dev boards are FPGA based (the 2+ has a bigger FPGA but is otherwise the same as the 2). Since the CPU itself and most of the devices are in the FPGA, the details of the board as seen by the guest depend significantly on the FPGA image.

QEMU models the following FPGA images:

`mps2-an385` Cortex-M3 as documented in Arm Application Note AN385

`mps2-an386` Cortex-M4 as documented in Arm Application Note AN386

`mps2-an500` Cortex-M7 as documented in Arm Application Note AN500

`mps2-an505` Cortex-M33 as documented in Arm Application Note AN505

`mps2-an511` Cortex-M3 'DesignStart' as documented in Arm Application Note AN511

`mps2-an521` Dual Cortex-M33 as documented in Arm Application Note AN521

Differences between QEMU and real hardware:

- AN385/AN386 remapping of low 16K of memory to either ZBT SSRAM1 or to block RAM is unimplemented (QEMU always maps this to ZBT SSRAM1, as if `zbt_boot_ctrl` is always zero)
- QEMU provides a LAN9118 ethernet rather than LAN9220; the only guest visible difference is that the LAN9118 doesn't support checksum offloading

Arm Musca boards (`musca-a`, `musca-b1`)

The Arm Musca development boards are a reference implementation of a system using the SSE-200 Subsystem for Embedded. They are dual Cortex-M33 systems.

QEMU provides models of the A and B1 variants of this board.

Unimplemented devices:

- SPI
- I²C
- I²S
- PWM
- QSPI
- Timer
- SCC
- GPIO
- eFlash
- MHU
- PVT
- SDIO
- CryptoCell

Note that (like the real hardware) the Musca-A machine is asymmetric: CPU 0 does not have the FPU or DSP extensions, but CPU 1 does. Also like the real hardware, the memory maps for the A and B1 variants differ significantly, so guest software must be built for the right variant.

Arm Realview boards (`realview-eb`, `realview-eb-mpcore`, `realview-pb-a8`, `realview-pbx-a9`)

Several variants of the Arm RealView baseboard are emulated, including the EB, PB-A8 and PBX-A9. Due to interactions with the bootloader, only certain Linux kernel configurations work out of the box on these boards.

Kernels for the PB-A8 board should have `CONFIG_REALVIEW_HIGH_PHYS_OFFSET` enabled in the kernel, and expect 512M RAM. Kernels for The PBX-A9 board should have `CONFIG_SPARSEMEM` enabled, `CONFIG_REALVIEW_HIGH_PHYS_OFFSET` disabled and expect 1024M RAM.

The following devices are emulated:

- ARM926E, ARM1136, ARM11MPCore, Cortex-A8 or Cortex-A9 MPCore CPU
- Arm AMBA Generic/Distributed Interrupt Controller
- Four PL011 UARTs

- SMC 91c111 or SMSC LAN9118 Ethernet adapter
- PL110 LCD controller
- PL050 KMI with PS/2 keyboard and mouse
- PCI host bridge
- PCI OHCI USB controller
- LSI53C895A PCI SCSI Host Bus Adapter with hard disk and CD-ROM devices
- PL181 MultiMedia Card Interface with SD card.

Arm Server Base System Architecture Reference board (*sbsa-ref*)

While the *virt* board is a generic board platform that doesn't match any real hardware the *sbsa-ref* board intends to look like real hardware. The *Server Base System Architecture* <<https://developer.arm.com/documentation/den0029/latest>> defines a minimum base line of hardware support and importantly how the firmware reports that to any operating system. It is a static system that reports a very minimal DT to the firmware for non-discoverable information about components affected by the qemu command line (i.e. cpus and memory). As a result it must have a firmware specifically built to expect a certain hardware layout (as you would in a real machine).

It is intended to be a machine for developing firmware and testing standards compliance with operating systems.

Supported devices

The *sbsa-ref* board supports:

- A configurable number of AArch64 CPUs
- GIC version 3
- System bus AHCI controller
- System bus EHCI controller
- CDROM and hard disc on AHCI bus
- E1000E ethernet card on PCIe bus
- VGA display adaptor on PCIe bus
- A generic SBSA watchdog device

Arm Versatile boards (*versatileab*, *versatilepb*)

The Arm Versatile baseboard is emulated with the following devices:

- ARM926E, ARM1136 or Cortex-A8 CPU
- PL190 Vectored Interrupt Controller
- Four PL011 UARTs
- SMC 91c111 Ethernet adapter
- PL110 LCD controller
- PL050 KMI with PS/2 keyboard and mouse.

- PCI host bridge. Note the emulated PCI bridge only provides access to PCI memory space. It does not provide access to PCI IO space. This means some devices (eg. ne2k_pci NIC) are not usable, and others (eg. rtl8139 NIC) are only usable when the guest drivers use the memory mapped control registers.
- PCI OHCI USB controller.
- LSI53C895A PCI SCSI Host Bus Adapter with hard disk and CD-ROM devices.
- PL181 MultiMedia Card Interface with SD card.

Booting a Linux kernel

Building a current Linux kernel with `versatile_defconfig` should be enough to get something running. Nowadays an out-of-tree build is recommended (and also useful if you build a lot of different targets). In the following example `$BLD` points to the build directory and `$SRC` points to the root of the Linux source tree. You can drop `$SRC` if you are running from there.

```
$ make O=$BLD -C $SRC ARCH=arm CROSS_COMPILE=arm-linux-gnueabihf- versatile_defconfig
$ make O=$BLD -C $SRC ARCH=arm CROSS_COMPILE=arm-linux-gnueabihf-
```

You may want to enable some additional modules if you want to boot something from the SCSI interface:

```
CONFIG_PCI=y
CONFIG_PCI_VERSATILE=y
CONFIG_SCSI=y
CONFIG_SCSI_SYM53C8XX_2=y
```

You can then boot with a command line like:

```
$ qemu-system-arm -machine type=versatilepb \
  -serial mon:stdio \
  -drive if=scsi,driver=file,filename=debian-buster-armel-rootfs.ext4 \
  -kernel zImage \
  -dtb versatile-pb.dtb \
  -append "console=ttyAMA0 ro root=/dev/sda"
```

Arm Versatile Express boards (`vexpress-a9`, `vexpress-a15`)

QEMU models two variants of the Arm Versatile Express development board family:

- `vexpress-a9` models the combination of the Versatile Express motherboard and the CoreTile Express A9x4 daughterboard
- `vexpress-a15` models the combination of the Versatile Express motherboard and the CoreTile Express A15x2 daughterboard

Note that as this hardware does not have PCI, IDE or SCSI, the only available storage option is emulated SD card.

Implemented devices:

- PL041 audio
- PL181 SD controller
- PL050 keyboard and mouse
- PL011 UARTs
- SP804 timers

- I2C controller
- PL031 RTC
- PL111 LCD display controller
- Flash memory
- LAN9118 ethernet

Unimplemented devices:

- SP810 system control block
- PCI-express
- USB controller (Philips ISP1761)
- Local DAP ROM
- CoreSight interfaces
- PL301 AXI interconnect
- SCC
- System counter
- HDLCD controller (vexpress-a15)
- SP805 watchdog
- PL341 dynamic memory controller
- DMA330 DMA controller
- PL354 static memory controller
- BP147 TrustZone Protection Controller
- TrustZone Address Space Controller

Other differences between the hardware and the QEMU model:

- QEMU will default to creating one CPU unless you pass a different `-smp` argument
- QEMU allows the amount of RAM provided to be specified with the `-m` argument
- QEMU defaults to providing a CPU which does not provide either TrustZone or the Virtualization Extensions: if you want these you must enable them with `-machine secure=on` and `-machine virtualization=on`
- QEMU provides 4 virtio-mmio virtio transports; these start at address 0x10013000 for vexpress-a9 and at 0x1c130000 for vexpress-a15, and have IRQs from 40 upwards. If a dtb is provided on the command line then QEMU will edit it to include suitable entries describing these transports for the guest.

Booting a Linux kernel

Building a current Linux kernel with `multi_v7_defconfig` should be enough to get something running. Nowadays an out-of-tree build is recommended (and also useful if you build a lot of different targets). In the following example `$BLD` points to the build directory and `$SRC` points to the root of the Linux source tree. You can drop `$SRC` if you are running from there.

```
$ make O=$BLD -C $SRC ARCH=arm CROSS_COMPILE=arm-linux-gnueabihf- multi_v7_defconfig
$ make O=$BLD -C $SRC ARCH=arm CROSS_COMPILE=arm-linux-gnueabihf-
```

By default you will want to boot your rootfs off the sdcard interface. Your rootfs will need to be padded to the right size. With a suitable DTB you could also add devices to the virtio-mmio bus.

```
$ qemu-system-arm -cpu cortex-a15 -smp 4 -m 4096 \  
-machine type=vexpress-a15 -serial mon:stdio \  
-drive if=sd,driver=file,filename=armel-rootfs.ext4 \  
-kernel zImage \  
-dtb vexpress-v2p-ca15-tc1.dtb \  
-append "console=ttyAMA0 root=/dev/mmcblk0 ro"
```

Aspeed family boards (*-bmc, ast2500-evb, ast2600-evb)

The QEMU Aspeed machines model BMCs of various OpenPOWER systems and Aspeed evaluation boards. They are based on different releases of the Aspeed SoC : the AST2400 integrating an ARM926EJ-S CPU (400MHz), the AST2500 with an ARM1176JZS CPU (800MHz) and more recently the AST2600 with dual cores ARM Cortex A7 CPUs (1.2GHz).

The SoC comes with RAM, Gigabit ethernet, USB, SD/MMC, USB, SPI, I2C, etc.

AST2400 SoC based machines :

- palmetto-bmc OpenPOWER Palmetto POWER8 BMC

AST2500 SoC based machines :

- ast2500-evb Aspeed AST2500 Evaluation board
- romulus-bmc OpenPOWER Romulus POWER9 BMC
- witherspoon-bmc OpenPOWER Witherspoon POWER9 BMC
- sonorapass-bmc OCP SonoraPass BMC
- swift-bmc OpenPOWER Swift BMC POWER9

AST2600 SoC based machines :

- ast2600-evb Aspeed AST2600 Evaluation board (Cortex A7)
- tacoma-bmc OpenPOWER Witherspoon POWER9 AST2600 BMC

Supported devices

- SMP (for the AST2600 Cortex-A7)
- Interrupt Controller (VIC)
- Timer Controller
- RTC Controller
- I2C Controller
- System Control Unit (SCU)
- SRAM mapping
- X-DMA Controller (basic interface)
- Static Memory Controller (SMC or FMC) - Only SPI Flash support
- SPI Memory Controller

- USB 2.0 Controller
- SD/MMC storage controllers
- SDRAM controller (dummy interface for basic settings and training)
- Watchdog Controller
- GPIO Controller (Master only)
- UART
- Ethernet controllers
- Front LEDs (PCA9552 on I2C bus)

Missing devices

- Coprocessor support
- ADC (out of tree implementation)
- PWM and Fan Controller
- LPC Bus Controller
- Slave GPIO Controller
- Super I/O Controller
- Hash/Crypto Engine
- PCI-Express 1 Controller
- Graphic Display Controller
- PECE Controller
- MCTP Controller
- Mailbox Controller
- Virtual UART
- eSPI Controller
- I3C Controller

Boot options

The Aspeed machines can be started using the `-kernel` option to load a Linux kernel or from a firmware image which can be downloaded from the OpenPOWER jenkins :

<https://openpower.xyz/>

The image should be attached as an MTD drive. Run :

```
$ qemu-system-arm -M romulus-bmc -nic user \
    -drive file=flash-romulus,format=raw,if=mtd -nographic
```

Options specific to Aspeed machines are :

- `execute-in-place` which emulates the boot from the CE0 flash device by using the FMC controller to load the instructions, and not simply from RAM. This takes a little longer.

- `fmc-model` to change the FMC Flash model. FW needs support for the chip model to boot.
- `spi-model` to change the SPI Flash model.

For instance, to start the `ast2500-evb` machine with a different FMC chip and a bigger (64M) SPI chip, use :

```
-M ast2500-evb, fmc-model=mx25l25635e, spi-model=mx66u51235f
```

Boundary Devices SABRE Lite (`sabrelite`)

Boundary Devices SABRE Lite i.MX6 Development Board is a low-cost development platform featuring the powerful Freescale / NXP Semiconductor's i.MX 6 Quad Applications Processor.

Supported devices

The SABRE Lite machine supports the following devices:

- Up to 4 Cortex A9 cores
- Generic Interrupt Controller
- 1 Clock Controller Module
- 1 System Reset Controller
- 5 UARTs
- 2 EPIC timers
- 1 GPT timer
- 2 Watchdog timers
- 1 FEC Ethernet controller
- 3 I2C controllers
- 7 GPIO controllers
- 4 SDHC storage controllers
- 4 USB 2.0 host controllers
- 5 ECSPI controllers
- 1 SST 25VF016B flash

Please note above list is a complete superset the QEMU SABRE Lite machine can support. For a normal use case, a device tree blob that represents a real world SABRE Lite board, only exposes a subset of devices to the guest software.

Boot options

The SABRE Lite machine can start using the standard -kernel functionality for loading a Linux kernel, U-Boot boot-loader or ELF executable.

Running Linux kernel

Linux mainline v5.10 release is tested at the time of writing. To build a Linux mainline kernel that can be booted by the SABRE Lite machine, simply configure the kernel using the `imx_v6_v7_defconfig` configuration:

```
$ export ARCH=arm
$ export CROSS_COMPILE=arm-linux-gnueabihf-
$ make imx_v6_v7_defconfig
$ make
```

To boot the newly built Linux kernel in QEMU with the SABRE Lite machine, use:

```
$ qemu-system-arm -M sabrelite -smp 4 -m 1G \
  -display none -serial null -serial stdio \
  -kernel arch/arm/boot/zImage \
  -dtb arch/arm/boot/dts/imx6q-sabrelite.dtb \
  -initrd /path/to/rootfs.ext4 \
  -append "root=/dev/ram"
```

Running U-Boot

U-Boot mainline v2020.10 release is tested at the time of writing. To build a U-Boot mainline bootloader that can be booted by the SABRE Lite machine, use the `mx6qsabrelite_defconfig` with similar commands as described above for Linux:

```
$ export CROSS_COMPILE=arm-linux-gnueabihf-
$ make mx6qsabrelite_defconfig
```

Note we need to adjust settings by:

```
$ make menuconfig
```

then manually select the following configuration in U-Boot:

Device Tree Control > Provider of DTB for DT Control > Embedded DTB

To start U-Boot using the SABRE Lite machine, provide the u-boot binary to the `-kernel` argument, along with an SD card image with rootfs:

```
$ qemu-system-arm -M sabrelite -smp 4 -m 1G \
  -display none -serial null -serial stdio \
  -kernel u-boot
```

The following example shows booting Linux kernel from dhcp, and uses the rootfs on an SD card. This requires some additional command line parameters for QEMU:

```
-nic user,tftp=/path/to/kernel/zImage \
-drive file=sdcard.img,id=rootfs -device sd-card,drive=rootfs
```

The directory for the built-in TFTP server should also contain the device tree blob of the SABRE Lite board. The sample SD card image was populated with the root file system with one single partition. You may adjust the kernel “root=” boot parameter accordingly.

After U-Boot boots, type the following commands in the U-Boot command shell to boot the Linux kernel:

```
=> setenv ethaddr 00:11:22:33:44:55
=> setenv bootfile zImage
=> dhcp
=> tftpboot 14000000 imx6q-sabrelite.dtb
=> setenv bootargs root=/dev/mmcblk3p1
=> bootz 12000000 - 14000000
```

Canon A1100 (canon-a1100)

This machine is a model of the Canon PowerShot A1100 camera, which uses the DIGIC SoC. This model is based on reverse engineering efforts by the contributors to the [CHDK](#) and [Magic Lantern](#) projects.

The emulation is incomplete. In particular it can't be used to run the original camera firmware, but it can successfully run an experimental version of the [barebox bootloader](#).

Freecom MusicPal (musicpal)

The Freecom MusicPal internet radio emulation includes the following elements:

- Marvell MV88W8618 Arm core.
- 32 MB RAM, 256 KB SRAM, 8 MB flash.
- Up to 2 16550 UARTs
- MV88W8xx8 Ethernet controller
- MV88W8618 audio controller, WM8750 CODEC and mixer
- 128x64 display with brightness control
- 2 buttons, 2 navigation wheels with button function

Gumstix Connex and Verdex (connex, verdex)

These machines model the Gumstix Connex and Verdex boards. The Connex has a PXA255 CPU and the Verdex has a PXA270.

Implemented devices:

- NOR flash
- SMC91C111 ethernet
- Interrupt controller
- DMA
- Timer
- GPIO
- MMC/SD card
- Fast infra-red communications port (FIR)
- LCD controller
- Synchronous serial ports (SPI)

- PCMCIA interface
- I2C
- I2S

Nokia N800 and N810 tablets (`n800`, `n810`)

Nokia N800 and N810 internet tablets (known also as RX-34 and RX-44 / 48) emulation supports the following elements:

- Texas Instruments OMAP2420 System-on-chip (ARM1136 core)
- RAM and non-volatile OneNAND Flash memories
- Display connected to EPSON remote framebuffer chip and OMAP on-chip display controller and a LS041y3 MIPI DBI-C controller
- TI TSC2301 (in N800) and TI TSC2005 (in N810) touchscreen controllers driven through SPI bus
- National Semiconductor LM8323-controlled qwerty keyboard driven through I²C bus
- Secure Digital card connected to OMAP MMC/SD host
- Three OMAP on-chip UARTs and on-chip STI debugging console
- Mentor Graphics “Inventra” dual-role USB controller embedded in a TI TUSB6010 chip - only USB host mode is supported
- TI TMP105 temperature sensor driven through I²C bus
- TI TWL92230C power management companion with an RTC on I²C bus
- Nokia RETU and TAHVO multi-purpose chips with an RTC, connected through CBUS

Nuvoton iBMC boards (`npcm750-evb`, `quanta-gsj`)

The [Nuvoton iBMC](#) chips (NPCM7xx) are a family of ARM-based SoCs that are designed to be used as Baseboard Management Controllers (BMCs) in various servers. They all feature one or two ARM Cortex A9 CPU cores, as well as an assortment of peripherals targeted for either Enterprise or Data Center / Hyperscale applications. The former is a superset of the latter, so NPCM750 has all the peripherals of NPCM730 and more.

The NPCM750 SoC has two Cortex A9 cores and is targeted for the Enterprise segment. The following machines are based on this chip :

- `npcm750-evb` Nuvoton NPCM750 Evaluation board

The NPCM730 SoC has two Cortex A9 cores and is targeted for Data Center and Hyperscale applications. The following machines are based on this chip :

- `quanta-gsj` Quanta GSJ server BMC

There are also two more SoCs, NPCM710 and NPCM705, which are single-core variants of NPCM750 and NPCM730, respectively. These are currently not supported by QEMU.

Supported devices

- SMP (Dual Core Cortex-A9)
- Cortex-A9MPCore built-in peripherals: SCU, GIC, Global Timer, Private Timer and Watchdog.

- SRAM, ROM and DRAM mappings
- System Global Control Registers (GCR)
- Clock and reset controller (CLK)
- Timer controller (TIM)
- Serial ports (16550-based)
- DDR4 memory controller (dummy interface indicating memory training is done)
- OTP controllers (no protection features)
- Flash Interface Unit (FIU; no protection features)
- Random Number Generator (RNG)
- USB host (USBH)
- GPIO controller
- Analog to Digital Converter (ADC)
- Pulse Width Modulation (PWM)
- SMBus controller (SMBF)

Missing devices

- LPC/eSPI host-to-BMC interface, including
 - Keyboard and mouse controller interface (KBCI)
 - Keyboard Controller Style (KCS) channels
 - BIOS POST code FIFO
 - System Wake-up Control (SWC)
 - Shared memory (SHM)
 - eSPI slave interface
- Ethernet controllers (GMAC and EMC)
- USB device (USBD)
- Peripheral SPI controller (PSPI)
- SD/MMC host
- PECE interface
- Tachometer
- PCI and PCIe root complex and bridges
- VDM and MCTP support
- Serial I/O expansion
- LPC/eSPI host
- Coprocessor
- Graphics
- Video capture

- Encoding compression engine
- Security features

Boot options

The Nuvoton machines can boot from an OpenBMC firmware image, or directly into a kernel using the `-kernel` option. OpenBMC images for *quanta-gsj* and possibly others can be downloaded from the OpenPOWER jenkins :

<https://openpower.xyz/>

The firmware image should be attached as an MTD drive. Example :

```
$ qemu-system-arm -machine quanta-gsj -nographic \
  -drive file=image-bmc,if=mtd,bus=0,unit=0,format=raw
```

The default root password for test images is usually `OpenBmc`.

Orange Pi PC (*orange-pi-pc*)

The Xunlong Orange Pi PC is an Allwinner H3 System on Chip based embedded computer with mainline support in both U-Boot and Linux. The board comes with a Quad Core Cortex-A7 @ 1.3GHz, 1GiB RAM, 100Mbit ethernet, USB, SD/MMC, USB, HDMI and various other I/O.

Supported devices

The Orange Pi PC machine supports the following devices:

- SMP (Quad Core Cortex-A7)
- Generic Interrupt Controller configuration
- SRAM mappings
- SDRAM controller
- Real Time Clock
- Timer device (re-used from Allwinner A10)
- UART
- SD/MMC storage controller
- EMAC ethernet
- USB 2.0 interfaces
- Clock Control Unit
- System Control module
- Security Identifier device

Limitations

Currently, Orange Pi PC does *not* support the following features:

- Graphical output via HDMI, GPU and/or the Display Engine
- Audio output
- Hardware Watchdog

Also see the ‘unimplemented’ array in the Allwinner H3 SoC module for a complete list of unimplemented I/O devices:
`./hw/arm/allwinner-h3.c`

Boot options

The Orange Pi PC machine can start using the standard `-kernel` functionality for loading a Linux kernel or ELF executable. Additionally, the Orange Pi PC machine can also emulate the BootROM which is present on an actual Allwinner H3 based SoC, which loads the bootloader from a SD card, specified via the `-sd` argument to `qemu-system-arm`.

Machine-specific options

The following machine-specific options are supported:

- `allwinner-rtc.base-year=YYYY`

The Allwinner RTC device is automatically created by the Orange Pi PC machine and uses a default base year value which can be overridden using the ‘base-year’ property. The base year is the actual represented year when the RTC year value is zero. This option can be used in case the target operating system driver uses a different base year value. The minimum value for the base year is 1900.

- `allwinner-sid.identifier=abcd1122-a000-b000-c000-12345678ffff`

The Security Identifier value can be read by the guest. For example, U-Boot uses it to determine a unique MAC address.

The above machine-specific options can be specified in `qemu-system-arm` via the ‘-global’ argument, for example:

```
$ qemu-system-arm -M orangepi-pc -sd mycard.img \
    -global allwinner-rtc.base-year=2000
```

Running mainline Linux

Mainline Linux kernels from 4.19 up to latest master are known to work. To build a Linux mainline kernel that can be booted by the Orange Pi PC machine, simply configure the kernel using the `sunxi_defconfig` configuration:

```
$ ARCH=arm CROSS_COMPILE=arm-linux-gnueabi- make mrproper
$ ARCH=arm CROSS_COMPILE=arm-linux-gnueabi- make sunxi_defconfig
```

To be able to use USB storage, you need to manually enable the corresponding configuration item. Start the `kconfig` configuration tool:

```
$ ARCH=arm CROSS_COMPILE=arm-linux-gnueabi- make menuconfig
```

Navigate to the following item, enable it and save your configuration:

Device Drivers > USB support > USB Mass Storage support

Build the Linux kernel with:

```
$ ARCH=arm CROSS_COMPILE=arm-linux-gnueabi- make
```

To boot the newly build linux kernel in QEMU with the Orange Pi PC machine, use:

```
$ qemu-system-arm -M orangepi-pc -nic user -nographic \
  -kernel /path/to/linux/arch/arm/boot/zImage \
  -append 'console=ttyS0,115200' \
  -dtb /path/to/linux/arch/arm/boot/dts/sun8i-h3-orangepi-pc.dtb
```

Orange Pi PC images

Note that the mainline kernel does not have a root filesystem. You may provide it with an official Orange Pi PC image from the official website:

<http://www.orangepi.org/downloadresources/>

Another possibility is to run an Armbian image for Orange Pi PC which can be downloaded from:

<https://www.armbian.com/orange-pi-pc/>

Alternatively, you can also choose to build you own image with buildroot using the orangepi_pc_defconfig. Also see <https://buildroot.org> for more information.

When using an image as an SD card, it must be resized to a power of two. This can be done with the qemu-img command. It is recommended to only increase the image size instead of shrinking it to a power of two, to avoid loss of data. For example, to prepare a downloaded Armbian image, first extract it and then increase its size to one gigabyte as follows:

```
$ qemu-img resize Armbian_19.11.3_Orangepipc_bionic_current_5.3.9.img 1G
```

You can choose to attach the selected image either as an SD card or as USB mass storage. For example, to boot using the Orange Pi PC Debian image on SD card, simply add the -sd argument and provide the proper root= kernel parameter:

```
$ qemu-system-arm -M orangepi-pc -nic user -nographic \
  -kernel /path/to/linux/arch/arm/boot/zImage \
  -append 'console=ttyS0,115200 root=/dev/mmcblk0p2' \
  -dtb /path/to/linux/arch/arm/boot/dts/sun8i-h3-orangepi-pc.dtb \
  -sd OrangePi_pc_debian_stretch_server_linux5.3.5_v1.0.img
```

To attach the image as an USB mass storage device to the machine, simply append to the command:

```
-drive if=none,id=stick,file=myimage.img \
-device usb-storage,bus=usb-bus.0,drive=stick
```

Instead of providing a custom Linux kernel via the -kernel command you may also choose to let the Orange Pi PC machine load the bootloader from SD card, just like a real board would do using the BootROM. Simply pass the selected image via the -sd argument and remove the -kernel, -append, -dbt and -initrd arguments:

```
$ qemu-system-arm -M orangepi-pc -nic user -nographic \
  -sd Armbian_19.11.3_Orangepipc_buster_current_5.3.9.img
```

Note that both the official Orange Pi PC images and Armbian images start a lot of userland programs via systemd. Depending on the host hardware and OS, they may be slow to emulate, especially due to emulating the 4 cores. To

help reduce the performance slow down due to emulating the 4 cores, you can give the following kernel parameters via U-Boot (or via -append):

```
=> setenv extraargs 'systemd.default_timeout_start_sec=9000 loglevel=7 nosmp_  
↪console=ttyS0,115200'
```

Running U-Boot

U-Boot mainline can be build and configured using the `orange_pi_pc_defconfig` using similar commands as describe above for Linux. Note that it is recommended for development/testing to select the following configuration setting in U-Boot:

Device Tree Control > Provider for DTB for DT Control > Embedded DTB

To start U-Boot using the Orange Pi PC machine, provide the u-boot binary to the -kernel argument:

```
$ qemu-system-arm -M orange_pi_pc -nic user -nographic \  
-kernel /path/to/u-boot/u-boot -sd disk.img
```

Use the following U-boot commands to load and boot a Linux kernel from SD card:

```
=> setenv bootargs console=ttyS0,115200  
=> ext2load mmc 0 0x42000000 zImage  
=> ext2load mmc 0 0x43000000 sun8i-h3-orange_pi_pc.dtb  
=> bootz 0x42000000 - 0x43000000
```

Running NetBSD

The NetBSD operating system also includes support for Allwinner H3 based boards, including the Orange Pi PC. NetBSD 9.0 is known to work best for the Orange Pi PC board and provides a fully working system with serial console, networking and storage. For the Orange Pi PC machine, get the ‘evbarm-earmv7hf’ based image from:

<https://cdn.netbsd.org/pub/NetBSD/NetBSD-9.0/evbarm-earmv7hf/binary/gzimg/armv7.img.gz>

The image requires manually installing U-Boot in the image. Build U-Boot with the `orange_pi_pc_defconfig` configuration as described in the previous section. Next, unzip the NetBSD image and write the U-Boot binary including SPL using:

```
$ gunzip armv7.img.gz  
$ dd if=/path/to/u-boot-sunxi-with-spl.bin of=armv7.img bs=1024 seek=8 conv=notrunc
```

Finally, before starting the machine the SD image must be extended such that the size of the SD image is a power of two and that the NetBSD kernel will not conclude the NetBSD partition is larger than the emulated SD card:

```
$ qemu-img resize armv7.img 2G
```

Start the machine using the following command:

```
$ qemu-system-arm -M orange_pi_pc -nic user -nographic \  
-sd armv7.img -global allwinner-rtc.base-year=2000
```

At the U-Boot stage, interrupt the automatic boot process by pressing a key and set the following environment variables before booting:


```
=> setenv bootargs root=ld0a
=> setenv kernel netbsd-GENERIC.ub
=> setenv fdtfile dtb/sun8i-h3-orangepi-pc.dtb
=> setenv bootcmd 'fatload mmc 0:1 ${kernel_addr_r} ${kernel}; fatload mmc 0:1 ${fdt_
↪addr_r} ${fdtfile}; fdt addr ${fdt_addr_r}; bootm ${kernel_addr_r} - ${fdt_addr_r}'
```

Optionally you may save the environment variables to SD card with ‘saveenv’. To continue booting simply give the ‘boot’ command and NetBSD boots.

Orange Pi PC acceptance tests

The Orange Pi PC machine has several acceptance tests included. To run the whole set of tests, build QEMU from source and simply provide the following command:

```
$ AVOCADO_ALLOW_LARGE_STORAGE=yes avocado --show=app,console run \
  -t machine:orangepi-pc tests/acceptance/boot_linux_console.py
```

Palm Tungsten|E PDA (cheetah)

The Palm Tungsten|E PDA (codename “Cheetah”) emulation includes the following elements:

- Texas Instruments OMAP310 System-on-chip (ARM925T core)
- ROM and RAM memories (ROM firmware image can be loaded with -option-rom)
- On-chip LCD controller
- On-chip Real Time Clock
- TI TSC2102i touchscreen controller / analog-digital converter / Audio CODEC, connected through MicroWire and I²S busses
- GPIO-connected matrix keypad
- Secure Digital card connected to OMAP MMC/SD host
- Three on-chip UARTs

Raspberry Pi boards (raspi0, raspilap, raspi2b, raspi3ap, raspi3b)

QEMU provides models of the following Raspberry Pi boards:

raspi0 and raspilap ARM1176JZF-S core, 512 MiB of RAM

raspi2b Cortex-A7 (4 cores), 1 GiB of RAM

raspi3ap Cortex-A53 (4 cores), 512 MiB of RAM

raspi3b Cortex-A53 (4 cores), 1 GiB of RAM

Implemented devices

- ARM1176JZF-S, Cortex-A7 or Cortex-A53 CPU
- Interrupt controller
- DMA controller

- Clock and reset controller (CPRMAN)
- System Timer
- GPIO controller
- Serial ports (BCM2835 AUX - 16550 based - and PL011)
- Random Number Generator (RNG)
- Frame Buffer
- USB host (USBH)
- GPIO controller
- SD/MMC host controller
- SoC thermal sensor
- USB2 host controller (DWC2 and MPHI)
- MailBox controller (MBOX)
- VideoCore firmware (property)

Missing devices

- Peripheral SPI controller (SPI)
- Analog to Digital Converter (ADC)
- Pulse Width Modulation (PWM)

Sharp XScale-based PDA models (*akita*, *borzoi*, *spitz*, *terrier*, *tosa*)

The Sharp Zaurus are PDAs based on XScale, able to run Linux ('SL series').

The SL-6000 ("Tosa"), released in 2005, uses a PXA255 System-on-chip.

The SL-C3000 ("Spitz"), SL-C1000 ("Akita"), SL-C3100 ("Borzoi") and SL-C3200 ("Terrier") use a PXA270.

The clamshell PDA models emulation includes the following peripherals:

- Intel PXA255/PXA270 System-on-chip (ARMv5TE core)
- NAND Flash memory - not in "Tosa"
- IBM/Hitachi DSCM microdrive in a PXA PCMCIA slot - not in "Akita"
- On-chip OHCI USB controller - not in "Tosa"
- On-chip LCD controller
- On-chip Real Time Clock
- TI ADS7846 touchscreen controller on SSP bus
- Maxim MAX1111 analog-digital converter on I²C bus
- GPIO-connected keyboard controller and LEDs
- Secure Digital card connected to PXA MMC/SD host
- Three on-chip UARTs
- WM8750 audio CODEC on I²C and I²S busses

Sharp Zaurus SL-5500 (`collie`)

This machine is a model of the Sharp Zaurus SL-5500, which was a 1990s PDA based on the StrongARM SA1110.

Implemented devices:

- NOR flash
- Interrupt controller
- Timer
- RTC
- GPIO
- Peripheral Pin Controller (PPC)
- UARTs
- Synchronous Serial Ports (SSP)

Siemens SX1 (`sx1`, `sx1-v1`)

The Siemens SX1 models v1 and v2 (default) basic emulation. The emulation includes the following elements:

- Texas Instruments OMAP310 System-on-chip (ARM925T core)
- ROM and RAM memories (ROM firmware image can be loaded with `-pflash`) V1 1 Flash of 16MB and 1 Flash of 8MB V2 1 Flash of 32MB
- On-chip LCD controller
- On-chip Real Time Clock
- Secure Digital card connected to OMAP MMC/SD host
- Three on-chip UARTs

Stellaris boards (`lm3s6965evb`, `lm3s811evb`)

The Luminary Micro Stellaris LM3S811EVB emulation includes the following devices:

- Cortex-M3 CPU core.
- 64k Flash and 8k SRAM.
- Timers, UARTs, ADC and I²C interface.
- OSRAM Pictiva 96x16 OLED with SSD0303 controller on I²C bus.

The Luminary Micro Stellaris LM3S6965EVB emulation includes the following devices:

- Cortex-M3 CPU core.
- 256k Flash and 64k SRAM.
- Timers, UARTs, ADC, I²C and SSI interfaces.
- OSRAM Pictiva 128x64 OLED with SSD0323 controller connected via SSI.

‘virt’ generic virtual platform (`virt`)

The *virt* board is a platform which does not correspond to any real hardware; it is designed for use in virtual machines. It is the recommended board type if you simply want to run a guest such as Linux and do not care about reproducing the idiosyncrasies and limitations of a particular bit of real-world hardware.

This is a “versioned” board model, so as well as the `virt` machine type itself (which may have improvements, bugfixes and other minor changes between QEMU versions) a version is provided that guarantees to have the same behaviour as that of previous QEMU releases, so that VM migration will work between QEMU versions. For instance the `virt-5.0` machine type will behave like the `virt` machine from the QEMU 5.0 release, and migration should work between `virt-5.0` of the 5.0 release and `virt-5.0` of the 5.1 release. Migration is not guaranteed to work between different QEMU releases for the non-versioned `virt` machine type.

Supported devices

The *virt* board supports:

- PCI/PCIe devices
- Flash memory
- One PL011 UART
- An RTC
- The `fw_cfg` device that allows a guest to obtain data from QEMU
- A PL061 GPIO controller
- An optional SMMUv3 IOMMU
- hotpluggable DIMMs
- hotpluggable NVDIMMs
- An MSI controller (GICv2M or ITS). GICv2M is selected by default along with GICv2. ITS is selected by default with GICv3 (\geq `virt-2.7`). Note that ITS is not modeled in TCG mode.
- 32 virtio-mmio transport devices
- running guests using the KVM accelerator on aarch64 hardware
- large amounts of RAM (at least 255GB, and more if using `highmem`)
- many CPUs (up to 512 if using a GICv3 and `highmem`)
- Secure-World-only devices if the CPU has TrustZone:
 - A second PL011 UART
 - A second PL061 GPIO controller, with GPIO lines for triggering a system reset or system poweroff
 - A secure flash memory
 - 16MB of secure RAM

Supported guest CPU types:

- `cortex-a7` (32-bit)
- `cortex-a15` (32-bit; the default)
- `cortex-a53` (64-bit)
- `cortex-a57` (64-bit)

- `cortex-a72` (64-bit)
- `host` (with KVM only)
- `max` (same as `host` for KVM; best possible emulation with TCG)

Note that the default is `cortex-a15`, so for an AArch64 guest you must specify a CPU type.

Graphics output is available, but unlike the x86 PC machine types there is no default display device enabled: you should select one from the Display devices section of “-device help”. The recommended option is `virtio-gpu-pci`; this is the only one which will work correctly with KVM. You may also need to ensure your guest kernel is configured with support for this; see below.

Machine-specific options

The following machine-specific options are supported:

secure Set `on/off` to enable/disable emulating a guest CPU which implements the Arm Security Extensions (TrustZone). The default is `off`.

virtualization Set `on/off` to enable/disable emulating a guest CPU which implements the Arm Virtualization Extensions. The default is `off`.

mte Set `on/off` to enable/disable emulating a guest CPU which implements the Arm Memory Tagging Extensions. The default is `off`.

highmem Set `on/off` to enable/disable placing devices and RAM in physical address space above 32 bits. The default is `on` for machine types later than `virt-2.12`.

gic-version Specify the version of the Generic Interrupt Controller (GIC) to provide. Valid values are:

2 GICv2

3 GICv3

host Use the same GIC version the host provides, when using KVM

max Use the best GIC version possible (same as host when using KVM; currently same as 3` for TCG, but this may change in future)

its Set `on/off` to enable/disable ITS instantiation. The default is `on` for machine types later than `virt-2.7`.

iommu Set the IOMMU type to create for the guest. Valid values are:

none Don't create an IOMMU (the default)

smmu-v3 Create an SMMUv3

ras Set `on/off` to enable/disable reporting host memory errors to a guest using ACPI and guest external abort exceptions. The default is `off`.

Linux guest kernel configuration

The ‘defconfig’ for Linux arm and arm64 kernels should include the right device drivers for virtio and the PCI controller; however some older kernel versions, especially for 32-bit Arm, did not have everything enabled by default. If you're not seeing PCI devices that you expect, then check that your guest config has:

```
CONFIG_PCI=y
CONFIG_VIRTIO_PCI=y
CONFIG_PCI_HOST_GENERIC=y
```

If you want to use the `virtio-gpu-pci` graphics device you will also need:

```
CONFIG_DRM=y
CONFIG_DRM_VIRTIO_GPU=y
```

Hardware configuration information for bare-metal programming

The `virt` board automatically generates a device tree blob (“dtb”) which it passes to the guest. This provides information about the addresses, interrupt lines and other configuration of the various devices in the system. Guest code can rely on and hard-code the following addresses:

- Flash memory starts at address `0x0000_0000`
- RAM starts at `0x4000_0000`

All other information about device locations may change between QEMU versions, so guest code must look in the DTB.

QEMU supports two types of guest image boot for `virt`, and the way for the guest code to locate the dtb binary differs:

- For guests using the Linux kernel boot protocol (this means any non-ELF file passed to the QEMU `-kernel` option) the address of the DTB is passed in a register (`⌢2` for 32-bit guests, or `⌢0` for 64-bit guests)
- For guests booting as “bare-metal” (any other kind of boot), the DTB is at the start of RAM (`0x4000_0000`)

Xilinx Versal Virt (`xlnx-versal-virt`)

Xilinx Versal is a family of heterogeneous multi-core SoCs (System on Chip) that combine traditional hardened CPUs and I/O peripherals in a Processing System (PS) with runtime programmable FPGA logic (PL) and an Artificial Intelligence Engine (AIE).

More details here: <https://www.xilinx.com/products/silicon-devices/acap/versal.html>

The family of Versal SoCs share a single architecture but come in different parts with different speed grades, amounts of PL and other differences.

The Xilinx Versal Virt board in QEMU is a model of a virtual board (does not exist in reality) with a virtual Versal SoC without I/O limitations. Currently, we support the following cores and devices:

Implemented CPU cores:

- 2 ACPUs (ARM Cortex-A72)

Implemented devices:

- Interrupt controller (ARM GICv3)
- 2 UARTs (ARM PL011)
- An RTC (Versal built-in)
- 2 GEMs (Cadence MACB Ethernet MACs)
- 8 ADMA (Xilinx zDMA) channels
- 2 SD Controllers
- OCM (256KB of On Chip Memory)
- DDR memory

QEMU does not yet model any other devices, including the PL and the AI Engine.

Other differences between the hardware and the QEMU model:

- QEMU allows the amount of DDR memory provided to be specified with the `-m` argument. If a DTB is provided on the command line then QEMU will edit it to include suitable entries describing the Versal DDR memory ranges.
- QEMU provides 8 virtio-mmio virtio transports; these start at address `0xa0000000` and have IRQs from 111 and upwards.

Running

If the user provides an Operating System to be loaded, we expect users to use the `-kernel` command line option.

Users can load firmware or boot-loaders with the `-device loader` options.

When loading an OS, QEMU generates a DTB and selects an appropriate address where it gets loaded. This DTB will be passed to the kernel in register `x0`.

If there's no `-kernel` option, we generate a DTB and place it at `0x1000` for boot-loaders or firmware to pick it up.

If users want to provide their own DTB, they can use the `-dtb` option. These DTBs will have their memory nodes modified to match QEMU's selected `ram_size` option before they get passed to the kernel or FW.

When loading an OS, we turn on QEMU's PSCI implementation with SMC as the PSCI conduit. When there's no `-kernel` option, we assume the user provides EL3 firmware to handle PSCI.

A few examples:

Direct Linux boot of a generic ARM64 upstream Linux kernel:

```
$ qemu-system-aarch64 -M xlnx-versal-virt -m 2G \
  -serial mon:stdio -display none \
  -kernel arch/arm64/boot/Image \
  -nic user -nic user \
  -device virtio-rng-device,bus=virtio-mmio-bus.0 \
  -drive if=none,index=0,file=hd0.qcow2,id=hd0,snapshot \
  -drive file=qemu_sd.qcow2,if=sd,index=0,snapshot \
  -device virtio-blk-device,drive=hd0 -append root=/dev/vda
```

Direct Linux boot of PetaLinux 2019.2:

```
$ qemu-system-aarch64 -M xlnx-versal-virt -m 2G \
  -serial mon:stdio -display none \
  -kernel petalinux-v2019.2/Image \
  -append "rdinit=/sbin/init console=ttyAMA0,115200n8 earlycon=pl011,mmio,
  ↪0xFF000000,115200n8" \
  -net nic,model=cadence_gem,netdev=net0 -netdev user,id=net0 \
  -device virtio-rng-device,bus=virtio-mmio-bus.0,rng=rng0 \
  -object rng-random,filename=/dev/urandom,id=rng0
```

Boot PetaLinux 2019.2 via ARM Trusted Firmware (2018.3 because the 2019.2 version of ATF tries to configure the CCI which we don't model) and U-boot:

```
$ qemu-system-aarch64 -M xlnx-versal-virt -m 2G \
  -serial stdio -display none \
  -device loader,file=petalinux-v2018.3/bl31.elf,cpu-num=0 \
  -device loader,file=petalinux-v2019.2/u-boot.elf \
```

(continues on next page)

(continued from previous page)

```
-device loader,addr=0x20000000,file=petalinux-v2019.2/Image \
-nic user -nic user \
-device virtio-rng-device,bus=virtio-mmio-bus.0,rng=rng0 \
-object rng-random,filename=/dev/urandom,id=rng0
```

Run the following at the U-Boot prompt:

```
Versal>
fdt addr $fdtcontroladdr
fdt move $fdtcontroladdr 0x40000000
fdt set /timer clock-frequency <0x3dfd240>
setenv bootargs "rdinit=/sbin/init maxcpus=1 console=ttyAMA0,115200n8 earlycon=pl011,
↪mmio,0xFF000000,115200n8"
booti 20000000 - 40000000
fdt addr $fdtcontroladdr
```

Boot Linux as DOM0 on Xen via U-Boot:

```
$ qemu-system-aarch64 -M xlnx-versal-virt -m 4G \
  -serial stdio -display none \
  -device loader,file=petalinux-v2019.2/u-boot.elf,cpu-num=0 \
  -device loader,addr=0x30000000,file=linux/2018-04-24/xen \
  -device loader,addr=0x40000000,file=petalinux-v2019.2/Image \
  -nic user -nic user \
  -device virtio-rng-device,bus=virtio-mmio-bus.0,rng=rng0 \
  -object rng-random,filename=/dev/urandom,id=rng0
```

Run the following at the U-Boot prompt:

```
Versal>
fdt addr $fdtcontroladdr
fdt move $fdtcontroladdr 0x20000000
fdt set /timer clock-frequency <0x3dfd240>
fdt set /chosen xen,xen-bootargs "console=dtuart dtuart=/uart@ff000000 dom0_mem=640M_
↪bootscrub=0 maxcpus=1 timer_slop=0"
fdt set /chosen xen,dm0-bootargs "rdinit=/sbin/init clk_ignore_unused console=hvc0_
↪maxcpus=1"
fdt mknode /chosen dom0
fdt set /chosen/dom0 compatible "xen,multiboot-module"
fdt set /chosen/dom0 reg <0x00000000 0x40000000 0x0 0x03100000>
booti 30000000 - 20000000
```

Boot Linux as Dom0 on Xen via ARM Trusted Firmware and U-Boot:

```
$ qemu-system-aarch64 -M xlnx-versal-virt -m 4G \
  -serial stdio -display none \
  -device loader,file=petalinux-v2018.3/bl31.elf,cpu-num=0 \
  -device loader,file=petalinux-v2019.2/u-boot.elf \
  -device loader,addr=0x30000000,file=linux/2018-04-24/xen \
  -device loader,addr=0x40000000,file=petalinux-v2019.2/Image \
  -nic user -nic user \
  -device virtio-rng-device,bus=virtio-mmio-bus.0,rng=rng0 \
  -object rng-random,filename=/dev/urandom,id=rng0
```

Run the following at the U-Boot prompt:


```
Versal>
fdt addr $fdtcontroladdr
fdt move $fdtcontroladdr 0x20000000
fdt set /timer clock-frequency <0x3dfd240>
fdt set /chosen xen,xen-bootargs "console=dtuart dtuart=/uart@ff000000 dom0_mem=640M_
↪bootscrub=0 maxcpus=1 timer_slop=0"
fdt set /chosen xen,dm0-bootargs "rdinit=/sbin/init clk_ignore_unused console=hvc0_
↪maxcpus=1"
fdt mkn0de /chosen dm0
fdt set /chosen/dm0 compatible "xen,multiboot-module"
fdt set /chosen/dm0 reg <0x00000000 0x40000000 0x0 0x03100000>
booti 30000000 - 20000000
```

Arm CPU features

Arm CPU Features

CPU features are optional features that a CPU of supporting type may choose to implement or not. In QEMU, optional CPU features have corresponding boolean CPU properties that, when enabled, indicate that the feature is implemented, and, conversely, when disabled, indicate that it is not implemented. An example of an Arm CPU feature is the Performance Monitoring Unit (PMU). CPU types such as the Cortex-A15 and the Cortex-A57, which respectively implement Arm architecture reference manuals ARMv7-A and ARMv8-A, may both optionally implement PMUs. For example, if a user wants to use a Cortex-A15 without a PMU, then the `-cpu` parameter should contain `pmu=off` on the QEMU command line, i.e. `-cpu cortex-a15,pmu=off`.

As not all CPU types support all optional CPU features, then whether or not a CPU property exists depends on the CPU type. For example, CPUs that implement the ARMv8-A architecture reference manual may optionally support the AArch32 CPU feature, which may be enabled by disabling the `aarch64` CPU property. A CPU type such as the Cortex-A15, which does not implement ARMv8-A, will not have the `aarch64` CPU property.

QEMU's support may be limited for some CPU features, only partially supporting the feature or only supporting the feature under certain configurations. For example, the `aarch64` CPU feature, which, when disabled, enables the optional AArch32 CPU feature, is only supported when using the KVM accelerator and when running on a host CPU type that supports the feature. While `aarch64` currently only works with KVM, it could work with TCG. CPU features that are specific to KVM are prefixed with "kvm-" and are described in "KVM VCPU Features".

CPU Feature Probing

Determining which CPU features are available and functional for a given CPU type is possible with the `query-cpu-model-expansion` QMP command. Below are some examples where `scripts/qmp/qmp-shell` (see the top comment block in the script for usage) is used to issue the QMP commands.

1. Determine which CPU features are available for the `max` CPU type (Note, we started QEMU with `qemu-system-aarch64`, so `max` is implementing the ARMv8-A reference manual in this case):

```
(QEMU) query-cpu-model-expansion type=full model={"name":"max"}
{ "return": {
  "model": { "name": "max", "props": {
    "sve1664": true, "pmu": true, "sve1792": true, "sve1920": true,
    "sve128": true, "aarch64": true, "sve1024": true, "sve": true,
    "sve640": true, "sve768": true, "sve1408": true, "sve256": true,
    "sve1152": true, "sve512": true, "sve384": true, "sve1536": true,
    "sve896": true, "sve1280": true, "sve2048": true
  }}}}
```

We see that the *max* CPU type has the *pmu*, *aarch64*, *sve*, and many *sve<N>* CPU features. We also see that all the CPU features are enabled, as they are all *true*. (The *sve<N>* CPU features are all optional SVE vector lengths (see “SVE CPU Properties”). While with TCG all SVE vector lengths can be supported, when KVM is in use it’s more likely that only a few lengths will be supported, if SVE is supported at all.)

(2) Let’s try to disable the PMU:

```
(QEMU) query-cpu-model-expansion type=full model={"name":"max","props":{"pmu
↪":false}}
{ "return": {
  "model": { "name": "max", "props": {
    "sve1664": true, "pmu": false, "sve1792": true, "sve1920": true,
    "sve128": true, "aarch64": true, "sve1024": true, "sve": true,
    "sve640": true, "sve768": true, "sve1408": true, "sve256": true,
    "sve1152": true, "sve512": true, "sve384": true, "sve1536": true,
    "sve896": true, "sve1280": true, "sve2048": true
  }}}}
```

We see it worked, as *pmu* is now *false*.

(3) Let’s try to disable *aarch64*, which enables the AArch32 CPU feature:

```
(QEMU) query-cpu-model-expansion type=full model={"name":"max","props":{"aarch64
↪":false}}
{"error": {
  "class": "GenericError", "desc":
    "'aarch64' feature cannot be disabled unless KVM is enabled and 32-bit EL1 is_
↪supported"
}}
```

It looks like this feature is limited to a configuration we do not currently have.

(4) Let’s disable *sve* and see what happens to all the optional SVE vector lengths:

```
(QEMU) query-cpu-model-expansion type=full model={"name":"max","props":{"sve
↪":false}}
{ "return": {
  "model": { "name": "max", "props": {
    "sve1664": false, "pmu": true, "sve1792": false, "sve1920": false,
    "sve128": false, "aarch64": true, "sve1024": false, "sve": false,
    "sve640": false, "sve768": false, "sve1408": false, "sve256": false,
    "sve1152": false, "sve512": false, "sve384": false, "sve1536": false,
    "sve896": false, "sve1280": false, "sve2048": false
  }}}}
```

As expected they are now all *false*.

(5) Let’s try probing CPU features for the Cortex-A15 CPU type:

```
(QEMU) query-cpu-model-expansion type=full model={"name":"cortex-a15"}
{"return": {"model": {"name": "cortex-a15", "props": {"pmu": true}}}}
```

Only the *pmu* CPU feature is available.

A note about CPU feature dependencies

It’s possible for features to have dependencies on other features. I.e. it may be possible to change one feature at a time without error, but when attempting to change all features at once an error could occur depending on the order they

are processed. It's also possible changing all at once doesn't generate an error, because a feature's dependencies are satisfied with other features, but the same feature cannot be changed independently without error. For these reasons callers should always attempt to make their desired changes all at once in order to ensure the collection is valid.

A note about CPU models and KVM

Named CPU models generally do not work with KVM. There are a few cases that do work, e.g. using the named CPU model *cortex-a57* with KVM on a seattle host, but mostly if KVM is enabled the *host* CPU type must be used. This means the guest is provided all the same CPU features as the host CPU type has. And, for this reason, the *host* CPU type should enable all CPU features that the host has by default. Indeed it's even a bit strange to allow disabling CPU features that the host has when using the *host* CPU type, but in the absence of CPU models it's the best we can do if we want to launch guests without all the host's CPU features enabled.

Enabling KVM also affects the *query-cpu-model-expansion* QMP command. The affect is not only limited to specific features, as pointed out in example (3) of "CPU Feature Probing", but also to which CPU types may be expanded. When KVM is enabled, only the *max*, *host*, and current CPU type may be expanded. This restriction is necessary as it's not possible to know all CPU types that may work with KVM, but it does impose a small risk of users experiencing unexpected errors. For example on a seattle, as mentioned above, the *cortex-a57* CPU type is also valid when KVM is enabled. Therefore a user could use the *host* CPU type for the current type, but then attempt to query *cortex-a57*, however that query will fail with our restrictions. This shouldn't be an issue though as management layers and users have been preferring the *host* CPU type for use with KVM for quite some time. Additionally, if the KVM-enabled QEMU instance running on a seattle host is using the *cortex-a57* CPU type, then querying *cortex-a57* will work.

Using CPU Features

After determining which CPU features are available and supported for a given CPU type, then they may be selectively enabled or disabled on the QEMU command line with that CPU type:

```
$ qemu-system-aarch64 -M virt -cpu max,pmu=off,sve=on,sve128=on,sve256=on
```

The example above disables the PMU and enables the first two SVE vector lengths for the *max* CPU type. Note, the *sve=on* isn't actually necessary, because, as we observed above with our probe of the *max* CPU type, *sve* is already on by default. Also, based on our probe of defaults, it would seem we need to disable many SVE vector lengths, rather than only enabling the two we want. This isn't the case, because, as disabling many SVE vector lengths would be quite verbose, the *sve<N>* CPU properties have special semantics (see "SVE CPU Property Parsing Semantics").

KVM VCPU Features

KVM VCPU features are CPU features that are specific to KVM, such as paravirt features or features that enable CPU virtualization extensions. The features' CPU properties are only available when KVM is enabled and are named with the prefix "kvm-". KVM VCPU features may be probed, enabled, and disabled in the same way as other CPU features. Below is the list of KVM VCPU features and their descriptions.

kvm-no-adjvtime By default **kvm-no-adjvtime** is disabled. This means that by default the virtual time adjustment is enabled (vtime is not *not* adjusted).

When virtual time adjustment is enabled each time the VM transitions back to running state the VCPU's virtual counter is updated to ensure stopped time is not counted. This avoids time jumps surprising guest OSes and applications, as long as they use the virtual counter for timekeeping. However it has the side effect of the virtual and physical counters diverging. All timekeeping based on the virtual counter will appear to lag behind any timekeeping that does not subtract VM stopped time. The guest may resynchronize its virtual counter with other time sources as needed.

Enable `kvm-no-adjvtime` to disable virtual time adjustment, also restoring the legacy (pre-5.0) behavior.

kvm-steal-time Since v5.2, **kvm-steal-time is enabled by default** when KVM is enabled, the feature is supported, and the guest is 64-bit.

When `kvm-steal-time` is enabled a 64-bit guest can account for time its CPUs were not running due to the host not scheduling the corresponding VCPU threads. The accounting statistics may influence the guest scheduler behavior and/or be exposed to the guest userspace.

TCG VCPU Features

TCG VCPU features are CPU features that are specific to TCG. Below is the list of TCG VCPU features and their descriptions.

pauth Enable or disable *FEAT_Pauth*, pointer authentication. By default, the feature is enabled with `-cpu max`.

pauth-impdef When *FEAT_Pauth* is enabled, either the *impdef* (Implementation Defined) algorithm is enabled or the *architected* QARMA algorithm is enabled. By default the *impdef* algorithm is disabled, and QARMA is enabled.

The architected QARMA algorithm has good cryptographic properties, but can be quite slow to emulate. The *impdef* algorithm used by QEMU is non-cryptographic but significantly faster.

SVE CPU Properties

There are two types of SVE CPU properties: *sve* and *sve<N>*. The first is used to enable or disable the entire SVE feature, just as the *pmu* CPU property completely enables or disables the PMU. The second type is used to enable or disable specific vector lengths, where *N* is the number of bits of the length. The *sve<N>* CPU properties have special dependencies and constraints, see “SVE CPU Property Dependencies and Constraints” below. Additionally, as we want all supported vector lengths to be enabled by default, then, in order to avoid overly verbose command lines (command lines full of *sve<N>=off*, for all *N* not wanted), we provide the parsing semantics listed in “SVE CPU Property Parsing Semantics”.

SVE CPU Property Dependencies and Constraints

- 1) At least one vector length must be enabled when *sve* is enabled.
- 2) If a vector length *N* is enabled, then, when KVM is enabled, all smaller, host supported vector lengths must also be enabled. If KVM is not enabled, then only all the smaller, power-of-two vector lengths must be enabled. E.g. with KVM if the host supports all vector lengths up to 512-bits (128, 256, 384, 512), then if *sve512* is enabled, the 128-bit vector length, 256-bit vector length, and 384-bit vector length must also be enabled. Without KVM, the 384-bit vector length would not be required.
- 3) If KVM is enabled then only vector lengths that the host CPU type support may be enabled. If SVE is not supported by the host, then no *sve** properties may be enabled.

SVE CPU Property Parsing Semantics

- 1) If SVE is disabled (*sve=off*), then which SVE vector lengths are enabled or disabled is irrelevant to the guest, as the entire SVE feature is disabled and that disables all vector lengths for the guest. However QEMU will still track any *sve<N>* CPU properties provided by the user. If later an *sve=on* is provided, then the guest will get

only the enabled lengths. If no *sve=on* is provided and there are explicitly enabled vector lengths, then an error is generated.

- 2) If SVE is enabled (*sve=on*), but no *sve<N>* CPU properties are provided, then all supported vector lengths are enabled, which when KVM is not in use means including the non-power-of-two lengths, and, when KVM is in use, it means all vector lengths supported by the host processor.
- 3) If SVE is enabled, then an error is generated when attempting to disable the last enabled vector length (see constraint (1) of “SVE CPU Property Dependencies and Constraints”).
- 4) If one or more vector lengths have been explicitly enabled and at least one of the dependency lengths of the maximum enabled length has been explicitly disabled, then an error is generated (see constraint (2) of “SVE CPU Property Dependencies and Constraints”).
- 5) When KVM is enabled, if the host does not support SVE, then an error is generated when attempting to enable any *sve** properties (see constraint (3) of “SVE CPU Property Dependencies and Constraints”).
- 6) When KVM is enabled, if the host does support SVE, then an error is generated when attempting to enable any vector lengths not supported by the host (see constraint (3) of “SVE CPU Property Dependencies and Constraints”).
- 7) If one or more *sve<N>* CPU properties are set *off*, but no *sve<N>*, CPU properties are set *on*, then the specified vector lengths are disabled but the default for any unspecified lengths remains enabled. When KVM is not enabled, disabling a power-of-two vector length also disables all vector lengths larger than the power-of-two length. When KVM is enabled, then disabling any supported vector length also disables all larger vector lengths (see constraint (2) of “SVE CPU Property Dependencies and Constraints”).
- 8) If one or more *sve<N>* CPU properties are set to *on*, then they are enabled and all unspecified lengths default to disabled, except for the required lengths per constraint (2) of “SVE CPU Property Dependencies and Constraints”, which will even be auto-enabled if they were not explicitly enabled.
- 9) If SVE was disabled (*sve=off*), allowing all vector lengths to be explicitly disabled (i.e. avoiding the error specified in (3) of “SVE CPU Property Parsing Semantics”), then if later an *sve=on* is provided an error will be generated. To avoid this error, one must enable at least one vector length prior to enabling SVE.

SVE CPU Property Examples

- 1) Disable SVE:

```
$ qemu-system-aarch64 -M virt -cpu max,sve=off
```

- 2) Implicitly enable all vector lengths for the *max* CPU type:

```
$ qemu-system-aarch64 -M virt -cpu max
```

- 3) When KVM is enabled, implicitly enable all host CPU supported vector lengths with the *host* CPU type:

```
$ qemu-system-aarch64 -M virt,accel=kvm -cpu host
```

- 4) Only enable the 128-bit vector length:

```
$ qemu-system-aarch64 -M virt -cpu max,sve128=on
```

- 5) Disable the 512-bit vector length and all larger vector lengths, since 512 is a power-of-two. This results in all the smaller, uninitialized lengths (128, 256, and 384) defaulting to enabled:

```
$ qemu-system-aarch64 -M virt -cpu max,sve512=off
```

- 6) Enable the 128-bit, 256-bit, and 512-bit vector lengths:

```
$ qemu-system-aarch64 -M virt -cpu max,sve128=on,sve256=on,sve512=on
```

- 7) The same as (6), but since the 128-bit and 256-bit vector lengths are required for the 512-bit vector length to be enabled, then allow them to be auto-enabled:

```
$ qemu-system-aarch64 -M virt -cpu max,sve512=on
```

- 8) Do the same as (7), but by first disabling SVE and then re-enabling it:

```
$ qemu-system-aarch64 -M virt -cpu max,sve=off,sve512=on,sve=on
```

- 9) Force errors regarding the last vector length:

```
$ qemu-system-aarch64 -M virt -cpu max,sve128=off
$ qemu-system-aarch64 -M virt -cpu max,sve=off,sve128=off,sve=on
```

SVE CPU Property Recommendations

The examples in “SVE CPU Property Examples” exhibit many ways to select vector lengths which developers may find useful in order to avoid overly verbose command lines. However, the recommended way to select vector lengths is to explicitly enable each desired length. Therefore only example’s (1), (4), and (6) exhibit recommended uses of the properties.

1.19.7 ColdFire System emulator

Use the executable `qemu-system-m68k` to simulate a ColdFire machine. The emulator is able to boot a uClinux kernel.

The M5208EVB emulation includes the following devices:

- MCF5208 ColdFire V2 Microprocessor (ISA A+ with EMAC).
- Three Two on-chip UARTs.
- Fast Ethernet Controller (FEC)

The AN5206 emulation includes the following devices:

- MCF5206 ColdFire V2 Microprocessor.
- Two on-chip UARTs.

1.19.8 Xtensa System emulator

Two executables cover simulation of both Xtensa endian options, `qemu-system-xtensa` and `qemu-system-xtensaeb`. Two different machine types are emulated:

- Xtensa emulator pseudo board “sim”
- Avnet LX60/LX110/LX200 board

The sim pseudo board emulation provides an environment similar to one provided by the proprietary Tensilica ISS. It supports:

- A range of Xtensa CPUs, default is the DC232B

- Console and filesystem access via semihosting calls

The Avnet LX60/LX110/LX200 emulation supports:

- A range of Xtensa CPUs, default is the DC232B
- 16550 UART
- OpenCores 10/100 Mbps Ethernet MAC

1.19.9 s390x System emulator

QEMU can emulate z/Architecture (in particular, 64 bit) s390x systems via the `qemu-system-s390x` binary. Only one machine type, `s390-ccw-virtio`, is supported (with versioning for compatibility handling).

When using KVM as accelerator, QEMU can emulate CPUs up to the generation of the host. When using the default cpu model with TCG as accelerator, QEMU will emulate a subset of z13 cpu features that should be enough to run distributions built for the z13.

Device support

QEMU will not emulate most of the traditional devices found under LPAR or z/VM; virtio devices (especially using `virtio-ccw`) make up the bulk of the available devices. Passthrough of host devices via `vfio-pci`, `vfio-ccw`, or `vfio-ap` is also available.

Adjunct Processor (AP) Device

Contents

- *Adjunct Processor (AP) Device*
 - *Introduction*
 - *AP Architectural Overview*
 - *Start Interpretive Execution (SIE) Instruction*
 - * *Example 1: Valid configuration*
 - * *Example 2: Valid configuration*
 - * *Example 3: Invalid configuration*
 - *AP Matrix Configuration on Linux Host*
 - * *Binding AP devices to device drivers*
 - * *Configuring an AP matrix for a linux guest*
 - * *Starting a Linux Guest Configured with an AP Matrix*
 - * *Hot plug a vfio-ap device into a running guest*
 - * *Hot unplug a vfio-ap device from a running guest*
 - *Example: Configure AP Matrices for Three Linux Guests*
 - *Limitations*

Introduction

The IBM Adjunct Processor (AP) Cryptographic Facility is comprised of three AP instructions and from 1 to 256 PCIe cryptographic adapter cards. These AP devices provide cryptographic functions to all CPUs assigned to a linux system running in an IBM Z system LPAR.

On s390x, AP adapter cards are exposed via the AP bus. This document describes how those cards may be made available to KVM guests using the VFIO mediated device framework.

AP Architectural Overview

In order understand the terminology used in the rest of this document, let's start with some definitions:

- AP adapter

An AP adapter is an IBM Z adapter card that can perform cryptographic functions. There can be from 0 to 256 adapters assigned to an LPAR depending on the machine model. Adapters assigned to the LPAR in which a linux host is running will be available to the linux host. Each adapter is identified by a number from 0 to 255; however, the maximum adapter number allowed is determined by machine model. When installed, an AP adapter is accessed by AP instructions executed by any CPU.

- AP domain

An adapter is partitioned into domains. Each domain can be thought of as a set of hardware registers for processing AP instructions. An adapter can hold up to 256 domains; however, the maximum domain number allowed is determined by machine model. Each domain is identified by a number from 0 to 255. Domains can be further classified into two types:

- Usage domains are domains that can be accessed directly to process AP commands
- Control domains are domains that are accessed indirectly by AP commands sent to a usage domain to control or change the domain; for example, to set a secure private key for the domain.

- AP Queue

An AP queue is the means by which an AP command-request message is sent to an AP usage domain inside a specific AP. An AP queue is identified by a tuple comprised of an AP adapter ID (APID) and an AP queue index (APQI). The APQI corresponds to a given usage domain number within the adapter. This tuple forms an AP Queue Number (APQN) uniquely identifying an AP queue. AP instructions include a field containing the APQN to identify the AP queue to which the AP command-request message is to be sent for processing.

- AP Instructions:

There are three AP instructions:

- NQAP: to enqueue an AP command-request message to a queue
- DQAP: to dequeue an AP command-reply message from a queue
- PQAP: to administer the queues

AP instructions identify the domain that is targeted to process the AP command; this must be one of the usage domains. An AP command may modify a domain that is not one of the usage domains, but the modified domain must be one of the control domains.

Start Interpretive Execution (SIE) Instruction

A KVM guest is started by executing the Start Interpretive Execution (SIE) instruction. The SIE state description is a control block that contains the state information for a KVM guest and is supplied as input to the SIE instruction.

The SIE state description contains a satellite control block called the Crypto Control Block (CRYCB). The CRYCB contains three fields to identify the adapters, usage domains and control domains assigned to the KVM guest:

- The AP Mask (APM) field is a bit mask that identifies the AP adapters assigned to the KVM guest. Each bit in the mask, from left to right, corresponds to an APID from 0-255. If a bit is set, the corresponding adapter is valid for use by the KVM guest.
- The AP Queue Mask (AQM) field is a bit mask identifying the AP usage domains assigned to the KVM guest. Each bit in the mask, from left to right, corresponds to an AP queue index (APQI) from 0-255. If a bit is set, the corresponding queue is valid for use by the KVM guest.
- The AP Domain Mask field is a bit mask that identifies the AP control domains assigned to the KVM guest. The ADM bit mask controls which domains can be changed by an AP command-request message sent to a usage domain from the guest. Each bit in the mask, from left to right, corresponds to a domain from 0-255. If a bit is set, the corresponding domain can be modified by an AP command-request message sent to a usage domain.

If you recall from the description of an AP Queue, AP instructions include an APQN to identify the AP adapter and AP queue to which an AP command-request message is to be sent (NQAP and PQAP instructions), or from which a command-reply message is to be received (DQAP instruction). The validity of an APQN is defined by the matrix calculated from the APM and AQM; it is the cross product of all assigned adapter numbers (APM) with all assigned queue indexes (AQM). For example, if adapters 1 and 2 and usage domains 5 and 6 are assigned to a guest, the APQNs (1,5), (1,6), (2,5) and (2,6) will be valid for the guest.

The APQNs can provide secure key functionality - i.e., a private key is stored on the adapter card for each of its domains - so each APQN must be assigned to at most one guest or the linux host.

Example 1: Valid configuration

	Guest1	Guest2
adapters	1, 2	1, 2
domains	5, 6	7

This is valid because both guests have a unique set of APQNs:

- Guest1 has APQNs (1,5), (1,6), (2,5) and (2,6);
- Guest2 has APQNs (1,7) and (2,7).

Example 2: Valid configuration

	Guest1	Guest2
adapters	1, 2	3, 4
domains	5, 6	5, 6

This is also valid because both guests have a unique set of APQNs:

- Guest1 has APQNs (1,5), (1,6), (2,5), (2,6);
- Guest2 has APQNs (3,5), (3,6), (4,5), (4,6)

Example 3: Invalid configuration

	Guest1	Guest2
adapters	1, 2	1
domains	5, 6	6, 7

This is an invalid configuration because both guests have access to APQN (1,6).

AP Matrix Configuration on Linux Host

A linux system is a guest of the LPAR in which it is running and has access to the AP resources configured for the LPAR. The LPAR's AP matrix is configured via its Activation Profile which can be edited on the HMC. When the linux system is started, the AP bus will detect the AP devices assigned to the LPAR and create the following in sysfs:

```
/sys/bus/ap
... [devices]
..... xx.yyyy
..... ...
..... cardxx
..... ...
```

Where:

cardxx is AP adapter number xx (in hex)

xx.yyyy is an APQN with xx specifying the APID and yyyy specifying the APQI

For example, if AP adapters 5 and 6 and domains 4, 71 (0x47), 171 (0xab) and 255 (0xff) are configured for the LPAR, the sysfs representation on the linux host system would look like this:

```
/sys/bus/ap
... [devices]
..... 05.0004
..... 05.0047
..... 05.00ab
..... 05.00ff
..... 06.0004
..... 06.0047
..... 06.00ab
..... 06.00ff
..... card05
..... card06
```

A set of default device drivers are also created to control each type of AP device that can be assigned to the LPAR on which a linux host is running:

```
/sys/bus/ap
... [drivers]
..... [cex2acard]      for Crypto Express 2/3 accelerator cards
..... [cex2aqueue]    for AP queues served by Crypto Express 2/3
                      accelerator cards
..... [cex4card]      for Crypto Express 4/5/6 accelerator and coprocessor
                      cards
..... [cex4queue]    for AP queues served by Crypto Express 4/5/6
                      accelerator and coprocessor cards
```

(continues on next page)

(continued from previous page)

```

..... [pcixcccard]      for Crypto Express 2/3 coprocessor cards
..... [pcixccqueue]    for AP queues served by Crypto Express 2/3
                        coprocessor cards

```

Binding AP devices to device drivers

There are two sysfs files that specify bitmasks marking a subset of the APQN range as ‘usable by the default AP queue device drivers’ or ‘not usable by the default device drivers’ and thus available for use by the alternate device driver(s). The sysfs locations of the masks are:

```

/sys/bus/ap/apmask
/sys/bus/ap/aqmask

```

The `apmask` is a 256-bit mask that identifies a set of AP adapter IDs (APID). Each bit in the mask, from left to right (i.e., from most significant to least significant bit in big endian order), corresponds to an APID from 0-255. If a bit is set, the APID is marked as usable only by the default AP queue device drivers; otherwise, the APID is usable by the `vfiio_ap` device driver.

The `aqmask` is a 256-bit mask that identifies a set of AP queue indexes (APQI). Each bit in the mask, from left to right (i.e., from most significant to least significant bit in big endian order), corresponds to an APQI from 0-255. If a bit is set, the APQI is marked as usable only by the default AP queue device drivers; otherwise, the APQI is usable by the `vfiio_ap` device driver.

Take, for example, the following mask:

```
0x7dffffffffffffffffffffffffffffffffffffffffffffffffffffffffffff
```

It indicates:

1, 2, 3, 4, 5, and 7-255 belong to the default drivers’ pool, and 0 and 6 belong to the `vfiio_ap` device driver’s pool.

The APQN of each AP queue device assigned to the linux host is checked by the AP bus against the set of APQNs derived from the cross product of APIDs and APQIs marked as usable only by the default AP queue device drivers. If a match is detected, only the default AP queue device drivers will be probed; otherwise, the `vfiio_ap` device driver will be probed.

By default, the two masks are set to reserve all APQNs for use by the default AP queue device drivers. There are two ways the default masks can be changed:

1. The sysfs mask files can be edited by echoing a string into the respective sysfs mask file in one of two formats:

- An absolute hex string starting with 0x - like “0x12345678” - sets the mask. If the given string is shorter than the mask, it is padded with 0s on the right; for example, specifying a mask value of 0x41 is the same as specifying:

```
0x4100000000000000000000000000000000000000000000000000000000000000
```

Keep in mind that the mask reads from left to right (i.e., most significant to least significant bit in big endian order), so the mask above identifies device numbers 1 and 7 (01000001).

If the string is longer than the mask, the operation is terminated with an error (EINVAL).

- Individual bits in the mask can be switched on and off by specifying each bit number to be switched in a comma separated list. Each bit number string must be prepended with a (+) or minus (-) to indicate the corresponding bit is to be switched on (+) or off (-). Some valid values are:

```
"+0"    switches bit 0 on
"-13"   switches bit 13 off
"+0x41" switches bit 65 on
"-0xff" switches bit 255 off
```

The following example:

+0, -6, +0x47, -0xf0

Switches bits 0 and 71 (0x47) on Switches bits 6 and 240 (0xf0) off

Note that the bits not specified in the list remain as they were before the operation.

2. The masks can also be changed at boot time via parameters on the kernel command line like this:

```
ap.apmask=0xffff ap.aqmask=0x40
```

This would create the following masks:

apmask:

```
0xffff000000000000000000000000000000000000000000000000000000000000
```

aqmask:

[illegible]

Resulting in these two pools:

```
default drivers pool:  adapter 0-15, domain 1
alternate drivers pool: adapter 16-255, domains 0, 2-255
```

Configuring an AP matrix for a linux quest

The sysfs interfaces for configuring an AP matrix for a guest are built on the VFIO mediated device framework. To configure an AP matrix for a guest, a mediated matrix device must first be created for the `/sys/devices/vfio_ap/matrix` device. When the `vfio_ap` device driver is loaded, it registers with the VFIO mediated device framework. When the driver registers, the sysfs interfaces for creating mediated matrix devices is created:

```
/sys/devices
... [vfi0_ap]
.....[matrix]
..... [mdev_supported_types]
..... [vfi0_ap-passthrough]
..... create
..... [devices]
```

A mediated AP matrix device is created by writing a UUID to the attribute file named `create`, for example:

```
uuidgen > create
```

or

```
echo $uuid > create
```

When a mediated AP matrix device is created, a `sysfs` directory named after the UUID is created in the `devices` subdirectory:

```
/sys/devices
... [vfi0_ap]
.....[matrix]
..... [mdev_supported_types]
..... [vfi0_ap-passthrough]
..... create
..... [devices]
..... [$uuid]
```

There will also be three sets of attribute files created in the mediated matrix device's sysfs directory to configure an AP matrix for the KVM guest:

```
/sys/devices
... [vfi0_ap]
.....[matrix]
..... [mdev_supported_types]
..... [vfi0_ap-passthrough]
..... create
..... [devices]
..... [$uuid]
..... assign_adapter
..... assign_control_domain
..... assign_domain
..... matrix
..... unassign_adapter
..... unassign_control_domain
..... unassign_domain
```

assign_adapter To assign an AP adapter to the mediated matrix device, its APID is written to the `assign_adapter` file. This may be done multiple times to assign more than one adapter. The APID may be specified using conventional semantics as a decimal, hexadecimal, or octal number. For example, to assign adapters 4, 5 and 16 to a mediated matrix device in decimal, hexadecimal and octal respectively:

```
echo 4 > assign_adapter
echo 0x5 > assign_adapter
echo 020 > assign_adapter
```

In order to successfully assign an adapter:

- The adapter number specified must represent a value from 0 up to the maximum adapter number allowed by the machine model. If an adapter number higher than the maximum is specified, the operation will terminate with an error (ENODEV).
- All APQNs that can be derived from the adapter ID being assigned and the IDs of the previously assigned domains must be bound to the `vfi0_ap` device driver. If no domains have yet been assigned, then there must be at least one APQN with the specified APID bound to the `vfi0_ap` driver. If no such APQNs are bound to the driver, the operation will terminate with an error (EADDRNOTAVAIL).
- No APQN that can be derived from the adapter ID and the IDs of the previously assigned domains can be assigned to another mediated matrix device. If an APQN is assigned to another mediated matrix device, the operation will terminate with an error (EADDRINUSE).

unassign_adapter To unassign an AP adapter, its APID is written to the `unassign_adapter` file. This may also be done multiple times to unassign more than one adapter.

assign_domain To assign a usage domain, the domain number is written into the `assign_domain` file. This may be done multiple times to assign more than one usage domain. The domain number is specified using conventional semantics as a decimal, hexadecimal, or octal number. For example, to assign usage domains 4, 8, and 71 to a mediated matrix device in decimal, hexadecimal and octal respectively:

```
echo 4 > assign_domain
echo 0x8 > assign_domain
echo 0107 > assign_domain
```

In order to successfully assign a domain:

- The domain number specified must represent a value from 0 up to the maximum domain number allowed by the machine model. If a domain number higher than the maximum is specified, the operation will terminate with an error (ENODEV).
- All APQNs that can be derived from the domain ID being assigned and the IDs of the previously assigned adapters must be bound to the vfio_ap device driver. If no domains have yet been assigned, then there must be at least one APQN with the specified APQI bound to the vfio_ap driver. If no such APQNs are bound to the driver, the operation will terminate with an error (EADDRNOTAVAIL).
- No APQN that can be derived from the domain ID being assigned and the IDs of the previously assigned adapters can be assigned to another mediated matrix device. If an APQN is assigned to another mediated matrix device, the operation will terminate with an error (EADDRINUSE).

unassign_domain To unassign a usage domain, the domain number is written into the `unassign_domain` file. This may be done multiple times to unassign more than one usage domain.

assign_control_domain To assign a control domain, the domain number is written into the `assign_control_domain` file. This may be done multiple times to assign more than one control domain. The domain number may be specified using conventional semantics as a decimal, hexadecimal, or octal number. For example, to assign control domains 4, 8, and 71 to a mediated matrix device in decimal, hexadecimal and octal respectively:

```
echo 4 > assign_domain
echo 0x8 > assign_domain
echo 0107 > assign_domain
```

In order to successfully assign a control domain, the domain number specified must represent a value from 0 up to the maximum domain number allowed by the machine model. If a control domain number higher than the maximum is specified, the operation will terminate with an error (ENODEV).

unassign_control_domain To unassign a control domain, the domain number is written into the `unassign_domain` file. This may be done multiple times to unassign more than one control domain.

Notes: No changes to the AP matrix will be allowed while a guest using the mediated matrix device is running. Attempts to assign an adapter, domain or control domain will be rejected and an error (EBUSY) returned.

Starting a Linux Guest Configured with an AP Matrix

To provide a mediated matrix device for use by a guest, the following option must be specified on the QEMU command line:

```
-device vfio_ap,sysfsdev=$path-to-mdev
```

The `sysfsdev` parameter specifies the path to the mediated matrix device. There are a number of ways to specify this path:

```
/sys/devices/vfio_ap/matrix/$uuid
/sys/bus/mdev/devices/$uuid
/sys/bus/mdev/drivers/vfio_mdev/$uuid
/sys/devices/vfio_ap/matrix/mdev_supported_types/vfio_ap-passthrough/devices/$uuid
```

When the linux guest is started, the guest will open the mediated matrix device's file descriptor to get information about the mediated matrix device. The `vfio_ap` device driver will update the APM, AQM, and ADM fields in the guest's CRYCB with the adapter, usage domain and control domains assigned via the mediated matrix device's sysfs attribute files. Programs running on the linux guest will then:

1. Have direct access to the APQNs derived from the cross product of the AP adapter numbers (APID) and queue indexes (APQI) specified in the APM and AQM fields of the guests's CRYCB respectively. These APQNs identify the AP queues that are valid for use by the guest; meaning, AP commands can be sent by the guest to any of these queues for processing.
2. Have authorization to process AP commands to change a control domain identified in the ADM field of the guest's CRYCB. The AP command must be sent to a valid APQN (see 1 above).

CPU model features:

Three CPU model features are available for controlling guest access to AP facilities:

1. AP facilities feature

The AP facilities feature indicates that AP facilities are installed on the guest. This feature will be exposed for use only if the AP facilities are installed on the host system. The feature is s390-specific and is represented as a parameter of the `-cpu` option on the QEMU command line:

```
qemu-system-s390x -cpu $model,ap=on|off
```

Where:

\$model is the CPU model defined for the guest (defaults to the model of the host system if not specified).

ap=on|off indicates whether AP facilities are installed (on) or not (off). The default for CPU models zEC12 or newer is `ap=on`. AP facilities must be installed on the guest if a `vfio-ap` device (`-device vfio-ap,sysfsdev=$path`) is configured for the guest, or the guest will fail to start.

2. Query Configuration Information (QCI) facility

The QCI facility is used by the AP bus running on the guest to query the configuration of the AP facilities. This facility will be available only if the QCI facility is installed on the host system. The feature is s390-specific and is represented as a parameter of the `-cpu` option on the QEMU command line:

```
qemu-system-s390x -cpu $model,apqci=on|off
```

Where:

\$model is the CPU model defined for the guest

apqci=on|off indicates whether the QCI facility is installed (on) or not (off). The default for CPU models zEC12 or newer is `apqci=on`; for older models, QCI will not be installed.

If QCI is installed (`apqci=on`) but AP facilities are not (`ap=off`), an error message will be logged, but the guest will be allowed to start. It makes no sense to have QCI installed if the AP facilities are not; this is considered an invalid configuration.

If the QCI facility is not installed, APQNs with an APQI greater than 15 will not be detected by the AP bus running on the guest.

3. Adjunct Process Facility Test (APFT) facility

The APFT facility is used by the AP bus running on the guest to test the AP facilities available for a given AP queue. This facility will be available only if the APFT facility is installed on the host system. The feature is s390-specific and is represented as a parameter of the `-cpu` option on the QEMU command line:

```
qemu-system-s390x -cpu $model,apft=on|off
```

Where:

\$model is the CPU model defined for the guest (defaults to the model of the host system if not specified).

apft=on|off indicates whether the APFT facility is installed (on) or not (off). The default for CPU models zEC12 and newer is `apft=on` for older models, APFT will not be installed.

If APFT is installed (`apft=on`) but AP facilities are not (`ap=off`), an error message will be logged, but the guest will be allowed to start. It makes no sense to have APFT installed if the AP facilities are not; this is considered an invalid configuration.

It also makes no sense to turn APFT off because the AP bus running on the guest will not detect CEX4 and newer devices without it. Since only CEX4 and newer devices are supported for guest usage, no AP devices can be made accessible to a guest started without APFT installed.

Hot plug a vfiio-ap device into a running guest

Only one vfiio-ap device can be attached to the virtual machine's ap-bus, so a vfiio-ap device can be hot plugged if and only if no vfiio-ap device is attached to the bus already, whether via the QEMU command line or a prior hot plug action.

To hot plug a vfiio-ap device, use the QEMU `device_add` command:

```
(qemu) device_add vfiio-ap,sysfsdev="$path-to-mdev",id="$id"
```

Where the `$path-to-mdev` value specifies the absolute path to a mediated device to which AP resources to be used by the guest have been assigned. `$id` is the name value for the optional `id` parameter.

Note that on Linux guests, the AP devices will be created in the `/sys/bus/ap/devices` directory when the AP bus subsequently performs its periodic scan, so there may be a short delay before the AP devices are accessible on the guest.

The command will fail if:

- A vfiio-ap device has already been attached to the virtual machine's ap-bus.
- The CPU model features for controlling guest access to AP facilities are not enabled (see 'CPU model features' subsection in the previous section).

Hot unplug a vfiio-ap device from a running guest

A vfiio-ap device can be unplugged from a running KVM guest if a vfiio-ap device has been attached to the virtual machine's ap-bus via the QEMU command line or a prior hot plug action.

To hot unplug a vfiio-ap device, use the QEMU `device_del` command:

```
(qemu) device_del "$id"
```

Where `$id` is the same `id` that was specified at device creation.

On a Linux guest, the AP devices will be removed from the `/sys/bus/ap/devices` directory on the guest when the AP bus subsequently performs its periodic scan, so there may be a short delay before the AP devices are no longer accessible by the guest.

The command will fail if the `$path-to-mdev` specified on the `device_del` command does not match the value specified when the `vfiio-ap` device was attached to the virtual machine's `ap-bus`.

Example: Configure AP Matrices for Three Linux Guests

Let's now provide an example to illustrate how KVM guests may be given access to AP facilities. For this example, we will show how to configure three guests such that executing the `lszcrypt` command on the guests would look like this:

Guest1:

CARD.DOMAIN	TYPE	MODE
05	CEX5C	CCA-Coproc
05.0004	CEX5C	CCA-Coproc
05.00ab	CEX5C	CCA-Coproc
06	CEX5A	Accelerator
06.0004	CEX5A	Accelerator
06.00ab	CEX5C	CCA-Coproc

Guest2:

CARD.DOMAIN	TYPE	MODE
05	CEX5A	Accelerator
05.0047	CEX5A	Accelerator
05.00ff	CEX5A	Accelerator

Guest3:

CARD.DOMAIN	TYPE	MODE
06	CEX5A	Accelerator
06.0047	CEX5A	Accelerator
06.00ff	CEX5A	Accelerator

These are the steps:

1. Install the `vfiio_ap` module on the linux host. The dependency chain for the `vfiio_ap` module is:

- `iommu`
- `s390`
- `zcrypt`
- `vfiio`
- `vfiio_mdev`
- `vfiio_mdev_device`
- `KVM`

To build the `vfiio_ap` module, the kernel build must be configured with the following Kconfig elements selected:

- `IOMMU_SUPPORT`
- `S390`
- `ZCRYPT`

- S390_AP_IOMMU
- VFIO
- VFIO_MDEV
- VFIO_MDEV_DEVICE
- KVM

If using make menuconfig select the following to build the vfio_ap module::

-> Device Drivers

-> IOMMU Hardware Support select S390 AP IOMMU Support

-> VFIO Non-Privileged userspace driver framework

-> Mediated device driver framework -> VFIO driver for Mediated devices

-> I/O subsystem -> VFIO support for AP devices

2. Secure the AP queues to be used by the three guests so that the host can not access them. To secure the AP queues 05.0004, 05.0047, 05.00ab, 05.00ff, 06.0004, 06.0047, 06.00ab, and 06.00ff for use by the vfio_ap device driver, the corresponding APQNs must be removed from the default queue drivers pool as follows:

```
echo -5,-6 > /sys/bus/ap/apmask  
  
echo -4,-0x47,-0xab,-0xff > /sys/bus/ap/aqmask
```

This will result in AP queues 05.0004, 05.0047, 05.00ab, 05.00ff, 06.0004, 06.0047, 06.00ab, and 06.00ff getting bound to the vfio_ap device driver. The sysfs directory for the vfio_ap device driver will now contain symbolic links to the AP queue devices bound to it:

```
/sys/bus/ap  
... [drivers]  
..... [vfio_ap]  
..... [05.0004]  
..... [05.0047]  
..... [05.00ab]  
..... [05.00ff]  
..... [06.0004]  
..... [06.0047]  
..... [06.00ab]  
..... [06.00ff]
```

Keep in mind that only type 10 and newer adapters (i.e., CEX4 and later) can be bound to the vfio_ap device driver. The reason for this is to simplify the implementation by not needlessly complicating the design by supporting older devices that will go out of service in the relatively near future, and for which there are few older systems on which to test.

The administrator, therefore, must take care to secure only AP queues that can be bound to the vfio_ap device driver. The device type for a given AP queue device can be read from the parent card's sysfs directory. For example, to see the hardware type of the queue 05.0004:

```
cat /sys/bus/ap/devices/card05/hwtype
```

The hwtype must be 10 or higher (CEX4 or newer) in order to be bound to the vfio_ap device driver.

3. Create the mediated devices needed to configure the AP matrixes for the three guests and to provide an interface to the vfio_ap driver for use by the guests:

```
/sys/devices/vfio_ap/matrix/
... [mdev_supported_types]
..... [vfio_ap-passthrough] (passthrough mediated matrix device type)
..... create
..... [devices]
```

To create the mediated devices for the three guests:

```
uuidgen > create
uuidgen > create
uuidgen > create
```

or

```
echo $uuid1 > create
echo $uuid2 > create
echo $uuid3 > create
```

This will create three mediated devices in the [devices] subdirectory named after the UUID used to create the mediated device. We'll call them \$uuid1, \$uuid2 and \$uuid3 and this is the sysfs directory structure after creation:

```
/sys/devices/vfio_ap/matrix/
... [mdev_supported_types]
..... [vfio_ap-passthrough]
..... [devices]
..... [$uuid1]
..... assign_adapter
..... assign_control_domain
..... assign_domain
..... matrix
..... unassign_adapter
..... unassign_control_domain
..... unassign_domain

..... [$uuid2]
..... assign_adapter
..... assign_control_domain
..... assign_domain
..... matrix
..... unassign_adapter
..... unassign_control_domain
..... unassign_domain

..... [$uuid3]
..... assign_adapter
..... assign_control_domain
..... assign_domain
..... matrix
..... unassign_adapter
..... unassign_control_domain
..... unassign_domain
```

4. The administrator now needs to configure the matrixes for the mediated devices \$uuid1 (for Guest1), \$uuid2 (for Guest2) and \$uuid3 (for Guest3).

This is how the matrix is configured for Guest1:

```
echo 5 > assign_adapter
echo 6 > assign_adapter
echo 4 > assign_domain
echo 0xab > assign_domain
```

Control domains can similarly be assigned using the `assign_control_domain` sysfs file.

If a mistake is made configuring an adapter, domain or control domain, you can use the `unassign_XXX` interfaces to unassign the adapter, domain or control domain.

To display the matrix configuration for Guest1:

```
cat matrix
```

The output will display the APQNs in the format `xx.yyyy`, where `xx` is the adapter number and `yyyy` is the domain number. The output for Guest1 will look like this:

```
05.0004
05.00ab
06.0004
06.00ab
```

This is how the matrix is configured for Guest2:

```
echo 5 > assign_adapter
echo 0x47 > assign_domain
echo 0xff > assign_domain
```

This is how the matrix is configured for Guest3:

```
echo 6 > assign_adapter
echo 0x47 > assign_domain
echo 0xff > assign_domain
```

5. Start Guest1:

```
/usr/bin/qemu-system-s390x ... -cpu host,ap=on,apqci=on,apft=on -device vfio-ap,
↪sysfsdev=/sys/devices/vfio_ap/matrix/$uuid1 ...
```

7. Start Guest2:

```
/usr/bin/qemu-system-s390x ... -cpu host,ap=on,apqci=on,apft=on -device vfio-ap,
↪sysfsdev=/sys/devices/vfio_ap/matrix/$uuid2 ...
```

7. Start Guest3:

```
/usr/bin/qemu-system-s390x ... -cpu host,ap=on,apqci=on,apft=on -device vfio-ap,
↪sysfsdev=/sys/devices/vfio_ap/matrix/$uuid3 ...
```

When the guest is shut down, the mediated matrix devices may be removed.

Using our example again, to remove the mediated matrix device `$uuid1`:

```
/sys/devices/vfio_ap/matrix/
... [mdev_supported_types]
..... [vfio_ap-passthrough]
..... [devices]
..... [$uuid1]
```

(continues on next page)

(continued from previous page)

```
..... remove

echo 1 > remove
```

This will remove all of the mdev matrix device's sysfs structures including the mdev device itself. To recreate and reconfigure the mdev matrix device, all of the steps starting with step 3 will have to be performed again. Note that the remove will fail if a guest using the mdev is still running.

It is not necessary to remove an mdev matrix device, but one may want to remove it if no guest will use it during the remaining lifetime of the linux host. If the mdev matrix device is removed, one may want to also reconfigure the pool of adapters and queues reserved for use by the default drivers.

Limitations

- The KVM/kernel interfaces do not provide a way to prevent restoring an APQN to the default drivers pool of a queue that is still assigned to a mediated device in use by a guest. It is incumbent upon the administrator to ensure there is no mediated device in use by a guest to which the APQN is assigned lest the host be given access to the private data of the AP queue device, such as a private key configured specifically for the guest.
- Dynamically assigning AP resources to or unassigning AP resources from a mediated matrix device - see [Configuring an AP matrix for a linux guest](#) section above - while a running guest is using it is currently not supported.
- Live guest migration is not supported for guests using AP devices. If a guest is using AP devices, the vfio-ap device configured for the guest must be unplugged before migrating the guest (see [Hot unplug a vfio-ap device from a running guest](#) section above.)

The virtual channel subsystem

QEMU implements a virtual channel subsystem with subchannels, (mostly functionless) channel paths, and channel devices (virtio-ccw, 3270, and devices passed via vfio-ccw). It supports multiple subchannel sets (MSS) and multiple channel subsystems extended (MCSS-E).

All channel devices support the `devno` property, which takes a parameter in the form `<cssid>.<ssid>.<device number>`.

The default channel subsystem image id (`<cssid>`) is `0xfe`. Devices in there will show up in channel subsystem image 0 to guests that do not enable MCSS-E. Note that devices with a different `cssid` will not be visible if the guest OS does not enable MCSS-E (which is true for all supported guest operating systems today).

Supported values for the subchannel set id (`<ssid>`) range from 0-3. Devices with a `ssid` that is not 0 will not be visible if the guest OS does not enable MSS (any Linux version that supports virtio also enables MSS). Any device may be put into any subchannel set, there is no restriction by device type.

The device number can range from 0-0xffff.

If the `devno` property is not specified for a device, QEMU will choose the next free device number in subchannel set 0, skipping to the next subchannel set if no more device numbers are free.

QEMU places a device at the first free subchannel in the specified subchannel set. If a device is hotunplugged and later replugged, it may appear at a different subchannel. (This is similar to how z/VM works.)

Examples

- a virtio-net device, `cssid/ssid/devno` automatically assigned:

```
-device virtio-net-ccw
```

In a Linux guest (without default devices and no other devices specified prior to this one), this will show up as 0.0.0000 under subchannel 0.0.0000.

The auto-assigned-properties in QEMU (as seen via e.g. `info qtree`) would be `dev_id = "fe.0.0000"` and `subch_id = "fe.0.0000"`.

- a virtio-rng device in subchannel set 0:

```
-device virtio-rng-ccw,devno=fe.0.0042
```

If added to the same Linux guest as above, it would show up as 0.0.0042 under subchannel 0.0.0001.

The properties for the device would be `dev_id = "fe.0.0042"` and `subch_id = "fe.0.0001"`.

- a virtio-gpu device in subchannel set 2:

```
-device virtio-gpu-ccw,devno=fe.2.1111
```

If added to the same Linux guest as above, it would show up as 0.2.1111 under subchannel 0.2.0000.

The properties for the device would be `dev_id = "fe.2.1111"` and `subch_id = "fe.2.0000"`.

- a virtio-mouse device in a non-standard channel subsystem image:

```
-device virtio-mouse-ccw,devno=2.0.2222
```

This would not show up in a standard Linux guest.

The properties for the device would be `dev_id = "2.0.2222"` and `subch_id = "2.0.0000"`.

- a virtio-keyboard device in another non-standard channel subsystem image:

```
-device virtio-keyboard-ccw,devno=0.0.1234
```

This would not show up in a standard Linux guest, either, as 0 is not the standard channel subsystem image id.

The properties for the device would be `dev_id = "0.0.1234"` and `subch_id = "0.0.0000"`.

3270 devices

The 3270 is the classic ‘green-screen’ console of the mainframes (see the [IBM 3270 Wikipedia article](#)).

The 3270 data stream is not implemented within QEMU; the device only provides TN3270 (a telnet extension; see [RFC 854](#) and [RFC 1576](#)) and leaves the heavy lifting to an external 3270 terminal emulator (such as `x3270`) to make a single 3270 device available to a guest. Note that this supports basic features only.

To provide a 3270 device to a guest, create a `x-terminal3270` linked to a `tn3270` chardev. The guest will see a 3270 channel device. In order to actually be able to use it, attach the `x3270` emulator to the chardev.

Example configuration

- Make sure that 3270 support is enabled in the guest’s Linux kernel. You need `CONFIG_TN3270` and at least one of `CONFIG_TN3270_TTY` (for additional ttys) or `CONFIG_TN3270_CONSOLE` (for a 3270 console).
- Add a `tn3270` chardev and a `x-terminal3270` to the QEMU command line:

```
-chardev socket,id=ch0,host=0.0.0.0,port=2300,nowait,server,tn3270
-device x-terminal3270,chardev=ch0,devno=fe.0.000a,id=terminal0
```

- Start the guest. In the guest, use `chccwdev -e 0.0.000a` to enable the device.
- On the host, start the `x3270` emulator:

```
x3270 <host>:2300
```

- In the guest, locate the 3270 device node under `/dev/3270/` (say, `ttty1`) and start a getty on it:

```
systemctl start serial-getty@3270-tty1.service
```

This should get you an additional tty for logging into the guest.

- If you want to use the 3270 device as the Linux kernel console instead of an additional tty, you can also append `conmode=3270 condev=000a` to the guest's kernel command line. The kernel then should use the 3270 as console after the next boot.

Restrictions

3270 support is very basic. In particular:

- Only one 3270 device is supported.
- It has only been tested with Linux guests and the `x3270` emulator.
- TLS/SSL is not supported.
- Resizing on reattach is not supported.
- Multiple commands in one inbound buffer (for example, when the reset key is pressed while the network is slow) are not supported.

Subchannel passthrough via vfio-ccw

vfio-ccw (based upon the mediated vfio device infrastructure) allows to make certain I/O subchannels and their devices available to a guest. The host will not interact with those subchannels/devices any more.

Note that while vfio-ccw should work with most non-QDIO devices, only ECKD DASDs have really been tested.

Example configuration

Step 1: configure the host device

As every mdev is identified by a uuid, the first step is to obtain one:

```
[root@host ~]# uuidgen
7e270a25-e163-4922-af60-757fc8ed48c6
```

Note: it is recommended to use the `mdevctl` tool for actually configuring the host device.

To define the same device as configured below to be started automatically, use

```
[root@host ~]# driverctl -b css set-override 0.0.0313 vfio_ccw
[root@host ~]# mdevctl define -u 7e270a25-e163-4922-af60-757fc8ed48c6 \
    -p 0.0.0313 -t vfio_ccw-io -a
```

If using `mdevctl` is not possible or wanted, follow the manual procedure below.

- Locate the subchannel for the device (in this example, `0.0.2b09`):

```
[root@host ~]# lscss | grep 0.0.2b09 | awk '{print $2}'
0.0.0313
```

- Unbind the subchannel (in this example, `0.0.0313`) from the standard I/O subchannel driver and bind it to the `vfio-ccw` driver:

```
[root@host ~]# echo 0.0.0313 > /sys/bus/css/devices/0.0.0313/driver/unbind
[root@host ~]# echo 0.0.0313 > /sys/bus/css/drivers/vfio_ccw/bind
```

- Create the mediated device (identified by the uuid):

```
[root@host ~]# echo "7e270a25-e163-4922-af60-757fc8ed48c6" > \
/sys/bus/css/devices/0.0.0313/mdev_supported_types/vfio_ccw-io/create
```

Step 2: configure QEMU

- Reference the created mediated device and (optionally) pick a device id to be presented in the guest (here, `fe.0.1234`, which will end up visible in the guest as `0.0.1234`):

```
-device vfio-ccw,devno=fe.0.1234,sysfsdev=\
/sys/bus/mdev/devices/7e270a25-e163-4922-af60-757fc8ed48c6
```

- Start the guest. The device (here, `0.0.1234`) should now be usable:

```
[root@guest ~]# lscss -d 0.0.1234
Device    Subchan.  DevType CU Type Use  PIM PAM POM  CHPID
-----
0.0.1234 0.0.0007  3390/0e 3990/e9      f0 f0 ff  1a2a3a0a 00000000
[root@guest ~]# chccwdev -e 0.0.1234
Setting device 0.0.1234 online
Done
[root@guest ~]# dmesg -t
(...)
dasd-eckd 0.0.1234: A channel path to the device has become operational
dasd-eckd 0.0.1234: New DASD 3390/0E (CU 3990/01) with 10017 cylinders, 15 heads,
↳ 224 sectors
dasd-eckd 0.0.1234: DASD with 4 KB/block, 7212240 KB total size, 48 KB/track,
↳ compatible disk layout
dasda:VOL1/ 0X2B09: dasda1
```

Architectural features

Boot devices on s390x

Booting with bootindex parameter

For classical mainframe guests (i.e. LPAR or z/VM installations), you always have to explicitly specify the disk where you want to boot from (or “IPL” from, in s390x-speak – IPL means “Initial Program Load”). In particular, there can also be only one boot device according to the architecture specification, thus specifying multiple boot devices is not possible (yet).

So for booting an s390x guest in QEMU, you should always mark the device where you want to boot from with the `bootindex` property, for example:

```
qemu-system-s390x -drive if=none,id=dr1,file=guest.qcow2 \
                  -device virtio-blk,drive=dr1,bootindex=1
```

For booting from a CD-ROM ISO image (which needs to include El-Torito boot information in order to be bootable), it is recommended to specify a `scsi-cd` device, for example like this:

```
qemu-system-s390x -blockdev file,node-name=c1,filename=... \
                  -device virtio-scsi \
                  -device scsi-cd,drive=c1,bootindex=1
```

Note that you really have to use the `bootindex` property to select the boot device. The old-fashioned `-boot order=...` command of QEMU (and also `-boot once=...`) is not supported on s390x.

Booting without bootindex parameter

The QEMU guest firmware (the so-called s390-ccw bios) has also some rudimentary support for scanning through the available block devices. So in case you did not specify a boot device with the `bootindex` property, there is still a chance that it finds a bootable device on its own and starts a guest operating system from it. However, this scanning algorithm is still very rough and may be incomplete, so that it might fail to detect a bootable device in many cases. It is really recommended to always specify the boot device with the `bootindex` property instead.

This also means that you should avoid the classical short-cut commands like `-hda`, `-cdrom` or `-drive if=virtio`, since it is not possible to specify the `bootindex` with these commands. Note that the convenience `-cdrom` option even does not give you a real (virtio-scsi) CD-ROM device on s390x. Due to technical limitations in the QEMU code base, you will get a virtio-blk device with this parameter instead, which might not be the right device type for installing a Linux distribution via ISO image. It is recommended to specify a CD-ROM device via `-device scsi-cd` (as mentioned above) instead.

Booting from a network device

Beside the normal guest firmware (which is loaded from the file `s390-ccw.img` in the data directory of QEMU, or via the `-bios` option), QEMU ships with a small TFTP network bootloader firmware for virtio-net-ccw devices, too. This firmware is loaded from a file called `s390-netboot.img` in the QEMU data directory. In case you want to load it from a different filename instead, you can specify it via the `-global s390-ipl.netboot_fw=filename` command line option.

The `bootindex` property is especially important for booting via the network. If you don't specify the `bootindex` property here, the network bootloader firmware code won't get loaded into the guest memory so that the network boot will fail. For a successful network boot, try something like this:

```
qemu-system-s390x -netdev user,id=n1,tftp=...,bootfile=... \
                  -device virtio-net-ccw,netdev=n1,bootindex=1
```

The network bootloader firmware also has basic support for pxelinux.cfg-style configuration files. See the [PXELINUX Configuration](#) page for details how to set up the configuration file on your TFTP server. The supported configuration file entries are `DEFAULT`, `LABEL`, `KERNEL`, `INITRD` and `APPEND` (see the [Syslinux Config file syntax](#) for more information).

Protected Virtualization on s390x

The memory and most of the registers of Protected Virtual Machines (PVMs) are encrypted or inaccessible to the hypervisor, effectively prohibiting VM introspection when the VM is running. At rest, PVMs are encrypted and can only be decrypted by the firmware, represented by an entity called Ultravisor, of specific IBM Z machines.

Prerequisites

To run PVMs, a machine with the Protected Virtualization feature, as indicated by the Ultravisor Call facility (stfle bit 158), is required. The Ultravisor needs to be initialized at boot by setting `prot_virt=1` on the host's kernel command line.

Running PVMs requires using the KVM hypervisor.

If those requirements are met, the capability `KVM_CAP_S390_PROTECTED` will indicate that KVM can support PVMs on that LPAR.

Running a Protected Virtual Machine

To run a PVM you will need to select a CPU model which includes the *Unpack facility* (stfle bit 161 represented by the feature `unpack/S390_FEAT_UNPACK`), and add these options to the command line:

```
-object s390-pv-guest,id=pv0 \  
-machine confidential-guest-support=pv0
```

Adding these options will:

- Ensure the *unpack* facility is available
- Enable the IOMMU by default for all I/O devices
- Initialize the PV mechanism

Passthrough (vfiio) devices are currently not supported.

Host huge page backings are not supported. However guests can use huge pages as indicated by its facilities.

Boot Process

A secure guest image can either be loaded from disk or supplied on the QEMU command line. Booting from disk is done by the unmodified s390-ccw BIOS. I.e., the bootmap is interpreted, multiple components are read into memory and control is transferred to one of the components (zipl stage3). Stage3 does some fixups and then transfers control to some program residing in guest memory, which is normally the OS kernel. The secure image has another component prepended (stage3a) that uses the new diag308 subcodes 8 and 10 to trigger the transition into secure mode.

Bootting from the image supplied on the QEMU command line requires that the file passed via `-kernel` has the same memory layout as would result from the disk boot. This memory layout includes the encrypted components (kernel, initrd, cmdline), the stage3a loader and metadata. In case this boot method is used, the command line options `-initrd` and `-cmdline` are ineffective. The preparation of a PVM image is done via the *genproting* tool from the s390-tools collection.

1.19.10 RX System emulator

Use the executable `qemu-system-rx` to simulate RX target (GDB simulator). This target emulated following devices.

- R5F562N8 MCU
 - On-chip memory (ROM 512KB, RAM 96KB)
 - Interrupt Control Unit (ICUa)
 - 8Bit Timer x 1CH (TMR0,1)
 - Compare Match Timer x 2CH (CMT0,1)
 - Serial Communication Interface x 1CH (SCI0)
- External memory 16MByte

Example of `qemu-system-rx` usage for RX is shown below:

Download `<u-boot_image_file>` from <https://osdn.net/users/ysato/pf/qemu/dl/u-boot.bin.gz>

Start emulation of rx-virt:: `qemu-system-rx -M gdbsim-r5f562n8 -bios <u-boot_image_file>`

Download `kernel_image_file` from <https://osdn.net/users/ysato/pf/qemu/dl/zImage>

Download `device_tree_blob` from <https://osdn.net/users/ysato/pf/qemu/dl/rx-virt.dtb>

Start emulation of rx-virt::

```
qemu-system-rx -M gdbsim-r5f562n8 -kernel <kernel_image_file> -dtb <device_tree_blob> -append "early-con"
```

1.19.11 AVR System emulator

Use the executable `qemu-system-avr` to emulate a AVR 8 bit based machine. These can have one of the following cores: `avr1`, `avr2`, `avr25`, `avr3`, `avr31`, `avr35`, `avr4`, `avr5`, `avr51`, `avr6`, `avrtiny`, `xmega2`, `xmega3`, `xmega4`, `xmega5`, `xmega6` and `xmega7`.

As for now it supports few Arduino boards for educational and testing purposes. These boards use a ATmega controller, which model is limited to USART & 16-bit timer devices, enough to run FreeRTOS based applications (like https://github.com/seharris/qemu-avr-tests/blob/master/free-rtos/Demo/AVR_ATMega2560_GCC/demo.elf).

Following are examples of possible usages, assuming `demo.elf` is compiled for AVR cpu

- Continuous non interrupted execution:

```
qemu-system-avr -machine mega2560 -bios demo.elf
```

- Continuous non interrupted execution with serial output into telnet window:

```
qemu-system-avr -M mega2560 -bios demo.elf -nographic \
    -serial tcp::5678,server,nowait
```

and then in another shell:

```
telnet localhost 5678
```

- Debugging with GDB debugger:

```
qemu-system-avr -machine mega2560 -bios demo.elf -s -S
```

and then in another shell:

```
avr-gdb demo.elf
```

and then within GDB shell:

```
target remote :1234
```

- Print out executed instructions (that have not been translated by the JIT compiler yet):

```
qemu-system-avr -machine mega2560 -bios demo.elf -d in_asm
```

1.20 Security

1.20.1 Overview

This chapter explains the security requirements that QEMU is designed to meet and principles for securely deploying QEMU.

1.20.2 Security Requirements

QEMU supports many different use cases, some of which have stricter security requirements than others. The community has agreed on the overall security requirements that users may depend on. These requirements define what is considered supported from a security perspective.

Virtualization Use Case

The virtualization use case covers cloud and virtual private server (VPS) hosting, as well as traditional data center and desktop virtualization. These use cases rely on hardware virtualization extensions to execute guest code safely on the physical CPU at close-to-native speed.

The following entities are untrusted, meaning that they may be buggy or malicious:

- Guest
- User-facing interfaces (e.g. VNC, SPICE, WebSocket)
- Network protocols (e.g. NBD, live migration)
- User-supplied files (e.g. disk images, kernels, device trees)
- Passthrough devices (e.g. PCI, USB)

Bugs affecting these entities are evaluated on whether they can cause damage in real-world use cases and treated as security bugs if this is the case.

Non-virtualization Use Case

The non-virtualization use case covers emulation using the Tiny Code Generator (TCG). In principle the TCG and device emulation code used in conjunction with the non-virtualization use case should meet the same security requirements as the virtualization use case. However, for historical reasons much of the non-virtualization use case code was not written with these security requirements in mind.

Bugs affecting the non-virtualization use case are not considered security bugs at this time. Users with non-virtualization use cases must not rely on QEMU to provide guest isolation or any security guarantees.

1.20.3 Architecture

This section describes the design principles that ensure the security requirements are met.

Guest Isolation

Guest isolation is the confinement of guest code to the virtual machine. When guest code gains control of execution on the host this is called escaping the virtual machine. Isolation also includes resource limits such as throttling of CPU, memory, disk, or network. Guests must be unable to exceed their resource limits.

QEMU presents an attack surface to the guest in the form of emulated devices. The guest must not be able to gain control of QEMU. Bugs in emulated devices could allow malicious guests to gain code execution in QEMU. At this point the guest has escaped the virtual machine and is able to act in the context of the QEMU process on the host.

Guests often interact with other guests and share resources with them. A malicious guest must not gain control of other guests or access their data. Disk image files and network traffic must be protected from other guests unless explicitly shared between them by the user.

Principle of Least Privilege

The principle of least privilege states that each component only has access to the privileges necessary for its function. In the case of QEMU this means that each process only has access to resources belonging to the guest.

The QEMU process should not have access to any resources that are inaccessible to the guest. This way the guest does not gain anything by escaping into the QEMU process since it already has access to those same resources from within the guest.

Following the principle of least privilege immediately fulfills guest isolation requirements. For example, guest A only has access to its own disk image file `a.img` and not guest B's disk image file `b.img`.

In reality certain resources are inaccessible to the guest but must be available to QEMU to perform its function. For example, host system calls are necessary for QEMU but are not exposed to guests. A guest that escapes into the QEMU process can then begin invoking host system calls.

New features must be designed to follow the principle of least privilege. Should this not be possible for technical reasons, the security risk must be clearly documented so users are aware of the trade-off of enabling the feature.

Isolation mechanisms

Several isolation mechanisms are available to realize this architecture of guest isolation and the principle of least privilege. With the exception of Linux seccomp, these mechanisms are all deployed by management tools that launch QEMU, such as libvirt. They are also platform-specific so they are only described briefly for Linux here.

The fundamental isolation mechanism is that QEMU processes must run as unprivileged users. Sometimes it seems more convenient to launch QEMU as root to give it access to host devices (e.g. `/dev/net/tun`) but this poses a huge security risk. File descriptor passing can be used to give an otherwise unprivileged QEMU process access to host devices without running QEMU as root. It is also possible to launch QEMU as a non-root user and configure UNIX groups for access to `/dev/kvm`, `/dev/net/tun`, and other device nodes. Some Linux distros already ship with UNIX groups for these devices by default.

- SELinux and AppArmor make it possible to confine processes beyond the traditional UNIX process and file permissions model. They restrict the QEMU process from accessing processes and files on the host system that are not needed by QEMU.
- Resource limits and cgroup controllers provide throughput and utilization limits on key resources such as CPU time, memory, and I/O bandwidth.
- Linux namespaces can be used to make process, file system, and other system resources unavailable to QEMU. A namespaced QEMU process is restricted to only those resources that were granted to it.
- Linux seccomp is available via the QEMU `--sandbox` option. It disables system calls that are not needed by QEMU, thereby reducing the host kernel attack surface.

1.20.4 Sensitive configurations

There are aspects of QEMU that can have security implications which users & management applications must be aware of.

Monitor console (QMP and HMP)

The monitor console (whether used with QMP or HMP) provides an interface to dynamically control many aspects of QEMU's runtime operation. Many of the commands exposed will instruct QEMU to access content on the host file system and/or trigger spawning of external processes.

For example, the `migrate` command allows for the spawning of arbitrary processes for the purpose of tunnelling the migration data stream. The `blockdev-add` command instructs QEMU to open arbitrary files, exposing their content to the guest as a virtual disk.

Unless QEMU is otherwise confined using technologies such as SELinux, AppArmor, or Linux namespaces, the monitor console should be considered to have privileges equivalent to those of the user account QEMU is running under.

It is further important to consider the security of the character device backend over which the monitor console is exposed. It needs to have protection against malicious third parties which might try to make unauthorized connections, or perform man-in-the-middle attacks. Many of the character device backends do not satisfy this requirement and so must not be used for the monitor console.

The general recommendation is that the monitor console should be exposed over a UNIX domain socket backend to the local host only. Use of the TCP based character device backend is inappropriate unless configured to use both TLS encryption and authorization control policy on client connections.

In summary, the monitor console is considered a privileged control interface to QEMU and as such should only be made accessible to a trusted management application or user.

1.21 Multi-process QEMU

This document describes how to configure and use multi-process qemu. For the design document refer to `docs/devel/qemu-multiprocess`.

1.21.1 1) Configuration

multi-process is enabled by default for targets that enable KVM

1.21.2 2) Usage

Multi-process QEMU requires an orchestrator to launch.

Following is a description of command-line used to launch mpqemu.

- Orchestrator:
 - The Orchestrator creates a unix socketpair
 - It launches the remote process and passes one of the sockets to it via command-line.
 - It then launches QEMU and specifies the other socket as an option to the Proxy device object
- Remote Process:
 - QEMU can enter remote process mode by using the “remote” machine option.
 - The orchestrator creates a “remote-object” with details about the device and the file descriptor for the device
 - The remaining options are no different from how one launches QEMU with devices.
 - Example command-line for the remote process is as follows:


```
/usr/bin/qemu-system-x86_64 -machine x-remote -device lsi53c895a,id=lsi0
-drive id=drive_image2,file=/build/ol7-nvme-test-1.qcow2 -device scsi-
hd,id=drive2,drive=drive_image2,bus=lsi0.0,scsi-id=0 -object x-remote-
object,id=robj1,dev=lsi1,fd=4,
```
- QEMU:
 - Since parts of the RAM are shared between QEMU & remote process, a memory-backend-memfd is required to facilitate this, as follows:


```
-object memory-backend-memfd,id=mem,size=2G
```
 - A “x-pci-proxy-dev” device is created for each of the PCI devices emulated in the remote process. A “socket” sub-option specifies the other end of unix channel created by orchestrator. The “id” sub-option must be specified and should be the same as the “id” specified for the remote PCI device
 - Example commandline for QEMU is as follows:


```
-device x-pci-proxy-dev,id=lsi0,socket=3
```

1.22 Deprecated features

In general features are intended to be supported indefinitely once introduced into QEMU. In the event that a feature needs to be removed, it will be listed in this section. The feature will remain functional for the release in which it was deprecated and one further release. After these two releases, the feature is liable to be removed. Deprecated features may also generate warnings on the console when QEMU starts up, or if activated via a monitor command, however, this is not a mandatory requirement.

Prior to the 2.10.0 release there was no official policy on how long features would be deprecated prior to their removal, nor any documented list of which features were deprecated. Thus any features deprecated prior to 2.10.0 will be treated as if they were first deprecated in the 2.10.0 release.

What follows is a list of all features currently marked as deprecated.

1.22.1 System emulator command line arguments

`-usbdevice` (since 2.10.0)

The `-usbdevice DEV` argument is now a synonym for setting the `-device usb-DEV` argument instead. The deprecated syntax would automatically enable USB support on the machine type. If using the new syntax, USB support must be explicitly enabled via the `-machine usb=on` argument.

`-drive file=json:{...{'driver':'file'}}` (since 3.0)

The ‘file’ driver for drives is no longer appropriate for character or host devices and will only accept regular files (S_IFREG). The correct driver for these file types is ‘host_cdrom’ or ‘host_device’ as appropriate.

`-vnc acl` (since 4.0.0)

The `acl` option to the `-vnc` argument has been replaced by the `tls-authz` and `sasl-authz` options.

`QEMU_AUDIO_` environment variables and `-audio-help` (since 4.0)

The `-audiodev` argument is now the preferred way to specify audio backend settings instead of environment variables. To ease migration to the new format, the `-audiodev-help` option can be used to convert the current values of the environment variables to `-audiodev` options.

Creating sound card devices and vnc without `audiodev=` property (since 4.2)

When not using the deprecated legacy audio config, each sound card should specify an `audiodev=` property. Additionally, when using `vnc`, you should specify an `audiodev=` property if you plan to transmit audio through the VNC protocol.

Creating sound card devices using `-soundhw` (since 5.1)

Sound card devices should be created using `-device` instead. The names are the same for most devices. The exceptions are `hda` which needs two devices (`-device intel-hda -device hda-duplex`) and `pcspk` which can be activated using `-machine pcspk-audiodev=<name>`.

`-mon ...,control=readline,pretty=on|off` (since 4.1)

The `pretty=on|off` switch has no effect for HMP monitors, but is silently ignored. Using the switch with HMP monitors will become an error in the future.

RISC-V `-bios` (since 5.1)

QEMU 4.1 introduced support for the `-bios` option in QEMU for RISC-V for the RISC-V virt machine and sifive_u machine. QEMU 4.1 had no changes to the default behaviour to avoid breakages.

QEMU 5.1 changes the default behaviour from `-bios none` to `-bios default`.

QEMU 5.1 has three options:

1. **-bios default** - This is the current default behavior if no **-bios option** is included. This option will load the default OpenSBI firmware automatically. The firmware is included with the QEMU release and no user interaction is required. All a user needs to do is specify the kernel they want to boot with the **-kernel** option
2. **-bios none** - QEMU will not automatically load any firmware. It is up to the user to load all the images they need.
3. **-bios <file>** - Tells QEMU to load the specified file as the firmware.

Configuring floppies with `--global`

Use `-device floppy,...` instead:

```
-global isa-fdc.driveA=...
-global sysbus-fdc.driveA=...
-global SUNW,fdtwo.drive=...
```

become

```
-device floppy,unit=0,drive=...
```

and

```
-global isa-fdc.driveB=...
-global sysbus-fdc.driveB=...
```

become

```
-device floppy,unit=1,drive=...
```

`-drive` with bogus interface type

Drives with interface types other than `if=none` are for onboard devices. It is possible to use drives the board doesn't pick up with `-device`. This usage is now deprecated. Use `if=none` instead.

Short-form boolean options (since 6.0)

Boolean options such as `share=on/share=off` could be written in short form as `share` and `noshare`. This is now deprecated and will cause a warning.

`--enable-fips` (since 6.0)

This option restricts usage of certain cryptographic algorithms when the host is operating in FIPS mode.

If FIPS compliance is required, QEMU should be built with the `libgcrypt` library enabled as a cryptography provider.

Neither the `nettle` library, or the built-in cryptography provider are supported on FIPS enabled hosts.

1.22.2 QEMU Machine Protocol (QMP) commands

blockdev-open-tray, blockdev-close-tray argument device (since 2.8.0)

Use argument `id` instead.

eject argument device (since 2.8.0)

Use argument `id` instead.

blockdev-change-medium argument device (since 2.8.0)

Use argument `id` instead.

block_set_io_throttle argument device (since 2.8.0)

Use argument `id` instead.

migrate_set_downtime and **migrate_set_speed** (since 2.8.0)

Use `migrate-set-parameters` instead.

query-named-block-nodes result encryption_key_missing (since 2.10.0)

Always false.

query-block result inserted.encryption_key_missing (since 2.10.0)

Always false.

blockdev-add empty string argument backing (since 2.10.0)

Use argument value `null` instead.

migrate-set-cache-size and **query-migrate-cache-size** (since 2.11.0)

Use `migrate-set-parameters` and `query-migrate-parameters` instead.

block-commit arguments **base** and **top** (since 3.1.0)

Use arguments `base-node` and `top-node` instead.

object-add option props (since 5.0)

Specify the properties for the object as top-level arguments instead.

query-named-block-nodes and query-block result dirty-bitmaps[i].status (since 4.0)

The `status` field of the `BlockDirtyInfo` structure, returned by these commands is deprecated. Two new boolean fields, `recording` and `busy` effectively replace it.

query-block result field dirty-bitmaps (Since 4.2)

The `dirty-bitmaps` field of the `BlockInfo` structure, returned by the `query-block` command is itself now deprecated. The `dirty-bitmaps` field of the `BlockDeviceInfo` struct should be used instead, which is the type of the `inserted` field in `query-block` replies, as well as the type of array items in `query-named-block-nodes`.

Since the `dirty-bitmaps` field is optionally present in both the old and new locations, clients must use introspection to learn where to anticipate the field if/when it does appear in command output.

query-cpus (since 2.12.0)

The `query-cpus` command is replaced by the `query-cpus-fast` command.

query-cpus-fast arch output member (since 3.0.0)

The `arch` output member of the `query-cpus-fast` command is replaced by the `target` output member.

query-events (since 4.0)

The `query-events` command has been superseded by the more powerful and accurate `query-qmp-schema` command.

chardev client socket with wait option (since 4.0)

Character devices creating sockets in client mode should not specify the `'wait'` field, which is only applicable to sockets in server mode

nbd-server-add and nbd-server-remove (since 5.2)

Use the more generic commands `block-export-add` and `block-export-del` instead. As part of this deprecation, where `nbd-server-add` used a single bitmap, the new `block-export-add` uses a list of bitmaps.

1.22.3 Human Monitor Protocol (HMP) commands

acl_show, acl_reset, acl_policy, acl_add, acl_remove (since 4.0.0)

The `acl_show`, `acl_reset`, `acl_policy`, `acl_add`, and `acl_remove` commands are deprecated with no replacement. Authorization for VNC should be performed using the pluggable `QAuthZ` objects.

1.22.4 System emulator CPUS

moxie CPU (since 5.2.0)

The `moxie` guest CPU support is deprecated and will be removed in a future version of QEMU. It's unclear whether anybody is still using CPU emulation in QEMU, and there are no test images available to make sure that the code is still working.

lm32 CPUs (since 5.2.0)

The `lm32` guest CPU support is deprecated and will be removed in a future version of QEMU. The only public user of this architecture was the `milkymist` project, which has been dead for years; there was never an upstream Linux port.

unicore32 CPUs (since 5.2.0)

The `unicore32` guest CPU support is deprecated and will be removed in a future version of QEMU. Support for this CPU was removed from the upstream Linux kernel, and there is no available upstream toolchain to build binaries for it.

Icelake-Client CPU Model (since 5.2.0)

Icelake-Client CPU Models are deprecated. Use Icelake-Server CPU Models instead.

MIPS I7200 CPU Model (since 5.2)

The `I7200` guest CPU relies on the nanoMIPS ISA, which is deprecated (the ISA has never been upstreamed to a compiler toolchain). Therefore this CPU is also deprecated.

1.22.5 System emulator devices

ide-drive (since 4.2)

The `'ide-drive'` device is deprecated. Users should use `'ide-hd'` or `'ide-cd'` as appropriate to get an IDE hard disk or CD-ROM as needed.

scsi-disk (since 4.2)

The `'scsi-disk'` device is deprecated. Users should use `'scsi-hd'` or `'scsi-cd'` as appropriate to get a SCSI hard disk or CD-ROM as needed.

1.22.6 System emulator machines

Raspberry Pi `raspi2` and `raspi3` machines (since 5.2)

The Raspberry Pi machines come in various models (A, A+, B, B+). To be able to distinguish which model QEMU is implementing, the `raspi2` and `raspi3` machines have been renamed `raspi2b` and `raspi3b`.

1.22.7 Device options

Emulated device options

-device virtio-blk,scsi=on|off (since 5.0.0)

The virtio-blk SCSI passthrough feature is a legacy VIRTIO feature. VIRTIO 1.0 and later do not support it because the virtio-scsi device was introduced for full SCSI support. Use virtio-scsi instead when SCSI passthrough is required.

Note this also applies to `-device virtio-blk-pci,scsi=on|off`, which is an alias.

Block device options

"backing": "" (since 2.12.0)

In order to prevent QEMU from automatically opening an image's backing chain, use `"backing": null` instead.

rbid keyvalue pair encoded filenames: "" (since 3.1.0)

Options for `rbid` should be specified according to its runtime options, like other block drivers. Legacy parsing of keyvalue pair encoded filenames is useful to open images with the old format for backing files; These image files should be updated to use the current format.

Example of legacy encoding:

```
json:{"file.driver":"rbid", "file.filename":"rbid:rbid/name"}
```

The above, converted to the current supported format:

```
json:{"file.driver":"rbid", "file.pool":"rbid", "file.image":"name"}
```

sheepdog driver (since 5.2.0)

The `sheepdog` block device driver is deprecated. The corresponding upstream server project is no longer actively maintained. Users are recommended to switch to an alternative distributed block device driver such as RBD. The `qemu-img convert` command can be used to liberate existing data by moving it out of `sheepdog` volumes into an alternative storage backend.

1.22.8 linux-user mode CPUs

tilegx CPUs (since 5.1.0)

The `tilegx` guest CPU support (which was only implemented in linux-user mode) is deprecated and will be removed in a future version of QEMU. Support for this CPU was removed from the upstream Linux kernel in 2018, and has also been dropped from glibc.

ppc64abi32 CPUs (since 5.2.0)

The `ppc64abi32` architecture has a number of issues which regularly trip up our CI testing and is suspected to be quite broken. For that reason the maintainers strongly suspect no one actually uses it.

MIPS I7200 CPU (since 5.2)

The I7200 guest CPU relies on the nanoMIPS ISA, which is deprecated (the ISA has never been upstreamed to a compiler toolchain). Therefore this CPU is also deprecated.

1.22.9 Related binaries

qemu-img amend to adjust backing file (since 5.1)

The use of `qemu-img amend` to modify the name or format of a qcow2 backing image is deprecated; this functionality was never fully documented or tested, and interferes with other amend operations that need access to the original backing image (such as deciding whether a v3 zero cluster may be left unallocated when converting to a v2 image). Rather, any changes to the backing chain should be performed with `qemu-img rebase -u` either before or after the remaining changes being performed by amend, as appropriate.

qemu-img backing file without format (since 5.1)

The use of `qemu-img create`, `qemu-img rebase`, or `qemu-img convert` to create or modify an image that depends on a backing file now recommends that an explicit backing format be provided. This is for safety: if QEMU probes a different format than what you thought, the data presented to the guest will be corrupt; similarly, presenting a raw image to a guest allows a potential security exploit if a future probe sees a non-raw image based on guest writes.

To avoid the warning message, or even future refusal to create an unsafe image, you must pass `-o backing_fmt=` (or the shorthand `-F` during create) to specify the intended backing format. You may use `qemu-img rebase -u` to retroactively add a backing format to an existing image. However, be aware that there are already potential security risks to blindly using `qemu-img info` to probe the format of an untrusted backing image, when deciding what format to add into an existing image.

1.22.10 Backwards compatibility

Runnability guarantee of CPU models (since 4.1.0)

Previous versions of QEMU never changed existing CPU models in ways that introduced additional host software or hardware requirements to the VM. This allowed management software to safely change the machine type of an existing VM without introducing new requirements (“runnability guarantee”). This prevented CPU models from being updated to include CPU vulnerability mitigations, leaving guests vulnerable in the default configuration.

The CPU model runnability guarantee won’t apply anymore to existing CPU models. Management software that needs runnability guarantees must resolve the CPU model aliases using the `alias-of` field returned by the `query-cpu-definitions QMP` command.

While those guarantees are kept, the return value of `query-cpu-definitions` will have existing CPU model aliases point to a version that doesn’t break runnability guarantees (specifically, version 1 of those CPU models). In future QEMU versions, aliases will point to newer CPU model versions depending on the machine type, so management software must resolve CPU model aliases before starting a virtual machine.

1.22.11 Guest Emulator ISAs

nanoMIPS ISA

The nanoMIPS ISA has never been upstreamed to any compiler toolchain. As it is hard to generate binaries for it, declare it deprecated.

1.23 Removed features

What follows is a record of recently removed, formerly deprecated features that serves as a record for users who have encountered trouble after a recent upgrade.

1.23.1 System emulator command line arguments

-net . . . , name=*name* (removed in 5.1)

The `name` parameter of the `-net` option was a synonym for the `id` parameter, which should now be used instead.

-no-kvm (removed in 5.2)

The `-no-kvm` argument was a synonym for setting `-machine accel=tcg`.

-realtime (removed in 6.0)

The `-realtime mlock=on|off` argument has been replaced by the `-overcommit mem-lock=on|off` argument.

-show-cursor option (since 5.0)

Use `-display sdl,show-cursor=on`, `-display gtk,show-cursor=on` or `-display default,show-cursor=on` instead.

-tb-size option (removed in 6.0)

QEMU 5.0 introduced an alternative syntax to specify the size of the translation block cache, `-accel tcg, tb-size=`.

1.23.2 QEMU Machine Protocol (QMP) commands

block-dirty-bitmap-add “autoload” parameter (removed in 4.2.0)

The “autoload” parameter has been ignored since 2.12.0. All bitmaps are automatically loaded from qcow2 images.

cpu-add (removed in 5.2)

Use `device_add` for hotplugging vCPUs instead of `cpu-add`. See documentation of `query-hotpluggable-cpus` for additional details.

change (removed in 6.0)

Use `blockdev-change-medium` or `change-vnc-password` instead.

1.23.3 Human Monitor Protocol (HMP) commands

The `hub_id` parameter of `hostfwd_add` / `hostfwd_remove` (removed in 5.0)

The `[hub_id name]` parameter tuple of the ‘`hostfwd_add`’ and ‘`hostfwd_remove`’ HMP commands has been replaced by `netdev_id`.

cpu-add (removed in 5.2)

Use `device_add` for hotplugging vCPUs instead of `cpu-add`. See documentation of `query-hotpluggable-cpus` for additional details.

change vnc TARGET (removed in 6.0)

No replacement. The `change vnc password` and `change DEVICE MEDIUM` commands are not affected.

1.23.4 Guest Emulator ISAs

RISC-V ISA privilege specification version 1.09.1 (removed in 5.1)

The RISC-V ISA privilege specification version 1.09.1 has been removed. QEMU supports both the newer version 1.10.0 and the ratified version 1.11.0, these should be used instead of the 1.09.1 version.

1.23.5 System emulator CPUS

KVM guest support on 32-bit Arm hosts (removed in 5.2)

The Linux kernel has dropped support for allowing 32-bit Arm systems to host KVM guests as of the 5.7 kernel. Accordingly, QEMU is deprecating its support for this configuration and will remove it in a future version. Running 32-bit guests on a 64-bit Arm host remains supported.

RISC-V ISA Specific CPUs (removed in 5.1)

The RISC-V cpus with the ISA version in the CPU name have been removed. The four CPUs are: `rv32gcsu-v1.9.1`, `rv32gcsu-v1.10.0`, `rv64gcsu-v1.9.1` and `rv64gcsu-v1.10.0`. Instead the version can be specified via the CPU `priv_spec` option when using the `rv32` or `rv64` CPUs.

RISC-V no MMU CPUs (removed in 5.1)

The RISC-V no MMU cpus have been removed. The two CPUs: `rv32imacu-nommu` and `rv64imacu-nommu` can no longer be used. Instead the MMU status can be specified via the CPU `mmu` option when using the `rv32` or `rv64` CPUs.

1.23.6 System emulator machines

`spike_v1.9.1` and `spike_v1.10` (removed in 5.1)

The version specific Spike machines have been removed in favour of the generic `spike` machine. If you need to specify an older machine version of the RISC-V spec you can use the `-cpu rv64gcsu,priv_spec=v1.10.0` command line argument.

`mips_r4k` platform (removed in 5.2)

This machine type was very old and unmaintained. Users should use the `malta` machine type instead.

`mips_fulong2e` machine alias (removed in 6.0)

This machine has been renamed `fuloong2e`.

`pc-1.0`, `pc-1.1`, `pc-1.2` and `pc-1.3` (removed in 6.0)

These machine types were very old and likely could not be used for live migration from old QEMU versions anymore. Use a newer machine type instead.

1.23.7 Related binaries

`qemu-nbd --partition` (removed in 5.0)

The `qemu-nbd --partition $digit` code (also spelled `-P`) could only handle MBR partitions, and never correctly handled logical partitions beyond partition 5. Exporting a partition can still be done by utilizing the `--image-opts` option with a raw blockdev using the `offset` and `size` parameters layered on top of any other existing blockdev. For example, if partition 1 is 100MiB long starting at 1MiB, the old command:

```
qemu-nbd -t -P 1 -f qcow2 file.qcow2
```

can be rewritten as:

```
qemu-nbd -t --image-opts driver=raw,offset=1M,size=100M,file.driver=qcow2,file.file.
↪driver=file,file.file.filename=file.qcow2
```

`qemu-img convert -n -o` (removed in 5.1)

All options specified in `-o` are image creation options, so they are now rejected when used with `-n` to skip image creation.

`qemu-img create -b bad file $size` (removed in 5.1)

When creating an image with a backing file that could not be opened, `qemu-img create` used to issue a warning about the failure but proceed with the image creation if an explicit size was provided. However, as the `-u` option exists for this purpose, it is safer to enforce that any failure to open the backing image (including if the backing file is missing or an incorrect format was specified) is an error when `-u` is not used.

1.23.8 Command line options

-smp (invalid topologies) (removed 5.2)

CPU topology properties should describe whole machine topology including possible CPUs.

However, historically it was possible to start QEMU with an incorrect topology where $n \leq \text{sockets} * \text{cores} * \text{threads} < \text{maxcpus}$, which could lead to an incorrect topology enumeration by the guest. Support for invalid topologies is removed, the user must ensure topologies described with -smp include all possible cpus, i.e. $\text{sockets} * \text{cores} * \text{threads} = \text{maxcpus}$.

-numa node (without memory specified) (removed 5.2)

Splitting RAM by default between NUMA nodes had the same issues as mem parameter with the difference that the role of the user plays QEMU using implicit generic or board specific splitting rule. Use memdev with *memory-backend-ram* backend or mem (if it's supported by used machine type) to define mapping explicitly instead. Users of existing VMs, wishing to preserve the same RAM distribution, should configure it explicitly using -numa node, memdev options. Current RAM distribution can be retrieved using HMP command info numa and if separate memory devices (pclnv-dimm) are present use info memory-device and subtract device memory from output of info numa.

-numa node, mem=size (removed in 5.1)

The parameter mem of -numa node was used to assign a part of guest RAM to a NUMA node. But when using it, it's impossible to manage a specified RAM chunk on the host side (like bind it to a host node, setting bind policy, ...), so the guest ends up with the fake NUMA configuration with suboptimal performance. However since 2014 there is an alternative way to assign RAM to a NUMA node using parameter memdev, which does the same as mem and adds means to actually manage node RAM on the host side. Use parameter memdev with *memory-backend-ram* backend as replacement for parameter mem to achieve the same fake NUMA effect or a properly configured *memory-backend-file* backend to actually benefit from NUMA configuration. New machine versions (since 5.1) will not accept the option but it will still work with old machine types. User can check the QAPI schema to see if the legacy option is supported by looking at MachineInfo::numa-mem-supported property.

-mem-path fallback to RAM (removed in 5.0)

If guest RAM allocation from file pointed by mem-path failed, QEMU was falling back to allocating from RAM, which might have resulted in unpredictable behavior since the backing file specified by the user as ignored. Currently, users are responsible for making sure the backing storage specified with -mem-path can actually provide the guest RAM configured with -m and QEMU fails to start up if RAM allocation is unsuccessful.

-smp (invalid topologies) (removed 5.2)

CPU topology properties should describe whole machine topology including possible CPUs.

However, historically it was possible to start QEMU with an incorrect topology where $n \leq \text{sockets} * \text{cores} * \text{threads} < \text{maxcpus}$, which could lead to an incorrect topology enumeration by the guest. Support for invalid topologies is removed, the user must ensure topologies described with -smp include all possible cpus, i.e. $\text{sockets} * \text{cores} * \text{threads} = \text{maxcpus}$.

`-machine enforce-config-section=on|off` (removed 5.2)

The `enforce-config-section` property was replaced by the `-global migration.send-configuration={on|off}` option.

1.23.9 Block devices

VXHS backend (removed in 5.1)

The VXHS code did not compile since v2.12.0. It was removed in 5.1.

1.24 Supported build platforms

QEMU aims to support building and executing on multiple host OS platforms. This appendix outlines which platforms are the major build targets. These platforms are used as the basis for deciding upon the minimum required versions of 3rd party software QEMU depends on. The supported platforms are the targets for automated testing performed by the project when patches are submitted for review, and tested before and after merge.

If a platform is not listed here, it does not imply that QEMU won't work. If an unlisted platform has comparable software versions to a listed platform, there is every expectation that it will work. Bug reports are welcome for problems encountered on unlisted platforms unless they are clearly older vintage than what is described here.

Note that when considering software versions shipped in distros as support targets, QEMU considers only the version number, and assumes the features in that distro match the upstream release with the same version. In other words, if a distro backports extra features to the software in their distro, QEMU upstream code will not add explicit support for those backports, unless the feature is auto-detectable in a manner that works for the upstream releases too.

The [Repology](#) site is a useful resource to identify currently shipped versions of software in various operating systems, though it does not cover all distros listed below.

1.24.1 Linux OS, macOS, FreeBSD, NetBSD, OpenBSD

The project aims to support the most recent major version at all times. Support for the previous major version will be dropped 2 years after the new major version is released or when the vendor itself drops support, whichever comes first. In this context, third-party efforts to extend the lifetime of a distro are not considered, even when they are endorsed by the vendor (eg. Debian LTS).

For the purposes of identifying supported software versions available on Linux, the project will look at CentOS, Debian, Fedora, openSUSE, RHEL, SLES and Ubuntu LTS. Other distros will be assumed to ship similar software versions.

For FreeBSD and OpenBSD, decisions will be made based on the contents of the respective ports repository, while NetBSD will use the pkgsrc repository.

For macOS, [HomeBrew](#) will be used, although [MacPorts](#) is expected to carry similar versions.

1.24.2 Windows

The project supports building with current versions of the MinGW toolchain, hosted on Linux (Debian/Fedora).

The version of the Windows API that's currently targeted is Vista / Server 2008.

1.25 License

QEMU is a trademark of Fabrice Bellard.

QEMU is released under the [GNU General Public License](#), version 2. Parts of QEMU have specific licenses, see file [LICENSE](#).

QEMU User Mode Emulation User's Guide

This manual is the overall guide for users using QEMU for user-mode emulation. In this mode, QEMU can launch processes compiled for one CPU on another CPU.

Contents:

2.1 QEMU User space emulator

2.1.1 Supported Operating Systems

The following OS are supported in user space emulation:

- Linux (referred as qemu-linux-user)
- BSD (referred as qemu-bsd-user)

2.1.2 Features

QEMU user space emulation has the following notable features:

System call translation: QEMU includes a generic system call translator. This means that the parameters of the system calls can be converted to fix endianness and 32/64-bit mismatches between hosts and targets. IOCTLs can be converted too.

POSIX signal handling: QEMU can redirect to the running program all signals coming from the host (such as SIGALRM), as well as synthesize signals from virtual CPU exceptions (for example SIGFPE when the program executes a division by zero).

QEMU relies on the host kernel to emulate most signal system calls, for example to emulate the signal mask. On Linux, QEMU supports both normal and real-time signals.

Threading: On Linux, QEMU can emulate the `clone` syscall and create a real host thread (with a separate virtual CPU) for each emulated thread. Note that not all targets currently emulate atomic operations correctly. x86 and Arm use a global lock in order to preserve their semantics.

QEMU was conceived so that ultimately it can emulate itself. Although it is not very useful, it is an important test to show the power of the emulator.

2.1.3 Linux User space emulator

Command line options

```
qemu-i386 [-h] [-d] [-L path] [-s size] [-cpu model] [-g port] [-B offset] [-R size]   
↳program [arguments...]
```

-h Print the help

-L path Set the x86 elf interpreter prefix (default=/usr/local/qemu-i386)

-s size Set the x86 stack size in bytes (default=524288)

-cpu model Select CPU model (-cpu help for list and additional feature selection)

-E var=value Set environment var to value.

-U var Remove var from the environment.

-B offset Offset guest address by the specified number of bytes. This is useful when the address region required by guest applications is reserved on the host. This option is currently only supported on some hosts.

-R size Pre-allocate a guest virtual address space of the given size (in bytes). “G”, “M”, and “k” suffixes may be used when specifying the size.

Debug options:

-d item1,... Activate logging of the specified items (use ‘-d help’ for a list of log items)

-p pagesize Act as if the host page size was ‘pagesize’ bytes

-g port Wait gdb connection to port

-singlestep Run the emulation in single step mode.

Environment variables:

QEMU_STRACE Print system calls and arguments similar to the ‘strace’ program (NOTE: the actual ‘strace’ program will not work because the user space emulator hasn’t implemented ptrace). At the moment this is incomplete. All system calls that don’t have a specific argument format are printed with information for six arguments. Many flag-style arguments don’t have decoders and will show up as numbers.

Other binaries

- user mode (Alpha)
 - qemu-alpha TODO.
- user mode (Arm)
 - qemu-armeb TODO.
 - qemu-arm is also capable of running Arm “Angel” semihosted ELF binaries (as implemented by the arm-elf and arm-eabi Newlib/GDB configurations), and arm-uclinux bFLT format binaries.
- user mode (ColdFire)
- user mode (M68K)

- `qemu-m68k` is capable of running semihosted binaries using the BDM (`m5xxx-ram-hosted.ld`) or `m68k-sim` (`sim.ld`) syscall interfaces, and coldfire uClinux bFLT format binaries.

The binary format is detected automatically.

- user mode (Cris)
 - `qemu-cris` TODO.
- user mode (i386)
 - `qemu-i386` TODO.
 - `qemu-x86_64` TODO.
- user mode (Microblaze)
 - `qemu-microblaze` TODO.
- user mode (MIPS)
 - `qemu-mips` executes 32-bit big endian MIPS binaries (MIPS O32 ABI).
 - `qemu-mipsel` executes 32-bit little endian MIPS binaries (MIPS O32 ABI).
 - `qemu-mips64` executes 64-bit big endian MIPS binaries (MIPS N64 ABI).
 - `qemu-mips64el` executes 64-bit little endian MIPS binaries (MIPS N64 ABI).
 - `qemu-mipsn32` executes 32-bit big endian MIPS binaries (MIPS N32 ABI).
 - `qemu-mipsn32el` executes 32-bit little endian MIPS binaries (MIPS N32 ABI).
- user mode (NiosII)
 - `qemu-nios2` TODO.
- user mode (PowerPC)
 - `qemu-ppc64abi32` TODO.
 - `qemu-ppc64` TODO.
 - `qemu-ppc` TODO.
- user mode (SH4)
 - `qemu-sh4eb` TODO.
 - `qemu-sh4` TODO.
- user mode (SPARC)
 - `qemu-sparc` can execute Sparc32 binaries (Sparc32 CPU, 32 bit ABI).
 - `qemu-sparc32plus` can execute Sparc32 and SPARC32PLUS binaries (Sparc64 CPU, 32 bit ABI).
 - `qemu-sparc64` can execute some Sparc64 (Sparc64 CPU, 64 bit ABI) and SPARC32PLUS binaries (Sparc64 CPU, 32 bit ABI).

2.1.4 BSD User space emulator

BSD Status

- target Sparc64 on Sparc64: Some trivial programs work.

Quick Start

In order to launch a BSD process, QEMU needs the process executable itself and all the target dynamic libraries used by it.

- On Sparc64, you can just try to launch any process by using the native libraries:

```
qemu-sparc64 /bin/ls
```

Command line options

```
qemu-sparc64 [-h] [-d] [-L path] [-s size] [-bsd type] program [arguments...]
```

-h Print the help

-L path Set the library root path (default=)

-s size Set the stack size in bytes (default=524288)

-ignore-environment Start with an empty environment. Without this option, the initial environment is a copy of the caller's environment.

-E var=value Set environment var to value.

-U var Remove var from the environment.

-bsd type Set the type of the emulated BSD Operating system. Valid values are FreeBSD, NetBSD and OpenBSD (default).

Debug options:

-d item1,... Activate logging of the specified items (use '-d help' for a list of log items)

-p pagesize Act as if the host page size was 'pagesize' bytes

-singlestep Run the emulation in single step mode.

Contents:

3.1 QEMU disk image utility

3.1.1 Synopsis

qemu-img [*standard options*] *command* [*command options*]

3.1.2 Description

qemu-img allows you to create, convert and modify images offline. It can handle all image formats supported by QEMU.

Warning: Never use qemu-img to modify images in use by a running virtual machine or any other process; this may destroy the image. Also, be aware that querying an image that is being modified by another process may encounter inconsistent state.

3.1.3 Options

Standard options:

-h, --help

Display this help and exit

-V, --version

Display version information and exit

-T, --trace [[enable=]PATTERN] [,events=FILE] [,file=FILE]

Specify tracing options.

[enable=]PATTERN

Immediately enable events matching *PATTERN* (either event name or a globbing pattern). This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend. To specify multiple events or patterns, specify the `-trace` option multiple times.

Use `-trace help` to print a list of names of trace points.

`events=FILE`

Immediately enable events listed in *FILE*. The file must contain one event name (as listed in the `trace-events-all` file) per line; globbing patterns are accepted too. This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend.

`file=FILE`

Log output traces to *FILE*. This option is only available if QEMU has been compiled with the `simple` tracing backend.

The following commands are supported:

amend [--object OBJECTDEF] [--image-opts] [-p] [-q] [-f FMT] [-t CACHE] [--force] -o OPTION
bench [-c COUNT] [-d DEPTH] [-f FMT] [--flush-interval=FLUSH_INTERVAL] [-i AIO] [-n] [--no-
bitmap (--merge SOURCE | --add | --remove | --clear | --enable | --disable)... [-b SOURCE_I
check [--object OBJECTDEF] [--image-opts] [-q] [-f FMT] [--output=OFMT] [-r [leaks | all]]
commit [--object OBJECTDEF] [--image-opts] [-q] [-f FMT] [-t CACHE] [-b BASE] [-r RATE_LIM
compare [--object OBJECTDEF] [--image-opts] [-f FMT] [-F FMT] [-T SRC_CACHE] [-p] [-q] [-s
convert [--object OBJECTDEF] [--image-opts] [--target-image-opts] [--target-is-zero] [--bit
create [--object OBJECTDEF] [-q] [-f FMT] [-b BACKING_FILE] [-F BACKING_FMT] [-u] [-o OPTIO
dd [--image-opts] [-U] [-f FMT] [-O OUTPUT_FMT] [bs=BLOCK_SIZE] [count=BLOCKS] [skip=BLOCK
info [--object OBJECTDEF] [--image-opts] [-f FMT] [--output=OFMT] [--backing-chain] [-U] F
map [--object OBJECTDEF] [--image-opts] [-f FMT] [--start-offset=OFFSET] [--max-length=LEN
measure [--output=OFMT] [-O OUTPUT_FMT] [-o OPTIONS] [--size N | [--object OBJECTDEF] [--in
snapshot [--object OBJECTDEF] [--image-opts] [-U] [-q] [-l | -a SNAPSHOT | -c SNAPSHOT | -c
rebase [--object OBJECTDEF] [--image-opts] [-U] [-q] [-f FMT] [-t CACHE] [-T SRC_CACHE] [-p
resize [--object OBJECTDEF] [--image-opts] [-f FMT] [--preallocation=PREALLOC] [-q] [--shr

Command parameters:

FILENAME is a disk image filename.

FMT is the disk image format. It is guessed automatically in most cases. See below for a description of the supported disk formats.

SIZE is the disk image size in bytes. Optional suffixes `k` or `K` (kilobyte, 1024) `M` (megabyte, 1024k) and `G` (gigabyte, 1024M) and `T` (terabyte, 1024G) are supported. `b` is ignored.

OUTPUT_FILENAME is the destination disk image filename.

OUTPUT_FMT is the destination format.

OPTIONS is a comma separated list of format specific options in a `name=value` format. Use `-o ?` for an overview of the options supported by the used format or see the format descriptions below for details.

SNAPSHOT_PARAM is param used for internal snapshot, format is 'snapshot.id=[ID],snapshot.name=[NAME]' or '[ID_OR_NAME]'.

--object OBJECTDEF

is a QEMU user creatable object definition. See the *qemu(1)* manual page for a description of the object properties. The most common object type is a *secret*, which is used to supply passwords and/or encryption keys.

--image-opts

Indicates that the source *FILENAME* parameter is to be interpreted as a full option string, not a plain filename. This parameter is mutually exclusive with the *-f* parameter.

--target-image-opts

Indicates that the *OUTPUT_FILENAME* parameter(s) are to be interpreted as a full option string, not a plain filename. This parameter is mutually exclusive with the *-O* parameters. It is currently required to also use the *-n* parameter to skip image creation. This restriction may be relaxed in a future release.

--force-share (-U)

If specified, *qemu-img* will open the image in shared mode, allowing other QEMU processes to open it in write mode. For example, this can be used to get the image information (with 'info' subcommand) when the image is used by a running guest. Note that this could produce inconsistent results because of concurrent metadata changes, etc. This option is only allowed when opening images in read-only mode.

--backing-chain

Will enumerate information about backing files in a disk image chain. Refer below for further description.

-c

Indicates that target image must be compressed (qcow format only).

-h

With or without a command, shows help and lists the supported formats.

-p

Display progress bar (compare, convert and rebase commands only). If the *-p* option is not used for a command that supports it, the progress is reported when the process receives a *SIGUSR1* or *SIGINFO* signal.

-q

Quiet mode - do not print any output (except errors). There's no progress bar in case both *-q* and *-p* options are used.

-S SIZE

Indicates the consecutive number of bytes that must contain only zeros for *qemu-img* to create a sparse image during conversion. This value is rounded down to the nearest 512 bytes. You may use the common size suffixes like *k* for kilobytes.

-t CACHE

Specifies the cache mode that should be used with the (destination) file. See the documentation of the emulator's *-drive cache=...* option for allowed values.

-T SRC_CACHE

Specifies the cache mode that should be used with the source file(s). See the documentation of the emulator's *-drive cache=...* option for allowed values.

Parameters to compare subcommand:

-f

First image format

-F

Second image format

-s

Strict mode - fail on different image size or sector allocation

Parameters to convert subcommand:

--bitmaps

Additionally copy all persistent bitmaps from the top layer of the source

-n

Skip the creation of the target volume

-m

Number of parallel coroutines for the convert process

-W

Allow out-of-order writes to the destination. This option improves performance, but is only recommended for preallocated devices like host devices or other raw block devices.

-C

Try to use copy offloading to move data from source image to target. This may improve performance if the data is remote, such as with NFS or iSCSI backends, but will not automatically sparsify zero sectors, and may result in a fully allocated target image depending on the host support for getting allocation information.

-r

Rate limit for the convert process

--salvage

Try to ignore I/O errors when reading. Unless in quiet mode (**-q**), errors will still be printed. Areas that cannot be read from the source will be treated as containing only zeroes.

--target-is-zero

Assume that reading the destination image will always return zeros. This parameter is mutually exclusive with a destination image that has a backing file. It is required to also use the **-n** parameter to skip image creation.

Parameters to dd subcommand:

bs=BLOCK_SIZE

Defines the block size

count=BLOCKS

Sets the number of input blocks to copy

if=INPUT

Sets the input file

of=OUTPUT

Sets the output file

skip=BLOCKS

Sets the number of input blocks to skip

Parameters to snapshot subcommand:

snapshot

Is the name of the snapshot to create, apply or delete

-a

Applies a snapshot (revert disk to saved state)

-c

Creates a snapshot

-d

Deletes a snapshot

-l

Lists all snapshots in the given image

Command description:

amend [--object OBJECTDEF] [--image-opts] [-p] [-q] [-f FMT] [-t CACHE] [--force] -o OPTION
Amends the image format specific *OPTIONS* for the image file *FILENAME*. Not all file formats support this operation.

The set of options that can be amended are dependent on the image format, but note that amending the backing chain relationship should instead be performed with `qemu-img rebase`.

--force allows some unsafe operations. Currently for -f luks, it allows to erase the last encryption key, and to overwrite an active encryption key.

bench [-c COUNT] [-d DEPTH] [-f FMT] [--flush-interval=FLUSH_INTERVAL] [-i AIO] [-n] [--no-
Run a simple sequential I/O benchmark on the specified image. If -w is specified, a write test is performed, otherwise a read test is performed.

A total number of *COUNT* I/O requests is performed, each *BUFFER_SIZE* bytes in size, and with *DEPTH* requests in parallel. The first request starts at the position given by *OFFSET*, each following request increases the current position by *STEP_SIZE*. If *STEP_SIZE* is not given, *BUFFER_SIZE* is used for its value.

If *FLUSH_INTERVAL* is specified for a write test, the request queue is drained and a flush is issued before new writes are made whenever the number of remaining requests is a multiple of *FLUSH_INTERVAL*. If additionally --no-drain is specified, a flush is issued without draining the request queue first.

if -i is specified, *AIO* option can be used to specify different AIO backends: threads, native or io_uring.

If -n is specified, the native AIO backend is used if possible. On Linux, this option only works if -t none or -t directsync is specified as well.

For write tests, by default a buffer filled with zeros is written. This can be overridden with a pattern byte specified by *PATTERN*.

bitmap (--merge SOURCE | --add | --remove | --clear | --enable | --disable)... [-b SOURCE_
Perform one or more modifications of the persistent bitmap *BITMAP* in the disk image *FILENAME*. The various modifications are:

--add to create *BITMAP*, enabled to record future edits.

--remove to remove *BITMAP*.

--clear to clear *BITMAP*.

--enable to change *BITMAP* to start recording future edits.

--disable to change *BITMAP* to stop recording future edits.

--merge to merge the contents of the *SOURCE* bitmap into *BITMAP*.

Additional options include -g which sets a non-default *GRANULARITY* for --add, and -b and -F which select an alternative source file for all *SOURCE* bitmaps used by --merge.

To see what bitmaps are present in an image, use `qemu-img info`.

check [--object OBJECTDEF] [--image-opts] [-q] [-f FMT] [--output=OFMT] [-r [leaks | all]]
Perform a consistency check on the disk image *FILENAME*. The command can output in the format *OFMT* which is either human or json. The JSON output is an object of QAPI type ImageCheck.

If -r is specified, qemu-img tries to repair any inconsistencies found during the check. -r leaks repairs only cluster leaks, whereas -r all fixes all kinds of errors, with a higher risk of choosing the wrong fix or hiding corruption that has already occurred.

Only the formats qcow2, qed and vdi support consistency checks.

In case the image does not have any inconsistencies, check exits with 0. Other exit codes indicate the kind of inconsistency found or if another error occurred. The following table summarizes all exit codes of the check subcommand:

- 0 Check completed, the image is (now) consistent
- 1 Check not completed because of internal errors
- 2 Check completed, image is corrupted
- 3 Check completed, image has leaked clusters, but is not corrupted
- 63 Checks are not supported by the image format

If `-r` is specified, exit codes representing the image state refer to the state after (the attempt at) repairing it. That is, a successful `-r all` will yield the exit code 0, independently of the image state before.

commit [--object OBJECTDEF] [--image-opts] [-q] [-f FMT] [-t CACHE] [-b BASE] [-r RATE_LIM]

Commit the changes recorded in *FILENAME* in its base image or backing file. If the backing file is smaller than the snapshot, then the backing file will be resized to be the same size as the snapshot. If the snapshot is smaller than the backing file, the backing file will not be truncated. If you want the backing file to match the size of the smaller snapshot, you can safely truncate it yourself once the commit operation successfully completes.

The image *FILENAME* is emptied after the operation has succeeded. If you do not need *FILENAME* afterwards and intend to drop it, you may skip emptying *FILENAME* by specifying the `-d` flag.

If the backing chain of the given image file *FILENAME* has more than one layer, the backing file into which the changes will be committed may be specified as *BASE* (which has to be part of *FILENAME*'s backing chain). If *BASE* is not specified, the immediate backing file of the top image (which is *FILENAME*) will be used. Note that after a commit operation all images between *BASE* and the top image will be invalid and may return garbage data when read. For this reason, `-b` implies `-d` (so that the top image stays valid).

The rate limit for the commit process is specified by `-r`.

compare [--object OBJECTDEF] [--image-opts] [-f FMT] [-F FMT] [-T SRC_CACHE] [-p] [-q] [-s]

Check if two images have the same content. You can compare images with different format or settings.

The format is probed unless you specify it by `-f` (used for *FILENAME1*) and/or `-F` (used for *FILENAME2*) option.

By default, images with different size are considered identical if the larger image contains only unallocated and/or zeroed sectors in the area after the end of the other image. In addition, if any sector is not allocated in one image and contains only zero bytes in the second one, it is evaluated as equal. You can use Strict mode by specifying the `-s` option. When compare runs in Strict mode, it fails in case image size differs or a sector is allocated in one image and is not allocated in the second one.

By default, compare prints out a result message. This message displays information that both images are same or the position of the first different byte. In addition, result message can report different image size in case Strict mode is used.

Compare exits with 0 in case the images are equal and with 1 in case the images differ. Other exit codes mean an error occurred during execution and standard error output should contain an error message. The following table summarizes all exit codes of the compare subcommand:

- 0 Images are identical
- 1 Images differ
- 2 Error on opening an image
- 3 Error on checking a sector allocation
- 4 Error on reading data

convert [--object OBJECTDEF] [--image-opts] [--target-image-opts] [--target-is-zero] [--bit

Convert the disk image *FILENAME* or a snapshot *SNAPSHOT_PARAM* to disk image *OUTPUT_FILENAME* using format *OUTPUT_FMT*. It can be optionally compressed (*-c* option) or use any format specific options like encryption (*-o* option).

Only the formats *qcow* and *qcow2* support compression. The compression is read-only. It means that if a compressed sector is rewritten, then it is rewritten as uncompressed data.

Image conversion is also useful to get smaller image when using a growable format such as *qcow*: the empty sectors are detected and suppressed from the destination image.

SPARSE_SIZE indicates the consecutive number of bytes (defaults to 4k) that must contain only zeros for *qemu-img* to create a sparse image during conversion. If *SPARSE_SIZE* is 0, the source will not be scanned for unallocated or zero sectors, and the destination image will always be fully allocated.

You can use the *BACKING_FILE* option to force the output image to be created as a copy on write image of the specified base image; the *BACKING_FILE* should have the same content as the input's base image, however the path, image format, etc may differ.

If a relative path name is given, the backing file is looked up relative to the directory containing *OUTPUT_FILENAME*.

If the *-n* option is specified, the target volume creation will be skipped. This is useful for formats such as *raw* if the target volume has already been created with site specific options that cannot be supplied through *qemu-img*.

Out of order writes can be enabled with *-W* to improve performance. This is only recommended for preallocated devices like host devices or other raw block devices. Out of order write does not work in combination with creating compressed images.

NUM_COROUTINES specifies how many coroutines work in parallel during the convert process (defaults to 8).

create [--object OBJECTDEF] [-q] [-f FMT] [-b BACKING_FILE] [-F BACKING_FMT] [-u] [-o OPTIO

Create the new disk image *FILENAME* of size *SIZE* and format *FMT*. Depending on the file format, you can add one or more *OPTIONS* that enable additional features of this format.

If the option *BACKING_FILE* is specified, then the image will record only the differences from *BACKING_FILE*. No size needs to be specified in this case. *BACKING_FILE* will never be modified unless you use the *commit* monitor command (or *qemu-img commit*).

If a relative path name is given, the backing file is looked up relative to the directory containing *FILENAME*.

Note that a given backing file will be opened to check that it is valid. Use the *-u* option to enable unsafe backing file mode, which means that the image will be created even if the associated backing file cannot be opened. A matching backing file must be created or additional options be used to make the backing file specification valid when you want to use an image created this way.

The size can also be specified using the *SIZE* option with *-o*, it doesn't need to be specified separately in this case.

dd [--image-opts] [-U] [-f FMT] [-O OUTPUT_FMT] [bs=BLOCK_SIZE] [count=BLOCKS] [skip=BLOCKS]
dd copies from *INPUT* file to *OUTPUT* file converting it from *FMT* format to *OUTPUT_FMT* format.

The data is by default read and written using blocks of 512 bytes but can be modified by specifying *BLOCK_SIZE*. If *count=BLOCKS* is specified *dd* will stop reading input after reading *BLOCKS* input blocks.

The size syntax is similar to *dd(1)*'s size syntax.

info [--object OBJECTDEF] [--image-opts] [-f FMT] [--output=OFMT] [--backing-chain] [-U] F

Give information about the disk image *FILENAME*. Use it in particular to know the size reserved on disk which can be different from the displayed size. If VM snapshots are stored in the disk image, they are displayed too.

If a disk image has a backing file chain, information about each disk image in the chain can be recursively enumerated by using the option *--backing-chain*.

For instance, if you have an image chain like:

```
base.qcow2 <- snap1.qcow2 <- snap2.qcow2
```

To enumerate information about each disk image in the above chain, starting from top to base, do:

```
qemu-img info --backing-chain snap2.qcow2
```

The command can output in the format *OFMT* which is either `human` or `json`. The JSON output is an object of QAPI type `ImageInfo`; with `--backing-chain`, it is an array of `ImageInfo` objects.

`--output=human` reports the following information (for every image in the chain):

image The image file name

file format The image format

virtual size The size of the guest disk

disk size How much space the image file occupies on the host file system (may be shown as 0 if this information is unavailable, e.g. because there is no file system)

cluster_size Cluster size of the image format, if applicable

encrypted Whether the image is encrypted (only present if so)

cleanly shut down This is shown as `no` if the image is dirty and will have to be auto-repaired the next time it is opened in qemu.

backing file The backing file name, if present

backing file format The format of the backing file, if the image enforces it

Snapshot list A list of all internal snapshots

Format specific information Further information whose structure depends on the image format. This section is a textual representation of the respective `ImageInfoSpecific*` QAPI object (e.g. `ImageInfoSpecificQCow2` for qcow2 images).

map [`--object OBJECTDEF`] [`--image-opts`] [`-f FMT`] [`--start-offset=OFFSET`] [`--max-length=LEN`]
Dump the metadata of image *FILENAME* and its backing file chain. In particular, this commands dumps the allocation state of every sector of *FILENAME*, together with the topmost file that allocates it in the backing file chain.

Two option formats are possible. The default format (`human`) only dumps known-nonzero areas of the file. Known-zero parts of the file are omitted altogether, and likewise for parts that are not allocated throughout the chain. `qemu-img` output will identify a file from where the data can be read, and the offset in the file. Each line will include four fields, the first three of which are hexadecimal numbers. For example the first line of:

Offset	Length	Mapped to	File
0	0x20000	0x50000	/tmp/overlay.qcow2
0x100000	0x10000	0x95380000	/tmp/backing.qcow2

means that 0x20000 (131072) bytes starting at offset 0 in the image are available in `/tmp/overlay.qcow2` (opened in `raw` format) starting at offset 0x50000 (327680). Data that is compressed, encrypted, or otherwise not available in raw format will cause an error if `human` format is in use. Note that file names can include newlines, thus it is not safe to parse this output format in scripts.

The alternative format `json` will return an array of dictionaries in JSON format. It will include similar information in the `start`, `length`, `offset` fields; it will also include other more specific information:

- whether the sectors contain actual data or not (boolean field `data`; if false, the sectors are either unallocated or stored as optimized all-zero clusters);

- whether the data is known to read as zero (boolean field `zero`);
- in order to make the output shorter, the target file is expressed as a `depth`; for example, a depth of 2 refers to the backing file of the backing file of *FILENAME*.

In JSON format, the `offset` field is optional; it is absent in cases where human format would omit the entry or exit with an error. If `data` is false and the `offset` field is present, the corresponding sectors in the file are not yet in use, but they are preallocated.

For more information, consult `include/block/block.h` in QEMU's source code.

measure [--output=OFMT] [-O OUTPUT_FMT] [-o OPTIONS] [--size N | [--object OBJECTDEF] [--in

Calculate the file size required for a new image. This information can be used to size logical volumes or SAN LUNs appropriately for the image that will be placed in them. The values reported are guaranteed to be large enough to fit the image. The command can output in the format *OFMT* which is either `human` or `json`. The JSON output is an object of QAPI type `BlockMeasureInfo`.

If the size *N* is given then act as if creating a new empty image file using `qemu-img create`. If *FILENAME* is given then act as if converting an existing image file using `qemu-img convert`. The format of the new file is given by *OUTPUT_FMT* while the format of an existing file is given by *FMT*.

A snapshot in an existing image can be specified using *SNAPSHOT_PARAM*.

The following fields are reported:

```
required size: 524288
fully allocated size: 1074069504
bitmaps size: 0
```

The `required size` is the file size of the new image. It may be smaller than the virtual disk size if the image format supports compact representation.

The `fully allocated size` is the file size of the new image once data has been written to all sectors. This is the maximum size that the image file can occupy with the exception of internal snapshots, dirty bitmaps, vmstate data, and other advanced image format features.

The `bitmaps size` is the additional size required in order to copy bitmaps from a source image in addition to the guest-visible data; the line is omitted if either source or destination lacks bitmap support, or 0 if bitmaps are supported but there is nothing to copy.

snapshot [--object OBJECTDEF] [--image-opts] [-U] [-q] [-l | -a SNAPSHOT | -c SNAPSHOT | -o

List, apply, create or delete snapshots in image *FILENAME*.

rebase [--object OBJECTDEF] [--image-opts] [-U] [-q] [-f FMT] [-t CACHE] [-T SRC_CACHE] [-p

Changes the backing file of an image. Only the formats `qcow2` and `qed` support changing the backing file.

The backing file is changed to *BACKING_FILE* and (if the image format of *FILENAME* supports this) the backing file format is changed to *BACKING_FMT*. If *BACKING_FILE* is specified as "" (the empty string), then the image is rebased onto no backing file (i.e. it will exist independently of any backing file).

If a relative path name is given, the backing file is looked up relative to the directory containing *FILENAME*.

CACHE specifies the cache mode to be used for *FILENAME*, whereas *SRC_CACHE* specifies the cache mode for reading backing files.

There are two different modes in which `rebase` can operate:

Safe mode This is the default mode and performs a real rebase operation. The new backing file may differ from the old one and `qemu-img rebase` will take care of keeping the guest-visible content of *FILENAME* unchanged.

In order to achieve this, any clusters that differ between *BACKING_FILE* and the old backing file of *FILENAME* are merged into *FILENAME* before actually changing the backing file.

Note that the safe mode is an expensive operation, comparable to converting an image. It only works if the old backing file still exists.

Unsafe mode `qemu-img` uses the unsafe mode if `-u` is specified. In this mode, only the backing file name and format of *FILENAME* is changed without any checks on the file contents. The user must take care of specifying the correct new backing file, or the guest-visible content of the image will be corrupted.

This mode is useful for renaming or moving the backing file to somewhere else. It can be used without an accessible old backing file, i.e. you can use it to fix an image whose backing file has already been moved/renamed.

You can use `rebase` to perform a “diff” operation on two disk images. This can be useful when you have copied or cloned a guest, and you want to get back to a thin image on top of a template or base image.

Say that `base.img` has been cloned as `modified.img` by copying it, and that the `modified.img` guest has run so there are now some changes compared to `base.img`. To construct a thin image called `diff.qcow2` that contains just the differences, do:

```
qemu-img create -f qcow2 -b modified.img diff.qcow2
qemu-img rebase -b base.img diff.qcow2
```

At this point, `modified.img` can be discarded, since `base.img + diff.qcow2` contains the same information.

resize [--object OBJECTDEF] [--image-opts] [-f FMT] [--preallocation=PREALLOC] [-q] [--shrink] *SIZE*
Change the disk image as if it had been created with *SIZE*.

Before using this command to shrink a disk image, you **MUST** use file system and partitioning tools inside the VM to reduce allocated file systems and partition sizes accordingly. Failure to do so will result in data loss!

When shrinking images, the `--shrink` option must be given. This informs `qemu-img` that the user acknowledges all loss of data beyond the truncated image’s end.

After using this command to grow a disk image, you must use file system and partitioning tools inside the VM to actually begin using the new space on the device.

When growing an image, the `--preallocation` option may be used to specify how the additional image area should be allocated on the host. See the format description in the [Notes](#) section which values are allowed. Using this option may result in slightly more data being allocated than necessary.

3.1.4 Notes

Supported image file formats:

`raw`

Raw disk image format (default). This format has the advantage of being simple and easily exportable to all other emulators. If your file system supports *holes* (for example in ext2 or ext3 on Linux or NTFS on Windows), then only the written sectors will reserve space. Use `qemu-img info` to know the real size used by the image or `ls -ls` on Unix/Linux.

Supported options:

preallocation Preallocation mode (allowed values: `off`, `falloc`, `full`). `falloc` mode pre-allocates space for image by calling `posix_fallocate()`. `full` mode preallocates space for image by writing data to underlying storage. This data may or may not be zero, depending on the storage location.

`qcow2`

QEMU image format, the most versatile format. Use it to have smaller images (useful if your filesystem does not supports holes, for example on Windows), optional AES encryption, zlib based compression and support of multiple VM snapshots.

Supported options:

compat Determines the qcow2 version to use. `compat=0.10` uses the traditional image format that can be read by any QEMU since 0.10. `compat=1.1` enables image format extensions that only QEMU 1.1 and newer understand (this is the default). Amongst others, this includes zero clusters, which allow efficient copy-on-read for sparse images.

backing_file File name of a base image (see `create` subcommand)

backing_fmt Image format of the base image

encryption If this option is set to `on`, the image is encrypted with 128-bit AES-CBC.

The use of encryption in qcow and qcow2 images is considered to be flawed by modern cryptography standards, suffering from a number of design problems:

- The AES-CBC cipher is used with predictable initialization vectors based on the sector number. This makes it vulnerable to chosen plaintext attacks which can reveal the existence of encrypted data.
- The user passphrase is directly used as the encryption key. A poorly chosen or short passphrase will compromise the security of the encryption.
- In the event of the passphrase being compromised there is no way to change the passphrase to protect data in any qcow images. The files must be cloned, using a different encryption passphrase in the new file. The original file must then be securely erased using a program like `shred`, though even this is ineffective with many modern storage technologies.
- Initialization vectors used to encrypt sectors are based on the guest virtual sector number, instead of the host physical sector. When a disk image has multiple internal snapshots this means that data in multiple physical sectors is encrypted with the same initialization vector. With the CBC mode, this opens the possibility of watermarking attacks if the attack can collect multiple sectors encrypted with the same IV and some predictable data. Having multiple qcow2 images with the same passphrase also exposes this weakness since the passphrase is directly used as the key.

Use of qcow / qcow2 encryption is thus strongly discouraged. Users are recommended to use an alternative encryption technology such as the Linux `dm-crypt` / LUKS system.

cluster_size Changes the qcow2 cluster size (must be between 512 and 2M). Smaller cluster sizes can improve the image file size whereas larger cluster sizes generally provide better performance.

preallocation Preallocation mode (allowed values: `off`, `metadata`, `falloc`, `full`). An image with preallocated metadata is initially larger but can improve performance when the image needs to grow. `falloc` and `full` preallocations are like the same options of `raw` format, but sets up metadata also.

lazy_refcounts If this option is set to `on`, reference count updates are postponed with the goal of avoiding metadata I/O and improving performance. This is particularly interesting with `cache=writethrough` which doesn't batch metadata updates. The tradeoff is that after a host crash, the reference count tables must be rebuilt, i.e. on the next open an (automatic) `qemu-img check -r all` is required, which may take some time.

This option can only be enabled if `compat=1.1` is specified.

nocow If this option is set to `on`, it will turn off COW of the file. It's only valid on `btarfs`, no effect on other file systems.

Btrfs has low performance when hosting a VM image file, even more when the guest on the VM also using btrfs as file system. Turning off COW is a way to mitigate this bad performance. Generally there are two ways to turn off COW on btrfs:

- Disable it by mounting with `nodatacow`, then all newly created files will be NOCOW
- For an empty file, add the NOCOW file attribute. That's what this option does.

Note: this option is only valid to new or empty files. If there is an existing file which is COW and has data blocks already, it couldn't be changed to NOCOW by setting `nocow=on`. One can issue `lsattr filename` to check if the NOCOW flag is set or not (Capital 'C' is NOCOW flag).

Other

QEMU also supports various other image file formats for compatibility with older QEMU versions or other hypervisors, including VMDK, VDI, VHD (vpc), VHDX, qcow1 and QED. For a full list of supported formats see `qemu-img --help`. For a more detailed description of these formats, see the QEMU block drivers reference documentation.

The main purpose of the block drivers for these formats is image conversion. For running VMs, it is recommended to convert the disk images to either raw or qcow2 in order to achieve good performance.

3.2 QEMU Storage Daemon

3.2.1 Synopsis

qemu-storage-daemon [options]

3.2.2 Description

`qemu-storage-daemon` provides disk image functionality from QEMU, `qemu-img`, and `qemu-nbd` in a long-running process controlled via QMP commands without running a virtual machine. It can export disk images, run block job operations, and perform other disk-related operations. The daemon is controlled via a QMP monitor and initial configuration from the command-line.

The daemon offers the following subset of QEMU features:

- Block nodes
- Block jobs
- Block exports
- Throttle groups
- Character devices
- Crypto and secrets
- QMP
- IOThreads

Commands can be sent over a QEMU Monitor Protocol (QMP) connection. See the *qemu-storage-daemon-qmp-ref(7)* manual page for a description of the commands.

The daemon runs until it is stopped using the `quit` QMP command or `SIGINT`/`SIGHUP`/`SIGTERM`.

Warning: Never modify images in use by a running virtual machine or any other process; this may destroy the image. Also, be aware that querying an image that is being modified by another process may encounter inconsistent state.

3.2.3 Options

Standard options:

-h, --help

Display help and exit

-V, --version

Display version information and exit

-T, --trace [[enable=]PATTERN] [,events=FILE] [,file=FILE]

Specify tracing options.

[enable=]PATTERN

Immediately enable events matching *PATTERN* (either event name or a globbing pattern). This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend. To specify multiple events or patterns, specify the `-trace` option multiple times.

Use `-trace help` to print a list of names of trace points.

events=FILE

Immediately enable events listed in *FILE*. The file must contain one event name (as listed in the `trace-events-all` file) per line; globbing patterns are accepted too. This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend.

file=FILE

Log output traces to *FILE*. This option is only available if QEMU has been compiled with the `simple` tracing backend.

--blockdev BLOCKDEVDEF

is a block node definition. See the *qemu* (1) manual page for a description of block node properties and the *qemu-block-drivers* (7) manual page for a description of driver-specific parameters.

--chardev CHARDEVDEF

is a character device definition. See the *qemu* (1) manual page for a description of character device properties. A common character device definition configures a UNIX domain socket:

```
--chardev socket,id=char1,path=/tmp/qmp.sock,server,nowait
```

--export [type=]nbd,id=<id>,node-name=<node-name>[,name=<export-name>][,writable=on|off][,bitmap=]

--export [type=]vhost-user-blk,id=<id>,node-name=<node-name>,addr.type=unix,addr.path=<socket-path>[,writable=]

--export [type=]vhost-user-blk,id=<id>,node-name=<node-name>,addr.type=fd,addr.str=<fd>[,writable=]

is a block export definition. `node-name` is the block node that should be exported. `writable` determines whether or not the export allows write requests for modifying data (the default is off).

The `nbd` export type requires `--nbd-server` (see below). `name` is the NBD export name. `bitmap` is the name of a dirty bitmap reachable from the block node, so the NBD client can use `NBD_OPT_SET_META_CONTEXT` with the metadata context name “`qemu:dirty-bitmap:BITMAP`” to inspect the bitmap.

The `vhost-user-blk` export type takes a `vhost-user` socket address on which it accept incoming connections. Both `addr.type=unix,addr.path=<socket-path>` for UNIX domain sockets and `addr.type=fd,addr.str=<fd>` for file descriptor passing are supported. `logical-block-size` sets the logical block size in bytes (the default is 512). `num-queues` sets the number of virtqueues (the default is 1).

--monitor MONITORDEF

is a QMP monitor definition. See the *qemu* (1) manual page for a description of QMP monitor properties. A common QMP monitor definition configures a monitor on character device `char1`:

```
--monitor chardev=char1
```

--nbd-server *addr.type*=inet,*addr.host*=<host>,*addr.port*=<port>[,*tls-creds*=<id>][,*tls-authz*=<id>]
--nbd-server *addr.type*=unix,*addr.path*=<path>[,*tls-creds*=<id>][,*tls-authz*=<id>][,*max-connect*=<id>]
is a server for NBD exports. Both TCP and UNIX domain sockets are supported. TLS encryption can be configured using **--object** *tls-creds*-* and *authz*-* secrets (see below).

To configure an NBD server on UNIX domain socket path `/tmp/nbd.sock`:

```
--nbd-server addr.type=unix,addr.path=/tmp/nbd.sock
```

--object help

--object <type>[,help]

--object <type>[,<property>=<value>...]

is a QEMU user creatable object definition. List object types with help. List object properties with <type>[,help]. See the *qemu* (1) manual page for a description of the object properties.

3.2.4 Examples

Launch the daemon with QMP monitor socket `qmp.sock` so clients can execute QMP commands:

```
$ qemu-storage-daemon \  
  --chardev socket,path=qmp.sock,server,nowait,id=char1 \  
  --monitor chardev=char1
```

Export raw image file `disk.img` over NBD UNIX domain socket `nbd.sock`:

```
$ qemu-storage-daemon \  
  --blockdev driver=file,node-name=disk,filename=disk.img \  
  --nbd-server addr.type=unix,addr.path=nbd.sock \  
  --export type=nbd,id=export,node-name=disk,writable=on
```

Export a qcow2 image file `disk.qcow2` as a `vhosts-user-blk` device over UNIX domain socket `vhost-user-blk.sock`:

```
$ qemu-storage-daemon \  
  --blockdev driver=file,node-name=file,filename=disk.qcow2 \  
  --blockdev driver=qcow2,node-name=qcow2,file=file \  
  --export type=vhost-user-blk,id=export,addr.type=unix,addr.path=vhost-user-blk.  
↪sock,node-name=qcow2
```

3.2.5 See also

qemu (1), *qemu-block-drivers* (7), *qemu-storage-daemon-qmp-ref* (7)

3.3 QEMU Disk Network Block Device Server

3.3.1 Synopsis

qemu-nbd [*OPTION*]... *filename*

qemu-nbd -L [*OPTION*]...

qemu-nbd -d dev

3.3.2 Description

Export a QEMU disk image using the NBD protocol.

Other uses:

- Bind a /dev/nbdX block device to a QEMU server (on Linux).
- As a client to query exports of a remote NBD server.

3.3.3 Options

filename is a disk image filename, or a set of block driver options if `--image-opts` is specified.

dev is an NBD device.

--object *type*,*id*=*ID*,...*props*...

Define a new instance of the *type* object class identified by *ID*. See the *qemu (1)* manual page for full details of the properties supported. The common object types that it makes sense to define are the `secret` object, which is used to supply passwords and/or encryption keys, and the `tls-creds` object, which is used to supply TLS credentials for the `qemu-nbd` server or client.

-p, --port=*PORT*

TCP port to listen on as a server, or connect to as a client (default 10809).

-o, --offset=*OFFSET*

The offset into the image.

-b, --bind=*IFACE*

The interface to bind to as a server, or connect to as a client (default 0.0.0.0).

-k, --socket=*PATH*

Use a unix socket with path *PATH*.

--image-opts

Treat *filename* as a set of image options, instead of a plain filename. If this flag is specified, the `-f` flag should not be used, instead the `format=` option should be set.

-f, --format=*FMT*

Force the use of the block driver for format *FMT* instead of auto-detecting.

-r, --read-only

Export the disk as read-only.

-A, --allocation-depth

Expose allocation depth information via the `qemu:allocation-depth` metadata context accessible through `NBD_OPT_SET_META_CONTEXT`.

-B, --bitmap=*NAME*

If *filename* has a qcow2 persistent bitmap *NAME*, expose that bitmap via the `qemu:dirty-bitmap:NAME` metadata context accessible through `NBD_OPT_SET_META_CONTEXT`.

-s, --snapshot

Use *filename* as an external snapshot, create a temporary file with `backing_file=filename`, redirect the write to the temporary one.

- l, --load-snapshot=SNAPSHOT_PARAM**
Load an internal snapshot inside *filename* and export it as an read-only device, SNAPSHOT_PARAM format is `snapshot.id=[ID], snapshot.name=[NAME]` or `[ID_OR_NAME]`
- cache=CACHE**
The cache mode to be used with the file. See the documentation of the emulator's `-drive cache=...` option for allowed values.
- n, --nocache**
Equivalent to `--cache=none`.
- aio=AIO**
Set the asynchronous I/O mode between threads (the default), `native` (Linux only), and `io_uring` (Linux 5.1+).
- discard=DISCARD**
Control whether `discard` (also known as `trim` or `unmap`) requests are ignored or passed to the filesystem. *DISCARD* is one of `ignore` (or `off`), `unmap` (or `on`). The default is `ignore`.
- detect-zeroes=DETECT_ZEROES**
Control the automatic conversion of plain zero writes by the OS to driver-specific optimized zero write commands. *DETECT_ZEROES* is one of `off`, `on`, or `unmap`. `unmap` converts a zero write to an unmap operation and can only be used if *DISCARD* is set to `unmap`. The default is `off`.
- c, --connect=DEV**
Connect *filename* to NBD device *DEV* (Linux only).
- d, --disconnect**
Disconnect the device *DEV* (Linux only).
- e, --shared=NUM**
Allow up to *NUM* clients to share the device (default 1), 0 for unlimited. Safe for readers, but for now, consistency is not guaranteed between multiple writers.
- t, --persistent**
Don't exit on the last connection.
- x, --export-name=NAME**
Set the NBD volume export name (default of a zero-length string).
- D, --description=DESCRIPTION**
Set the NBD volume export description, as a human-readable string.
- L, --list**
Connect as a client and list all details about the exports exposed by a remote NBD server. This enables list mode, and is incompatible with options that change behavior related to a specific export (such as `--export-name`, `--offset`, ...).
- tls-creds=ID**
Enable mandatory TLS encryption for the server by setting the ID of the TLS credentials object previously created with the `-object` option; or provide the credentials needed for connecting as a client in list mode.
- fork**
Fork off the server process and exit the parent once the server is running.
- pid-file=PATH**
Store the server's process ID in the given file.
- tls-authz=ID**
Specify the ID of a `qauthz` object previously created with the `--object` option. This will be used to authorize connecting users against their x509 distinguished name.

- v, --verbose**
Display extra debugging information.
- h, --help**
Display this help and exit.
- V, --version**
Display version information and exit.
- T, --trace** [[enable=]PATTERN] [,events=FILE] [,file=FILE]
Specify tracing options.
- [enable=]PATTERN

Immediately enable events matching *PATTERN* (either event name or a globbing pattern). This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend. To specify multiple events or patterns, specify the `-trace` option multiple times.

Use `-trace help` to print a list of names of trace points.

events=FILE

Immediately enable events listed in *FILE*. The file must contain one event name (as listed in the `trace-events-all` file) per line; globbing patterns are accepted too. This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend.

file=FILE

Log output traces to *FILE*. This option is only available if QEMU has been compiled with the `simple` tracing backend.

3.3.4 Examples

Start a server listening on port 10809 that exposes only the guest-visible contents of a qcow2 file, with no TLS encryption, and with the default export name (an empty string). The command is one-shot, and will block until the first successful client disconnects:

```
qemu-nbd -f qcow2 file.qcow2
```

Start a long-running server listening with encryption on port 10810, and whitelist clients with a specific X.509 certificate to connect to a 1 megabyte subset of a raw file, using the export name 'subset':

```
qemu-nbd \
--object tls-creds-x509,id=tls0,endpoint=server,dir=/path/to/qemutls \
--object 'authz-simple,id=auth0,identity=CN=laptop.example.com,,\
O=Example Org,,L=London,,ST=London,,C=GB' \
--tls-creds tls0 --tls-authz auth0 \
-t -x subset -p 10810 \
--image-opts driver=raw,offset=1M,size=1M,file.driver=file,file.filename=file.raw
```

Serve a read-only copy of a guest image over a Unix socket with as many as 5 simultaneous readers, with a persistent process forked as a daemon:

```
qemu-nbd --fork --persistent --shared=5 --socket=/path/to/sock \
--read-only --format=qcow2 file.qcow2
```

Expose the guest-visible contents of a qcow2 file via a block device `/dev/nbd0` (and possibly creating `/dev/nbd0p1` and friends for partitions found within), then disconnect the device when done. Access to bind `qemu-nbd` to an `/dev/nbd` device generally requires root privileges, and may also require the execution of `modprobe nbd` to enable the kernel

NBD client module. *CAUTION:* Do not use this method to mount filesystems from an untrusted guest image - a malicious guest may have prepared the image to attempt to trigger kernel bugs in partition probing or file system mounting.

```
qemu-nbd -c /dev/nbd0 -f qcow2 file.qcow2
qemu-nbd -d /dev/nbd0
```

Query a remote server to see details about what export(s) it is serving on port 10809, and authenticating via PSK:

```
qemu-nbd \
  --object tls-creds-psk,id=tls0,dir=/tmp/keys,username=eblake,endpoint=client \
  --tls-creds tls0 -L -b remote.example.com
```

3.3.5 See also

qemu (1), *qemu-img (1)*

3.4 QEMU persistent reservation helper

3.4.1 Synopsis

qemu-pr-helper [*OPTION*]

3.4.2 Description

Implements the persistent reservation helper for QEMU.

SCSI persistent reservations allow restricting access to block devices to specific initiators in a shared storage setup. When implementing clustering of virtual machines, it is a common requirement for virtual machines to send persistent reservation SCSI commands. However, the operating system restricts sending these commands to unprivileged programs because incorrect usage can disrupt regular operation of the storage fabric. QEMU's SCSI passthrough devices `scsi-block` and `scsi-generic` support passing guest persistent reservation requests to a privileged external helper program. **qemu-pr-helper** is that external helper; it creates a socket which QEMU can connect to to communicate with it.

If you want to run VMs in a setup like this, this helper should be started as a system service, and you should read the QEMU manual section on “persistent reservation managers” to find out how to configure QEMU to connect to the socket created by **qemu-pr-helper**.

After connecting to the socket, **qemu-pr-helper** can optionally drop root privileges, except for those capabilities that are needed for its operation.

qemu-pr-helper can also use the systemd socket activation protocol. In this case, the systemd socket unit should specify a Unix stream socket, like this:

```
[Socket]
ListenStream=/var/run/qemu-pr-helper.sock
```

3.4.3 Options

-d, **--daemon**
run in the background (and create a PID file)

- q, --quiet**
decrease verbosity
- v, --verbose**
increase verbosity
- f, --pidfile=PATH**
PID file when running as a daemon. By default the PID file is created in the system runtime state directory, for example `/var/run/qemu-pr-helper.pid`.
- k, --socket=PATH**
path to the socket. By default the socket is created in the system runtime state directory, for example `/var/run/qemu-pr-helper.sock`.
- T, --trace** `[[enable=]PATTERN] [,events=FILE] [,file=FILE]`
Specify tracing options.
- `[enable=]PATTERN`
Immediately enable events matching *PATTERN* (either event name or a globbing pattern). This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend. To specify multiple events or patterns, specify the `-trace` option multiple times.
- Use `-trace help` to print a list of names of trace points.
- `events=FILE`
Immediately enable events listed in *FILE*. The file must contain one event name (as listed in the `trace-events-all` file) per line; globbing patterns are accepted too. This option is only available if QEMU has been compiled with the `simple`, `log` or `ftrace` tracing backend.
- `file=FILE`
Log output traces to *FILE*. This option is only available if QEMU has been compiled with the `simple` tracing backend.
- u, --user=USER**
user to drop privileges to
- g, --group=GROUP**
group to drop privileges to
- h, --help**
Display a help message and exit.
- V, --version**
Display version information and exit.

3.5 QEMU SystemTap trace tool

3.5.1 Synopsis

`qemu-trace-stap` `[GLOBAL-OPTIONS]` *COMMAND* `[COMMAND-OPTIONS]` *ARGS...*

3.5.2 Description

The `qemu-trace-stap` program facilitates tracing of the execution of QEMU emulators using SystemTap.

It is required to have the SystemTap runtime environment installed to use this program, since it is a wrapper around execution of the `stap` program.

3.5.3 Options

The following global options may be used regardless of which command is executed:

--verbose, -v

Display verbose information about command execution.

The following commands are valid:

list *BINARY* *PATTERN*...

List all the probe names provided by *BINARY* that match *PATTERN*.

If *BINARY* is not an absolute path, it will be located by searching the directories listed in the `$PATH` environment variable.

PATTERN is a plain string that is used to filter the results of this command. It may optionally contain a `*` wildcard to facilitate matching multiple probes without listing each one explicitly. Multiple *PATTERN* arguments may be given, causing listing of probes that match any of the listed names. If no *PATTERN* is given, the all possible probes will be listed.

For example, to list all probes available in the `qemu-system-x86_64` binary:

```
$ qemu-trace-stap list qemu-system-x86_64
```

To filter the list to only cover probes related to QEMU's cryptographic subsystem, in a binary outside `$PATH`

```
$ qemu-trace-stap list /opt/qemu/4.0.0/bin/qemu-system-x86_64 'qcrypto*'
```

run *OPTIONS* *BINARY* *PATTERN*...

Run a trace session, printing formatted output any time a process that is executing *BINARY* triggers a probe matching *PATTERN*.

If *BINARY* is not an absolute path, it will be located by searching the directories listed in the `$PATH` environment variable.

PATTERN is a plain string that matches a probe name shown by the *LIST* command. It may optionally contain a `*` wildcard to facilitate matching multiple probes without listing each one explicitly. Multiple *PATTERN* arguments may be given, causing all matching probes to be monitored. At least one *PATTERN* is required, since `stap` is not capable of tracing all known QEMU probes concurrently without overflowing its trace buffer.

Invocation of this command does not need to be synchronized with invocation of the QEMU process(es). It will match probes on all existing running processes and all future launched processes, unless told to only monitor a specific process.

Valid command specific options are:

--pid=PID, -p PID

Restrict the tracing session so that it only triggers for the process identified by *PID*.

For example, to monitor all processes executing `qemu-system-x86_64` as found on `$PATH`, displaying all I/O related probes:

```
$ qemu-trace-stap run qemu-system-x86_64 'qio*'
```

To monitor only the QEMU process with PID 1732

```
$ qemu-trace-stap run --pid=1732 qemu-system-x86_64 'qio*'

```

To monitor QEMU processes running an alternative binary outside of \$PATH, displaying verbose information about setup of the tracing environment:

```
$ qemu-trace-stap -v run /opt/qemu/4.0.0/qemu-system-x86_64 'qio*'

```

3.5.4 See also

qemu(1), *stap(1)*

3.6 QEMU 9p virtfs proxy filesystem helper

3.6.1 Synopsis

virtfs-proxy-helper [*OPTIONS*]

3.6.2 Description

Pass-through security model in QEMU 9p server needs root privilege to do few file operations (like chown, chmod to any mode/uid:gid). There are two issues in pass-through security model:

- TOCTTOU vulnerability: Following symbolic links in the server could provide access to files beyond 9p export path.
- Running QEMU with root privilege could be a security issue.

To overcome above issues, following approach is used: A new filesystem type ‘proxy’ is introduced. Proxy FS uses chroot + socket combination for securing the vulnerability known with following symbolic links. Intention of adding a new filesystem type is to allow qemu to run in non-root mode, but doing privileged operations using socket IO.

Proxy helper (a stand alone binary part of qemu) is invoked with root privileges. Proxy helper chroots into 9p export path and creates a socket pair or a named socket based on the command line parameter. QEMU and proxy helper communicate using this socket. QEMU proxy fs driver sends filesystem request to proxy helper and receives the response from it.

The proxy helper is designed so that it can drop root privileges except for the capabilities needed for doing filesystem operations.

3.6.3 Options

The following options are supported:

-h

Display help and exit

-p, --path PATH

Path to export for proxy filesystem driver

-f, --fd SOCKET_ID

Use given file descriptor as socket descriptor for communicating with qemu proxy fs driver. Usually a helper like libvirt will create socketpair and pass one of the fds as parameter to this option.

- s, --socket** SOCKET_FILE
Creates named socket file for communicating with qemu proxy fs driver
- u, --uid** UID
uid to give access to named socket file; used in combination with -g.
- g, --gid** GID
gid to give access to named socket file; used in combination with -u.
- n, --nodaemon**
Run as a normal program. By default program will run in daemon mode

3.7 QEMU virtio-fs shared file system daemon

3.7.1 Synopsis

virtiofsd [*OPTIONS*]

3.7.2 Description

Share a host directory tree with a guest through a virtio-fs device. This program is a vhost-user backend that implements the virtio-fs device. Each virtio-fs device instance requires its own virtiofsd process.

This program is designed to work with QEMU's `--device vhost-user-fs-pci` but should work with any virtual machine monitor (VMM) that supports vhost-user. See the Examples section below.

This program must be run as the root user. The program drops privileges where possible during startup although it must be able to create and access files with any uid/gid:

- The ability to invoke syscalls is limited using `seccomp(2)`.
- Linux capabilities(7) are dropped.

In “namespace” sandbox mode the program switches into a new file system namespace and invokes `pivot_root(2)` to make the shared directory tree its root. A new pid and net namespace is also created to isolate the process.

In “chroot” sandbox mode the program invokes `chroot(2)` to make the shared directory tree its root. This mode is intended for container environments where the container runtime has already set up the namespaces and the program does not have permission to create namespaces itself.

Both sandbox modes prevent “file system escapes” due to symlinks and other file system objects that might lead to files outside the shared directory.

3.7.3 Options

- h, --help**
Print help.
- V, --version**
Print version.
- d**
Enable debug output.
- syslog**
Print log messages to syslog instead of stderr.

-o OPTION

- **debug** - Enable debug output.
- **flockno_flock** - Enable/disable flock. The default is `no_flock`.
- **modcaps=CAPLIST** Modify the list of capabilities allowed; CAPLIST is a colon separated list of capabilities, each preceded by either + or -, e.g. `“+sys_admin:-chown”`.
- **log_level=LEVEL** - Print only log messages matching LEVEL or more severe. LEVEL is one of `err`, `warn`, `info`, or `debug`. The default is `info`.
- **posix_lockno_posix_lock** - Enable/disable remote POSIX locks. The default is `no_posix_lock`.
- **readdirplusno_readdirplus** - Enable/disable readdirplus. The default is `readdirplus`.
- **sandbox=namespace|chroot** - Sandbox mode: - **namespace**: Create mount, pid, and net namespaces and `pivot_root(2)` into the shared directory. - **chroot**: `chroot(2)` into shared directory (use in containers). The default is `“namespace”`.
- **source=PATH** - Share host directory tree located at PATH. This option is required.
- **timeout=TIMEOUT** - I/O timeout in seconds. The default depends on `cache` option.
- **writebackno_writeback** - Enable/disable writeback cache. The cache allows the FUSE client to buffer and merge write requests. The default is `no_writeback`.
- **xattrno_xattr** - Enable/disable extended attributes (xattr) on files and directories. The default is `no_xattr`.

--socket-path=PATH

Listen on vhost-user UNIX domain socket at PATH.

--socket-group=GROUP

Set the vhost-user UNIX domain socket gid to GROUP.

--fd=FDNUM

Accept connections from vhost-user UNIX domain socket file descriptor FDNUM. The file descriptor must already be listening for connections.

--thread-pool-size=NUM

Restrict the number of worker threads per request queue to NUM. The default is 64.

--cache=none|auto|always

Select the desired trade-off between coherency and performance. `none` forbids the FUSE client from caching to achieve best coherency at the cost of performance. `auto` acts similar to NFS with a 1 second metadata cache timeout. `always` sets a long cache lifetime at the expense of coherency. The default is `auto`.

3.7.4 xattr-mapping

By default the name of xattr's used by the client are passed through to the server file system. This can be a problem where either those xattr names are used by something on the server (e.g. selinux client/server confusion) or if the `virtiofsd` is running in a container with restricted privileges where it cannot access some attributes.

A mapping of xattr names can be made using `-o xattrmap=mapping` where the `mapping` string consists of a series of rules.

The first matching rule terminates the mapping. The set of rules must include a terminating rule to match any remaining attributes at the end.

Each rule consists of a number of fields separated with a separator that is the first non-white space character in the rule. This separator must then be used for the whole rule. White space may be added before and after each rule.

Using ‘:’ as the separator a rule is of the form:

```
:type:scope:key:prepend:
```

scope is:

- ‘client’ - match ‘key’ against a xattr name from the client for setxattr/getxattr/removexattr
- ‘server’ - match ‘prepend’ against a xattr name from the server for listxattr
- ‘all’ - can be used to make a single rule where both the server and client matches are triggered.

type is one of:

- ‘prefix’ - is designed to prepend and strip a prefix; the modified attributes then being passed on to the client/server.
- ‘ok’ - Causes the rule set to be terminated when a match is found while allowing matching xattr’s through unchanged. It is intended both as a way of explicitly terminating the list of rules, and to allow some xattr’s to skip following rules.
- ‘bad’ - If a client tries to use a name matching ‘key’ it’s denied using EPERM; when the server passes an attribute name matching ‘prepend’ it’s hidden. In many ways it’s use is very like ‘ok’ as either an explicit terminator or for special handling of certain patterns.

key is a string tested as a prefix on an attribute name originating on the client. It maybe empty in which case a ‘client’ rule will always match on client names.

prepend is a string tested as a prefix on an attribute name originating on the server, and used as a new prefix. It may be empty in which case a ‘server’ rule will always match on all names from the server.

e.g.:

```
:prefix:client:trusted.:user.virtiofs.:
```

will match ‘trusted.’ attributes in client calls and prefix them before passing them to the server.

```
:prefix:server::user.virtiofs.:
```

will strip ‘user.virtiofs.’ from all server replies.

```
:prefix:all:trusted.:user.virtiofs.:
```

combines the previous two cases into a single rule.

```
:ok:client:user.:
```

will allow get/set xattr for ‘user.’ xattr’s and ignore following rules.

```
:ok:server::security.:
```

will pass ‘security.’ xattr’s in listxattr from the server and ignore following rules.

```
:ok:all:::
```

will terminate the rule search passing any remaining attributes in both directions.

```
:bad:server::security.:
```

would hide ‘security.’ xattr’s in listxattr from the server.

A simpler ‘map’ type provides a shorter syntax for the common case:

```
:map:key:prepend:
```

The ‘map’ type adds a number of separate rules to add **prepend** as a prefix to the matched **key** (or all attributes if **key** is empty). There may be at most one ‘map’ rule and it must be the last rule in the set.

3.7.5 xattr-mapping Examples

- 1) Prefix all attributes with ‘user.virtiofs.’

```
-o xattrmap=":prefix:all::user.virtiofs.:bad:all:::"
```

This uses two rules, using `:` as the field separator; the first rule prefixes and strips ‘user.virtiofs.’, the second rule hides any non-prefixed attributes that the host set.

This is equivalent to the ‘map’ rule:

```
:: -o xattrmap=":map::user.virtiofs.:"
```

- 2) Prefix ‘trusted.’ attributes, allow others through

```
"/prefix/all/trusted./user.virtiofs./
/bad/server//trusted./
/bad/client/user.virtiofs./
/ok/all///"
```

Here there are four rules, using `/` as the field separator, and also demonstrating that new lines can be included between rules. The first rule is the prefixing of ‘trusted.’ and stripping of ‘user.virtiofs.’. The second rule hides unprefixed ‘trusted.’ attributes on the host. The third rule stops a guest from explicitly setting the ‘user.virtiofs.’ path directly. Finally, the fourth rule lets all remaining attributes through.

This is equivalent to the ‘map’ rule:

```
:: -o xattrmap=":map/trusted./user.virtiofs./"
```

- 3) Hide ‘security.’ attributes, and allow everything else

```
"/bad/all/security./security./
/ok/all///"
```

The first rule combines what could be separate client and server rules into a single ‘all’ rule, matching ‘security.’ in either client arguments or lists returned from the host. This stops the client seeing any ‘security.’ attributes on the server and stops it setting any.

3.7.6 Examples

Export `/var/lib/fs/vm001/` on vhost-user UNIX domain socket `/var/run/vm001-vhost-fs.sock`:

```
host# virtiofsd --socket-path=/var/run/vm001-vhost-fs.sock -o source=/var/lib/fs/vm001
host# qemu-system-x86_64 \
    -chardev socket,id=char0,path=/var/run/vm001-vhost-fs.sock \
    -device vhost-user-fs-pci,chardev=char0,tag=myfs \
    -object memory-backend-memfd,id=mem,size=4G,share=on \
    -numa node,memdev=mem \
    ...
guest# mount -t virtiofs myfs /mnt
```

QEMU System Emulation Management and Interoperability Guide

This manual contains documents and specifications that are useful for making QEMU interoperate with other software.

Contents:

4.1 Dirty Bitmaps and Incremental Backup

Dirty Bitmaps are in-memory objects that track writes to block devices. They can be used in conjunction with various block job operations to perform incremental or differential backup regimens.

This document explains the conceptual mechanisms, as well as up-to-date, complete and comprehensive documentation on the API to manipulate them. (Hopefully, the “why”, “what”, and “how”.)

The intended audience for this document is developers who are adding QEMU backup features to management applications, or power users who run and administer QEMU directly via QMP.

Contents

- *Dirty Bitmaps and Incremental Backup*
 - *Overview*
 - *Supported Image Formats*
 - *Dirty Bitmap Names*
 - *Bitmap Status*
 - *Basic QMP Usage*
 - * *Supported Commands*
 - * *Creation: block-dirty-bitmap-add*
 - * *Deletion: block-dirty-bitmap-remove*

- * *Resetting: block-dirty-bitmap-clear*
- * *Enabling: block-dirty-bitmap-enable*
- * *Enabling: block-dirty-bitmap-disable*
- * *Merging, Copying: block-dirty-bitmap-merge*
- * *Querying: query-block*
- *Bitmap Persistence*
- *Transactions*
 - * *Justification*
 - * *Supported Bitmap Transactions*
- *Incremental Backups - Push Model*
 - * *Example: New Incremental Backup Anchor Point*
 - * *Example: Resetting an Incremental Backup Anchor Point*
 - * *Example: First Incremental Backup*
 - * *Example: Second Incremental Backup*
 - * *Example: Incremental Push Backups without Backing Files*
 - * *Example: Multi-drive Incremental Backup*
- *Push Backup Errors & Recovery*
 - * *Example: Individual Failures*
 - * *Example: Partial Transactional Failures*
 - * *Example: Grouped Completion Mode*

4.1.1 Overview

Bitmaps are bit vectors where each ‘1’ bit in the vector indicates a modified (“dirty”) segment of the corresponding block device. The size of the segment that is tracked is the granularity of the bitmap. If the granularity of a bitmap is 64K, each ‘1’ bit means that a 64K region as a whole may have changed in some way, possibly by as little as one byte.

Smaller granularities mean more accurate tracking of modified disk data, but requires more computational overhead and larger bitmap sizes. Larger granularities mean smaller bitmap sizes, but less targeted backups.

The size of a bitmap (in bytes) can be computed as such: $\text{size} = \text{ceil}(\text{ceil}(\text{image_size} / \text{granularity}) / 8)$

e.g. the size of a 64KiB granularity bitmap on a 2TiB image is:

$$\text{size} = ((2147483648\text{K} / 64\text{K}) / 8) = 4194304\text{B} = 4\text{MiB}.$$

QEMU uses these bitmaps when making incremental backups to know which sections of the file to copy out. They are not enabled by default and must be explicitly added in order to begin tracking writes.

Bitmaps can be created at any time and can be attached to any arbitrary block node in the storage graph, but are most useful conceptually when attached to the root node attached to the guest’s storage device model.

That is to say: It’s likely most useful to track the guest’s writes to disk, but you could theoretically track things like qcow2 metadata changes by attaching the bitmap elsewhere in the storage graph. This is beyond the scope of this document.

QEMU supports persisting these bitmaps to disk via the qcow2 image format. Bitmaps which are stored or loaded in this way are called “persistent”, whereas bitmaps that are not are called “transient”.

QEMU also supports the migration of both transient bitmaps (tracking any arbitrary image format) or persistent bitmaps (qcow2) via live migration.

4.1.2 Supported Image Formats

QEMU supports all documented features below on the qcow2 image format.

However, qcow2 is only strictly necessary for the persistence feature, which writes bitmap data to disk upon close. If persistence is not required for a specific use case, all bitmap features excepting persistence are available for any arbitrary image format.

For example, Dirty Bitmaps can be combined with the ‘raw’ image format, but any changes to the bitmap will be discarded upon exit.

Warning: Transient bitmaps will not be saved on QEMU exit! Persistent bitmaps are available only on qcow2 images.

4.1.3 Dirty Bitmap Names

Bitmap objects need a method to reference them in the API. All API-created and managed bitmaps have a human-readable name chosen by the user at creation time.

- A bitmap’s name is unique to the node, but bitmaps attached to different nodes can share the same name. Therefore, all bitmaps are addressed via their (node, name) pair.
- The name of a user-created bitmap cannot be empty (“”).
- Transient bitmaps can have JSON unicode names that are effectively not length limited. (QMP protocol may restrict messages to less than 64MiB.)
- Persistent storage formats may impose their own requirements on bitmap names and namespaces. Presently, only qcow2 supports persistent bitmaps. See docs/interop/qcow2.txt for more details on restrictions. Notably:
 - qcow2 bitmap names are limited to between 1 and 1023 bytes long.
 - No two bitmaps saved to the same qcow2 file may share the same name.
- QEMU occasionally uses bitmaps for internal use which have no name. They are hidden from API query calls, cannot be manipulated by the external API, are never persistent, nor ever migrated.

4.1.4 Bitmap Status

Dirty Bitmap objects can be queried with the QMP command [query-block](#), and are visible via the [BlockDirtyInfo](#) QAPI structure.

This struct shows the name, granularity, and dirty byte count for each bitmap. Additionally, it shows several boolean status indicators:

- `recording`: This bitmap is recording writes.
- `busy`: This bitmap is in-use by an operation.
- `persistent`: This bitmap is a persistent type.

- `inconsistent`: This bitmap is corrupted and cannot be used.

The `+busy` status prohibits you from deleting, clearing, or otherwise modifying a bitmap, and happens when the bitmap is being used for a backup operation or is in the process of being loaded from a migration. Many of the commands documented below will refuse to work on such bitmaps.

The `+inconsistent` status similarly prohibits almost all operations, notably allowing only the `block-dirty-bitmap-remove` operation.

There is also a deprecated `status` field of type `DirtyBitmapStatus`. A bitmap historically had five visible states:

1. **Frozen**: This bitmap is currently in-use by an operation and is immutable. It can't be deleted, renamed, reset, etc.
(This is now `+busy`.)
2. **Disabled**: This bitmap is not recording new writes.
(This is now `-recording -busy`.)
3. **Active**: This bitmap is recording new writes.
(This is now `+recording -busy`.)
4. **Locked**: This bitmap is in-use by an operation, and is immutable. The difference from “Frozen” was primarily implementation details.
(This is now `+busy`.)
5. **Inconsistent**: This persistent bitmap was not saved to disk correctly, and can no longer be used. It remains in memory to serve as an indicator of failure.
(This is now `+inconsistent`.)

These states are directly replaced by the status indicators and should not be used. The difference between `Frozen` and `Locked` is an implementation detail and should not be relevant to external users.

4.1.5 Basic QMP Usage

The primary interface to manipulating bitmap objects is via the QMP interface. If you are not familiar, see `docs/interop/qmp-intro.txt` for a broad overview, and `qemu-qmp-ref` for a full reference of all QMP commands.

Supported Commands

There are six primary bitmap-management API commands:

- `block-dirty-bitmap-add`
- `block-dirty-bitmap-remove`
- `block-dirty-bitmap-clear`
- `block-dirty-bitmap-disable`
- `block-dirty-bitmap-enable`
- `block-dirty-bitmap-merge`

And one related query command:

- `query-block`

Creation: block-dirty-bitmap-add

block-dirty-bitmap-add:

Creates a new bitmap that tracks writes to the specified node. granularity, persistence, and recording state can be adjusted at creation time.

Example

to create a new, actively recording persistent bitmap:

```

-> { "execute": "block-dirty-bitmap-add",
    "arguments": {
        "node": "drive0",
        "name": "bitmap0",
        "persistent": true,
    }
}

<- { "return": {} }

```

- This bitmap will have a default granularity that matches the cluster size of its associated drive, if available, clamped to between [4KiB, 64KiB]. The current default for qcow2 is 64KiB.

Example

To create a new, disabled (`-recording`), transient bitmap that tracks changes in 32KiB segments:

```

-> { "execute": "block-dirty-bitmap-add",
    "arguments": {
        "node": "drive0",
        "name": "bitmap1",
        "granularity": 32768,
        "disabled": true
    }
}

<- { "return": {} }

```

Deletion: block-dirty-bitmap-remove

block-dirty-bitmap-remove:

Deletes a bitmap. Bitmaps that are `+busy` cannot be removed.

- Deleting a bitmap does not impact any other bitmaps attached to the same node, nor does it affect any backups already created from this bitmap or node.
- Because bitmaps are only unique to the node to which they are attached, you must specify the node/drive name here, too.
- Deleting a persistent bitmap will remove it from the qcow2 file.

Example

Remove a bitmap named `bitmap0` from node `drive0`:

```
-> { "execute": "block-dirty-bitmap-remove",  
    "arguments": {  
        "node": "drive0",  
        "name": "bitmap0"  
    }  
}  
  
<- { "return": {} }
```

Resetting: `block-dirty-bitmap-clear`

`block-dirty-bitmap-clear`:

Clears all dirty bits from a bitmap. `+busy` bitmaps cannot be cleared.

- An incremental backup created from an empty bitmap will copy no data, as if nothing has changed.

Example

Clear all dirty bits from bitmap `bitmap0` on node `drive0`:

```
-> { "execute": "block-dirty-bitmap-clear",  
    "arguments": {  
        "node": "drive0",  
        "name": "bitmap0"  
    }  
}  
  
<- { "return": {} }
```

Enabling: `block-dirty-bitmap-enable`

`block-dirty-bitmap-enable`:

“Enables” a bitmap, setting the `recording` bit to true, causing writes to begin being recorded. `+busy` bitmaps cannot be enabled.

- Bitmaps default to being enabled when created, unless configured otherwise.
- Persistent enabled bitmaps will remember their `+recording` status on load.

Example

To set `+recording` on bitmap `bitmap0` on node `drive0`:

```
-> { "execute": "block-dirty-bitmap-enable",  
    "arguments": {  
        "node": "drive0",  
        "name": "bitmap0"  
    }  
}
```

(continues on next page)

(continued from previous page)

```
<- { "return": {} }
```

Enabling: block-dirty-bitmap-disable

block-dirty-bitmap-disable:

“Disables” a bitmap, setting the `recording` bit to false, causing further writes to begin being ignored. `+busy` bitmaps cannot be disabled.

Warning: This is potentially dangerous: QEMU makes no effort to stop any writes if there are disabled bitmaps on a node, and will not mark any disabled bitmaps as `+inconsistent` if any such writes do happen. Backups made from such bitmaps will not be able to be used to reconstruct a coherent image.

- Disabling a bitmap may be useful for examining which sectors of a disk changed during a specific time period, or for explicit management of differential backup windows.
- Persistent disabled bitmaps will remember their `-recording` status on load.

Example

To set `-recording` on bitmap `bitmap0` on node `drive0`:

```
-> { "execute": "block-dirty-bitmap-disable",
    "arguments": {
        "node": "drive0",
        "name": "bitmap0"
    }
}

<- { "return": {} }
```

Merging, Copying: block-dirty-bitmap-merge

block-dirty-bitmap-merge:

Merges one or more bitmaps into a target bitmap. For any segment that is dirty in any one source bitmap, the target bitmap will mark that segment dirty.

- Merge takes one or more bitmaps as a source and merges them together into a single destination, such that any segment marked as dirty in any source bitmap(s) will be marked dirty in the destination bitmap.
- Merge does not create the destination bitmap if it does not exist. A blank bitmap can be created beforehand to achieve the same effect.
- The destination is not cleared prior to merge, so subsequent merge operations will continue to cumulatively mark more segments as dirty.
- If the merge operation should fail, the destination bitmap is guaranteed to be unmodified. The operation may fail if the source or destination bitmaps are busy, or have different granularities.

- Bitmaps can only be merged on the same node. There is only one “node” argument, so all bitmaps must be attached to that same node.
- Copy can be achieved by merging from a single source to an empty destination.

Example

Merge the data from `bitmap0` into the bitmap `new_bitmap` on node `drive0`. If `new_bitmap` was empty prior to this command, this achieves a copy.

```
-> { "execute": "block-dirty-bitmap-merge",
      "arguments": {
        "node": "drive0",
        "target": "new_bitmap",
        "bitmaps": [ "bitmap0" ]
      }
}

<- { "return": {} }
```

Querying: query-block

`query-block`:

Not strictly a bitmaps command, but will return information about any bitmaps attached to nodes serving as the root for guest devices.

- The “inconsistent” bit will not appear when it is false, appearing only when the value is true to indicate there is a problem.

Example

Query the block sub-system of QEMU. The following json has trimmed irrelevant keys from the response to highlight only the bitmap-relevant portions of the API. This result highlights a bitmap `bitmap0` attached to the root node of device `drive0`.

```
-> {
  "execute": "query-block",
  "arguments": {}
}

<- {
  "return": [ {
    "dirty-bitmaps": [ {
      "status": "active",
      "count": 0,
      "busy": false,
      "name": "bitmap0",
      "persistent": false,
      "recording": true,
      "granularity": 65536
    } ],
    "device": "drive0",
  } ]
}
```

4.1.6 Bitmap Persistence

As outlined in *Supported Image Formats*, QEMU can persist bitmaps to qcow2 files. Demonstrated in *Creation: block-dirty-bitmap-add*, passing `persistent: true` to `block-dirty-bitmap-add` will persist that bitmap to disk.

Persistent bitmaps will be automatically loaded into memory upon load, and will be written back to disk upon close. Their usage should be mostly transparent.

However, if QEMU does not get a chance to close the file cleanly, the bitmap will be marked as `+inconsistent` at next load and considered unsafe to use for any operation. At this point, the only valid operation on such bitmaps is `block-dirty-bitmap-remove`.

Losing a bitmap in this way does not invalidate any existing backups that have been made from this bitmap, but no further backups will be able to be issued for this chain.

4.1.7 Transactions

Transactions are a QMP feature that allows you to submit multiple QMP commands at once, being guaranteed that they will all succeed or fail atomically, together. The interaction of bitmaps and transactions are demonstrated below.

See `transaction` in the QMP reference for more details.

Justification

Bitmaps can generally be modified at any time, but certain operations often only make sense when paired directly with other commands. When a VM is paused, it's easy to ensure that no guest writes occur between individual QMP commands. When a VM is running, this is difficult to accomplish with individual QMP commands that may allow guest writes to occur between each command.

For example, using only individual QMP commands, we could:

1. Boot the VM in a paused state.
2. Create a full drive backup of `drive0`.
3. Create a new bitmap attached to `drive0`, confident that nothing has been written to `drive0` in the meantime.
4. Resume execution of the VM.
5. At a later point, issue incremental backups from `bitmap0`.

At this point, the bitmap and drive backup would be correctly in sync, and incremental backups made from this point forward would be correctly aligned to the full drive backup.

This is not particularly useful if we decide we want to start incremental backups after the VM has been running for a while, for which we would want to perform actions such as the following:

1. Boot the VM and begin execution.
2. Using a single transaction, perform the following operations:
 - Create `bitmap0`.
 - Create a full drive backup of `drive0`.
3. At a later point, issue incremental backups from `bitmap0`.

Note: As a consideration, if `bitmap0` is created prior to the full drive backup, incremental backups can still be authored from this bitmap, but they will copy extra segments reflecting writes that occurred prior to the backup operation. Transactions allow us to narrow critical points in time to reduce waste, or, in the other direction, to ensure that no segments are omitted.

Supported Bitmap Transactions

- `block-dirty-bitmap-add`
- `block-dirty-bitmap-clear`
- `block-dirty-bitmap-enable`
- `block-dirty-bitmap-disable`
- `block-dirty-bitmap-merge`

The usages for these commands are identical to their respective QMP commands, but see the sections below for concrete examples.

4.1.8 Incremental Backups - Push Model

Incremental backups are simply partial disk images that can be combined with other partial disk images on top of a base image to reconstruct a full backup from the point in time at which the incremental backup was issued.

The “Push Model” here references the fact that QEMU is “pushing” the modified blocks out to a destination. We will be using the `drive-backup` and `blockdev-backup` QMP commands to create both full and incremental backups.

Both of these commands are jobs, which have their own QMP API for querying and management documented in [Background jobs](#).

Example: New Incremental Backup Anchor Point

As outlined in the Transactions - *Justification* section, perhaps we want to create a new incremental backup chain attached to a drive.

This example creates a new, full backup of “drive0” and accompanies it with a new, empty bitmap that records writes from this point in time forward.

Note: Any new writes that happen after this command is issued, even while the backup job runs, will be written locally and not to the backup destination. These writes will be recorded in the bitmap accordingly.

```
-> {
  "execute": "transaction",
  "arguments": {
    "actions": [
      {
        "type": "block-dirty-bitmap-add",
        "data": {
          "node": "drive0",
          "name": "bitmap0"
        }
      }
    ]
  }
}
```

(continues on next page)

(continued from previous page)

```

    },
    {
      "type": "drive-backup",
      "data": {
        "device": "drive0",
        "target": "/path/to/drive0.full.qcow2",
        "sync": "full",
        "format": "qcow2"
      }
    }
  ]
}

```

```
<- { "return": {} }
```

```
<- {
  "timestamp": {
    "seconds": 1555436945,
    "microseconds": 179620
  },
  "data": {
    "status": "created",
    "id": "drive0"
  },
  "event": "JOB_STATUS_CHANGE"
}

```

```
...
```

```
<- {
  "timestamp": {...},
  "data": {
    "device": "drive0",
    "type": "backup",
    "speed": 0,
    "len": 68719476736,
    "offset": 68719476736
  },
  "event": "BLOCK_JOB_COMPLETED"
}

```

```
<- {
  "timestamp": {...},
  "data": {
    "status": "concluded",
    "id": "drive0"
  },
  "event": "JOB_STATUS_CHANGE"
}

```

```
<- {
  "timestamp": {...},
  "data": {
    "status": "null",
    "id": "drive0"
  },
}

```

(continues on next page)

(continued from previous page)

```

    "event": "JOB_STATUS_CHANGE"
  }

```

A full explanation of the job transition semantics and the JOB_STATUS_CHANGE event are beyond the scope of this document and will be omitted in all subsequent examples; above, several more events have been omitted for brevity.

Note: Subsequent examples will omit all events except BLOCK_JOB_COMPLETED except where necessary to illustrate workflow differences.

Omitted events and json objects will be represented by ellipses: . . .

Example: Resetting an Incremental Backup Anchor Point

If we want to start a new backup chain with an existing bitmap, we can also use a transaction to reset the bitmap while making a new full backup:

```

-> {
  "execute": "transaction",
  "arguments": {
    "actions": [
      {
        "type": "block-dirty-bitmap-clear",
        "data": {
          "node": "drive0",
          "name": "bitmap0"
        }
      },
      {
        "type": "drive-backup",
        "data": {
          "device": "drive0",
          "target": "/path/to/drive0.new_full.qcow2",
          "sync": "full",
          "format": "qcow2"
        }
      }
    ]
  }
}

<- { "return": {} }

...

<- {
  "timestamp": {...},
  "data": {
    "device": "drive0",
    "type": "backup",
    "speed": 0,
    "len": 68719476736,
    "offset": 68719476736
  },
  "event": "BLOCK_JOB_COMPLETED"
}

```

(continues on next page)

(continued from previous page)

}

The result of this example is identical to the first, but we clear an existing bitmap instead of adding a new one.

Tip: In both of these examples, “bitmap0” is tied conceptually to the creation of new, full backups. This relationship is not saved or remembered by QEMU; it is up to the operator or management layer to remember which bitmaps are associated with which backups.

Example: First Incremental Backup

1. Create a full backup and sync it to a dirty bitmap using any method:
 - Either of the two live backup method demonstrated above,
 - Using QMP commands with the VM paused as in the *Justification* section, or
 - With the VM offline, manually copy the image and start the VM in a paused state, careful to add a new bitmap before the VM begins execution.

Whichever method is chosen, let’s assume that at the end of this step:

- The full backup is named `drive0.full.qcow2`.
- The bitmap we created is named `bitmap0`, attached to `drive0`.

2. Create a destination image for the incremental backup that utilizes the full backup as a backing image.
 - Let’s assume the new incremental image is named `drive0.inc0.qcow2`:

```
$ qemu-img create -f qcow2 drive0.inc0.qcow2 \
-b drive0.full.qcow2 -F qcow2
```

3. Issue an incremental backup command:

```
-> {
  "execute": "drive-backup",
  "arguments": {
    "device": "drive0",
    "bitmap": "bitmap0",
    "target": "drive0.inc0.qcow2",
    "format": "qcow2",
    "sync": "incremental",
    "mode": "existing"
  }
}

<- { "return": {} }

...

<- {
  "timestamp": {...},
  "data": {
    "device": "drive0",
```

(continues on next page)

(continued from previous page)

```

    "type": "backup",
    "speed": 0,
    "len": 68719476736,
    "offset": 68719476736
  },
  "event": "BLOCK_JOB_COMPLETED"
}

...

```

This copies any blocks modified since the full backup was created into the `drive0.inc0.qcow2` file. During the operation, `bitmap0` is marked `+busy`. If the operation is successful, `bitmap0` will be cleared to reflect the “incremental” backup regimen, which only copies out new changes from each incremental backup.

Note: Any new writes that occur after the backup operation starts do not get copied to the destination. The backup’s “point in time” is when the backup starts, not when it ends. These writes are recorded in a special bitmap that gets re-added to `bitmap0` when the backup ends so that the next incremental backup can copy them out.

Example: Second Incremental Backup

1. Create a new destination image for the incremental backup that points to the previous one, e.g.: `drive0.inc1.qcow2`

```

$ qemu-img create -f qcow2 drive0.inc1.qcow2 \
  -b drive0.inc0.qcow2 -F qcow2

```

2. Issue a new incremental backup command. The only difference here is that we have changed the target image below.

```

-> {
  "execute": "drive-backup",
  "arguments": {
    "device": "drive0",
    "bitmap": "bitmap0",
    "target": "drive0.inc1.qcow2",
    "format": "qcow2",
    "sync": "incremental",
    "mode": "existing"
  }
}

<- { "return": {} }

...

<- {
  "timestamp": {...},
  "data": {
    "device": "drive0",
    "type": "backup",
    "speed": 0,
    "len": 68719476736,
    "offset": 68719476736
  }
}

```

(continues on next page)

(continued from previous page)

```

    },
    "event": "BLOCK_JOB_COMPLETED"
}
...

```

Because the first incremental backup from the previous example completed successfully, `bitmap0` was synchronized with `drive0.inc0.qcow2`. Here, we use `bitmap0` again to create a new incremental backup that targets the previous one, creating a chain of three images:

Diagram

```

+-----+ +-----+ +-----+
| drive0.full.qcow2 |<--| drive0.inc0.qcow2 |<--| drive0.inc1.qcow2 |
+-----+ +-----+ +-----+

```

Each new incremental backup re-synchronizes the bitmap to the latest backup authored, allowing a user to continue to “consume” it to create new backups on top of an existing chain.

In the above diagram, neither `drive0.inc1.qcow2` nor `drive0.inc0.qcow2` are complete images by themselves, but rely on their backing chain to reconstruct a full image. The dependency terminates with each full backup.

Each backup in this chain remains independent, and is unchanged by new entries made later in the chain. For instance, `drive0.inc0.qcow2` remains a perfectly valid backup of the disk as it was when that backup was issued.

Example: Incremental Push Backups without Backing Files

Backup images are best kept off-site, so we often will not have the preceding backups in a chain available to link against. This is not a problem at backup time; we simply do not set the backing image when creating the destination image:

1. Create a new destination image with no backing file set. We will need to specify the size of the base image, because the backing file isn’t available for QEMU to use to determine it.

```
$ qemu-img create -f qcow2 drive0.inc2.qcow2 64G
```

Note: Alternatively, you can omit `mode: "existing"` from the push backup commands to have QEMU create an image without a backing file for you, but you lose control over format options like compatibility and preallocation presets.

2. Issue a new incremental backup command. Apart from the new destination image, there is no difference from the last two examples.

```

-> {
  "execute": "drive-backup",
  "arguments": {
    "device": "drive0",
    "bitmap": "bitmap0",
    "target": "drive0.inc2.qcow2",
    "format": "qcow2",
    "sync": "incremental",
    "mode": "existing"
  }
}

```

(continues on next page)

(continued from previous page)

```

    }
  }

  <- { "return": {} }

  ...

  <- {
    "timestamp": {...},
    "data": {
      "device": "drive0",
      "type": "backup",
      "speed": 0,
      "len": 68719476736,
      "offset": 68719476736
    },
    "event": "BLOCK_JOB_COMPLETED"
  }

  ...

```

The only difference from the perspective of the user is that you will need to set the backing image when attempting to restore the backup:

```
$ qemu-img rebase drive0.inc2.qcow2 \
-u -b drive0.inc1.qcow2
```

This uses the “unsafe” rebase mode to simply set the backing file to a file that isn’t present.

It is also possible to use `--image-opts` to specify the entire backing chain by hand as an ephemeral property at runtime, but that is beyond the scope of this document.

Example: Multi-drive Incremental Backup

Assume we have a VM with two drives, “drive0” and “drive1” and we wish to back both of them up such that the two backups represent the same crash-consistent point in time.

1. For each drive, create an empty image:

```
$ qemu-img create -f qcow2 drive0.full.qcow2 64G
$ qemu-img create -f qcow2 drive1.full.qcow2 64G
```

2. Create a full (anchor) backup for each drive, with accompanying bitmaps:

```

-> {
  "execute": "transaction",
  "arguments": {
    "actions": [
      {
        "type": "block-dirty-bitmap-add",
        "data": {
          "node": "drive0",
          "name": "bitmap0"
        }
      }
    ],
  },
  {

```

(continues on next page)

(continued from previous page)

```

        "type": "block-dirty-bitmap-add",
        "data": {
            "node": "drive1",
            "name": "bitmap0"
        }
    },
    {
        "type": "drive-backup",
        "data": {
            "device": "drive0",
            "target": "/path/to/drive0.full.qcow2",
            "sync": "full",
            "format": "qcow2"
        }
    },
    {
        "type": "drive-backup",
        "data": {
            "device": "drive1",
            "target": "/path/to/drive1.full.qcow2",
            "sync": "full",
            "format": "qcow2"
        }
    }
]
}
}

<- { "return": {} }

...

<- {
    "timestamp": {...},
    "data": {
        "device": "drive0",
        "type": "backup",
        "speed": 0,
        "len": 68719476736,
        "offset": 68719476736
    },
    "event": "BLOCK_JOB_COMPLETED"
}

...

<- {
    "timestamp": {...},
    "data": {
        "device": "drive1",
        "type": "backup",
        "speed": 0,
        "len": 68719476736,
        "offset": 68719476736
    },
    "event": "BLOCK_JOB_COMPLETED"
}

```

(continues on next page)

(continued from previous page)

...

3. Later, create new destination images for each of the incremental backups that point to their respective full backups:

```
$ qemu-img create -f qcow2 drive0.inc0.qcow2 \
-b drive0.full.qcow2 -F qcow2
$ qemu-img create -f qcow2 drive1.inc0.qcow2 \
-b drive1.full.qcow2 -F qcow2
```

4. Issue a multi-drive incremental push backup transaction:

```
-> {
  "execute": "transaction",
  "arguments": {
    "actions": [
      {
        "type": "drive-backup",
        "data": {
          "device": "drive0",
          "bitmap": "bitmap0",
          "format": "qcow2",
          "mode": "existing",
          "sync": "incremental",
          "target": "drive0.inc0.qcow2"
        }
      },
      {
        "type": "drive-backup",
        "data": {
          "device": "drive1",
          "bitmap": "bitmap0",
          "format": "qcow2",
          "mode": "existing",
          "sync": "incremental",
          "target": "drive1.inc0.qcow2"
        }
      }
    ]
  }
}

<- { "return": {} }

...

<- {
  "timestamp": {...},
  "data": {
    "device": "drive0",
    "type": "backup",
    "speed": 0,
    "len": 68719476736,
    "offset": 68719476736
  },
```

(continues on next page)

(continued from previous page)

```

    "event": "BLOCK_JOB_COMPLETED"
  }
  ...
  <- {
    "timestamp": {...},
    "data": {
      "device": "drive1",
      "type": "backup",
      "speed": 0,
      "len": 68719476736,
      "offset": 68719476736
    },
    "event": "BLOCK_JOB_COMPLETED"
  }
  ...

```

4.1.9 Push Backup Errors & Recovery

In the event of an error that occurs after a push backup job is successfully launched, either by an individual QMP command or a QMP transaction, the user will receive a `BLOCK_JOB_COMPLETE` event with a failure message, accompanied by a `BLOCK_JOB_ERROR` event.

In the case of a job being cancelled, the user will receive a `BLOCK_JOB_CANCELLED` event instead of a pair of `COMPLETE` and `ERROR` events.

In either failure case, the bitmap used for the failed operation is not cleared. It will contain all of the dirty bits it did at the start of the operation, plus any new bits that got marked during the operation.

Effectively, the “point in time” that a bitmap is recording differences against is kept at the issuance of the last successful incremental backup, instead of being moved forward to the start of this now-failed backup.

Once the underlying problem is addressed (e.g. more storage space is allocated on the destination), the incremental backup command can be retried with the same bitmap.

Example: Individual Failures

Incremental Push Backup jobs that fail individually behave simply as described above. This example demonstrates the single-job failure case:

1. Create a target image:

```

$ qemu-img create -f qcow2 drive0.inc0.qcow2 \
  -b drive0.full.qcow2 -F qcow2

```

2. Attempt to create an incremental backup via QMP:

```

-> {
  "execute": "drive-backup",
  "arguments": {
    "device": "drive0",
    "bitmap": "bitmap0",
    "target": "drive0.inc0.qcow2",

```

(continues on next page)

(continued from previous page)

```

    "format": "qcow2",
    "sync": "incremental",
    "mode": "existing"
  }
}

<- { "return": {} }

```

3. Receive a pair of events indicating failure:

```

<- {
  "timestamp": {...},
  "data": {
    "device": "drive0",
    "action": "report",
    "operation": "write"
  },
  "event": "BLOCK_JOB_ERROR"
}

<- {
  "timestamp": {...},
  "data": {
    "speed": 0,
    "offset": 0,
    "len": 67108864,
    "error": "No space left on device",
    "device": "drive0",
    "type": "backup"
  },
  "event": "BLOCK_JOB_COMPLETED"
}

```

4. Delete the failed image, and re-create it.

```

$ rm drive0.inc0.qcow2
$ qemu-img create -f qcow2 drive0.inc0.qcow2 \
  -b drive0.full.qcow2 -F qcow2

```

5. Retry the command after fixing the underlying problem, such as freeing up space on the backup volume:

```

-> {
  "execute": "drive-backup",
  "arguments": {
    "device": "drive0",
    "bitmap": "bitmap0",
    "target": "drive0.inc0.qcow2",
    "format": "qcow2",
    "sync": "incremental",
    "mode": "existing"
  }
}

<- { "return": {} }

```

6. Receive confirmation that the job completed successfully:

```

<- {
  "timestamp": { ... },
  "data": {
    "device": "drive0",
    "type": "backup",
    "speed": 0,
    "len": 67108864,
    "offset": 67108864
  },
  "event": "BLOCK_JOB_COMPLETED"
}

```

Example: Partial Transactional Failures

QMP commands like `drive-backup` conceptually only start a job, and so transactions containing these commands may succeed even if the job it created later fails. This might have surprising interactions with notions of how a “transaction” ought to behave.

This distinction means that on occasion, a transaction containing such job launching commands may appear to succeed and return success, but later individual jobs associated with the transaction may fail. It is possible that a management application may have to deal with a partial backup failure after a “successful” transaction.

If multiple backup jobs are specified in a single transaction, if one of those jobs fails, it will not interact with the other backup jobs in any way by default. The job(s) that succeeded will clear the dirty bitmap associated with the operation, but the job(s) that failed will not. It is therefore not safe to delete any incremental backups that were created successfully in this scenario, even though others failed.

This example illustrates a transaction with two backup jobs, where one fails and one succeeds:

1. Issue the transaction to start a backup of both drives.

```

-> {
  "execute": "transaction",
  "arguments": {
    "actions": [
      {
        "type": "drive-backup",
        "data": {
          "device": "drive0",
          "bitmap": "bitmap0",
          "format": "qcow2",
          "mode": "existing",
          "sync": "incremental",
          "target": "drive0.inc0.qcow2"
        }
      },
      {
        "type": "drive-backup",
        "data": {
          "device": "drive1",
          "bitmap": "bitmap0",
          "format": "qcow2",
          "mode": "existing",
          "sync": "incremental",
          "target": "drive1.inc0.qcow2"
        }
      }
    ]
  }
}

```

(continues on next page)

(continued from previous page)

```
}
}
```

2. Receive notice that the Transaction was accepted, and jobs were launched:

```
<- { "return": {} }
```

3. Receive notice that the first job has completed:

```
<- {
  "timestamp": {...},
  "data": {
    "device": "drive0",
    "type": "backup",
    "speed": 0,
    "len": 67108864,
    "offset": 67108864
  },
  "event": "BLOCK_JOB_COMPLETED"
}
```

4. Receive notice that the second job has failed:

```
<- {
  "timestamp": {...},
  "data": {
    "device": "drive1",
    "action": "report",
    "operation": "read"
  },
  "event": "BLOCK_JOB_ERROR"
}

...

<- {
  "timestamp": {...},
  "data": {
    "speed": 0,
    "offset": 0,
    "len": 67108864,
    "error": "Input/output error",
    "device": "drive1",
    "type": "backup"
  },
  "event": "BLOCK_JOB_COMPLETED"
}
```

At the conclusion of the above example, `drive0.inc0.qcow2` is valid and must be kept, but `drive1.inc0.qcow2` is incomplete and should be deleted. If a VM-wide incremental backup of all drives at a point-in-time is to be made, new backups for both drives will need to be made, taking into account that a new incremental backup for `drive0` needs to be based on top of `drive0.inc0.qcow2`.

For this example, an incremental backup for `drive0` was created, but not for `drive1`. The last VM-wide crash-consistent backup that is available in this case is the full backup:


```
[drive0.full.qcow2] <-- [drive0.inc0.qcow2]
[drive1.full.qcow2]
```

To repair this, issue a new incremental backup across both drives. The result will be backup chains that resemble the following:

```
[drive0.full.qcow2] <-- [drive0.inc0.qcow2] <-- [drive0.inc1.qcow2]
[drive1.full.qcow2] <----- [drive1.inc1.qcow2]
```

Example: Grouped Completion Mode

While jobs launched by transactions normally complete or fail individually, it's possible to instruct them to complete or fail together as a group. QMP transactions take an optional properties structure that can affect the behavior of the transaction.

The `completion-mode` transaction property can be either `individual` which is the default legacy behavior described above, or `grouped`, detailed below.

In `grouped` completion mode, no jobs will report success until all jobs are ready to report success. If any job fails, all other jobs will be cancelled.

Regardless of if a participating incremental backup job failed or was cancelled, their associated bitmaps will all be held at their existing points-in-time, as in individual failure cases.

Here's the same multi-drive backup scenario from *Example: Partial Transactional Failures*, but with the `grouped` completion-mode property applied:

1. Issue the multi-drive incremental backup transaction:

```
-> {
  "execute": "transaction",
  "arguments": {
    "properties": {
      "completion-mode": "grouped"
    },
    "actions": [
      {
        "type": "drive-backup",
        "data": {
          "device": "drive0",
          "bitmap": "bitmap0",
          "format": "qcow2",
          "mode": "existing",
          "sync": "incremental",
          "target": "drive0.inc0.qcow2"
        }
      },
      {
        "type": "drive-backup",
        "data": {
          "device": "drive1",
          "bitmap": "bitmap0",
          "format": "qcow2",
          "mode": "existing",
          "sync": "incremental",
          "target": "drive1.inc0.qcow2"
        }
      }
    ]
  }
}
```

(continues on next page)

(continued from previous page)

```
    ]]  
  }  
}
```

2. Receive notice that the Transaction was accepted, and jobs were launched:

```
<- { "return": {} }
```

3. Receive notification that the backup job for drive1 has failed:

```
<- {  
  "timestamp": {...},  
  "data": {  
    "device": "drive1",  
    "action": "report",  
    "operation": "read"  
  },  
  "event": "BLOCK_JOB_ERROR"  
}  
  
<- {  
  "timestamp": {...},  
  "data": {  
    "speed": 0,  
    "offset": 0,  
    "len": 67108864,  
    "error": "Input/output error",  
    "device": "drive1",  
    "type": "backup"  
  },  
  "event": "BLOCK_JOB_COMPLETED"  
}
```

4. Receive notification that the job for drive0 has been cancelled:

```
<- {  
  "timestamp": {...},  
  "data": {  
    "device": "drive0",  
    "type": "backup",  
    "speed": 0,  
    "len": 67108864,  
    "offset": 16777216  
  },  
  "event": "BLOCK_JOB_CANCELLED"  
}
```

At the conclusion of *this* example, both jobs have been aborted due to a failure. Both destination images should be deleted and are no longer of use.

The transaction as a whole can simply be re-issued at a later time.

4.2 D-Bus

4.2.1 Introduction

QEMU may be running with various helper processes involved:

- vhost-user* processes (gpu, virtfs, input, etc...)
- TPM emulation (or other devices)
- user networking (slirp)
- network services (DHCP/DNS, samba/ftp etc)
- background tasks (compression, streaming etc)
- client UI
- admin & cli

Having several processes allows stricter security rules, as well as greater modularity.

While QEMU itself uses QMP as primary IPC (and Spice/VNC for remote display), D-Bus is the de facto IPC of choice on Unix systems. The wire format is machine friendly, good bindings exist for various languages, and there are various tools available.

Using a bus, helper processes can discover and communicate with each other easily, without going through QEMU. The bus topology is also easier to apprehend and debug than a mesh. However, it is wise to consider the security aspects of it.

4.2.2 Security

A QEMU D-Bus bus should be private to a single VM. Thus, only cooperative tasks are running on the same bus to serve the VM.

D-Bus, the protocol and standard, doesn't have mechanisms to enforce security between peers once the connection is established. Peers may have additional mechanisms to enforce security rules, based for example on UNIX credentials.

The daemon can control which peers can send/rcv messages using various metadata attributes, however, this is alone is not generally sufficient to make the deployment secure. The semantics of the actual methods implemented using D-Bus are just as critical. Peers need to carefully validate any information they received from a peer with a different trust level.

dbus-daemon policy

dbus-daemon can enforce various policies based on the UID/GID of the processes that are connected to it. It is thus a good idea to run helpers as different UID from QEMU and set appropriate policies.

Depending on the use case, you may choose different scenarios:

- Everything the same UID
 - Convenient for developers
 - Improved reliability - crash of one part doesn't take out entire VM
 - No security benefit over traditional QEMU, unless additional controls such as SELinux or AppArmor are applied
- Two UIDs, one for QEMU, one for dbus & helpers

- Moderately improved user based security isolation
- Many UIDs, one for QEMU one for dbus and one for each helpers
 - Best user based security isolation
 - Complex to manager distinct UIDs needed for each VM

For example, to allow only `qemu` user to talk to `qemu-helper` `org.qemu.Helper1` service, a `dbus-daemon` policy may contain:

```
<policy user="qemu">
  <allow send_destination="org.qemu.Helper1"/>
  <allow receive_sender="org.qemu.Helper1"/>
</policy>

<policy user="qemu-helper">
  <allow own="org.qemu.Helper1"/>
</policy>
```

`dbus-daemon` can also perform SELinux checks based on the security context of the source and the target. For example, `virtiofs_t` could be allowed to send a message to `svirt_t`, but `virtiofs_t` wouldn't be allowed to send a message to `virtiofs_t`.

See `dbus-daemon` man page for details.

4.2.3 Guidelines

When implementing new D-Bus interfaces, it is recommended to follow the “D-Bus API Design Guidelines”: <https://dbus.freedesktop.org/doc/dbus-api-design.html>

The “`org.qemu.*`” prefix is reserved for services implemented & distributed by the QEMU project.

4.2.4 QEMU Interfaces

D-Bus VMState

4.3 D-Bus VMState

4.3.1 Introduction

The QEMU `dbus-vmstate` object's aim is to migrate helpers' data running on a QEMU D-Bus bus. (refer to the *D-Bus* document for some recommendations on D-Bus usage)

Upon migration, QEMU will go through the queue of `org.qemu.VMState1` D-Bus name owners and query their `Id`. It must be unique among the helpers.

It will then save arbitrary data of each `Id` to be transferred in the migration stream and restored/loaded at the corresponding destination helper.

For now, the data amount to be transferred is arbitrarily limited to 1Mb. The state must be saved quickly (a fraction of a second). (D-Bus imposes a time limit on reply anyway, and migration would fail if data isn't given quickly enough.)

`dbus-vmstate` object can be configured with the expected list of helpers by setting its `id-list` property, with a comma-separated `Id` list.

4.3.2 Interface

On object path `/org/qemu/VMState1`, the following `org.qemu.VMState1` interface should be implemented:

```
<interface name="org.qemu.VMState1">
  <property name="Id" type="s" access="read"/>
  <method name="Load">
    <arg type="ay" name="data" direction="in"/>
  </method>
  <method name="Save">
    <arg type="ay" name="data" direction="out"/>
  </method>
</interface>
```

“Id” property

A string that identifies the helper uniquely. (maximum 256 bytes including terminating NUL byte)

Note: The helper ID namespace is a separate namespace. In particular, it is not related to QEMU “id” used in `-object/-device` objects.

Load(in u8[] bytes) method

The method called on destination with the state to restore.

The helper may be initially started in a waiting state (with an `-incoming` argument for example), and it may resume on success.

An error may be returned to the caller.

Save(out u8[] bytes) method

The method called on the source to get the current state to be migrated. The helper should continue to run normally.

An error may be returned to the caller.

4.4 Live Block Device Operations

QEMU Block Layer currently (as of QEMU 2.9) supports four major kinds of live block device jobs – stream, commit, mirror, and backup. These can be used to manipulate disk image chains to accomplish certain tasks, namely: live copy data from backing files into overlays; shorten long disk image chains by merging data from overlays into backing files; live synchronize data from a disk image chain (including current active disk) to another target image; and point-in-time (and incremental) backups of a block device. Below is a description of the said block (QMP) primitives, and some (non-exhaustive list of) examples to illustrate their use.

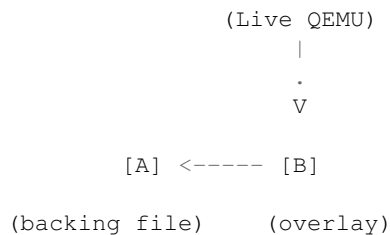
Note: The file `qapi/block-core.json` in the QEMU source tree has the canonical QEMU API (QAPI) schema documentation for the QMP primitives discussed here.

Contents

- *Live Block Device Operations*
 - *Disk image backing chain notation*
 - *Brief overview of live block QMP primitives*
 - *Interacting with a QEMU instance*
 - *Example disk image chain*
 - *A note on points-in-time vs file names*
 - *Live block streaming — block-stream*
 - * *QMP invocation for block-stream*
 - *Live block commit — block-commit*
 - * *QMP invocation for block-commit*
 - *Live disk synchronization — drive-mirror and blockdev-mirror*
 - * *QMP invocation for drive-mirror*
 - * *QMP invocation for live storage migration with drive-mirror + NBD*
 - * *Notes on blockdev-mirror*
 - * *QMP invocation for blockdev-mirror*
 - *Live disk backup — drive-backup and blockdev-backup*
 - * *QMP invocation for drive-backup*
 - * *Notes on blockdev-backup*
 - * *QMP invocation for blockdev-backup*

4.4.1 Disk image backing chain notation

A simple disk image chain. (This can be created live using QMP `blockdev-snapshot-sync`, or offline via `qemu-img`):



The arrow can be read as: Image [A] is the backing file of disk image [B]. And live QEMU is currently writing to image [B], consequently, it is also referred to as the “active layer”.

There are two kinds of terminology that are common when referring to files in a disk image backing chain:

- (1) Directional: ‘base’ and ‘top’. Given the simple disk image chain above, image [A] can be referred to as ‘base’, and image [B] as ‘top’. (This terminology can be seen in in QAPI schema file, `block-core.json`.)

- (2) Relational: ‘backing file’ and ‘overlay’. Again, taking the same simple disk image chain from the above, disk image [A] is referred to as the backing file, and image [B] as overlay.

Throughout this document, we will use the relational terminology.

Important: The overlay files can generally be any format that supports a backing file, although QCOW2 is the preferred format and the one used in this document.

4.4.2 Brief overview of live block QMP primitives

The following are the four different kinds of live block operations that QEMU block layer supports.

- (1) `block-stream`: Live copy of data from backing files into overlay files.

Note: Once the ‘stream’ operation has finished, three things to note:

- (a) QEMU rewrites the backing chain to remove reference to the now-streamed and redundant backing file;
 - (b) the streamed file *itself* won’t be removed by QEMU, and must be explicitly discarded by the user;
 - (c) the streamed file remains valid – i.e. further overlays can be created based on it. Refer the `block-stream` section further below for more details.
-

- (2) `block-commit`: Live merge of data from overlay files into backing files (with the optional goal of removing the overlay file from the chain). Since QEMU 2.0, this includes “active `block-commit`” (i.e. merge the current active layer into the base image).

Note: Once the ‘commit’ operation has finished, there are three things to note here as well:

- (a) QEMU rewrites the backing chain to remove reference to now-redundant overlay images that have been committed into a backing file;
 - (b) the committed file *itself* won’t be removed by QEMU – it ought to be manually removed;
 - (c) however, unlike in the case of `block-stream`, the intermediate images will be rendered invalid – i.e. no more further overlays can be created based on them. Refer the `block-commit` section further below for more details.
-

- (3) `drive-mirror` (and `blockdev-mirror`): Synchronize a running disk to another image.

- (4) `drive-backup` (and `blockdev-backup`): Point-in-time (live) copy of a block device to a destination.

4.4.3 Interacting with a QEMU instance

To show some example invocations of command-line, we will use the following invocation of QEMU, with a QMP server running over UNIX socket:

```
$ ./qemu-system-x86_64 -display none -no-user-config \
-M q35 -nodefaults -m 512 \
-blockdev node-name=node-A,driver=qcow2,file.driver=file,file.node-name=file,file.
↪filename=./a.qcow2 \
-device virtio-blk,drive=node-A,id=virtio0 \
-monitor stdio -qmp unix:/tmp/qmp-sock,server,nowait
```

The `-blockdev` command-line option, used above, is available from QEMU 2.9 onwards. In the above invocation, notice the `node-name` parameter that is used to refer to the disk image `a.qcow2` ('node-A') – this is a cleaner way to refer to a disk image (as opposed to referring to it by spelling out file paths). So, we will continue to designate a `node-name` to each further disk image created (either via `blockdev-snapshot-sync`, or `blockdev-add`) as part of the disk image chain, and continue to refer to the disks using their `node-name` (where possible, because `block-commit` does not yet, as of QEMU 2.9, accept `node-name` parameter) when performing various block operations.

To interact with the QEMU instance launched above, we will use the `qmp-shell` utility (located at: `qemu/scripts/qmp`, as part of the QEMU source directory), which takes key-value pairs for QMP commands. Invoke it as below (which will also print out the complete raw JSON syntax for reference – examples in the following sections):

```
$ ./qmp-shell -v -p /tmp/qmp-sock
(QEMU)
```

Note: In the event we have to repeat a certain QMP command, we will: for the first occurrence of it, show the `qmp-shell` invocation, *and* the corresponding raw JSON QMP syntax; but for subsequent invocations, present just the `qmp-shell` syntax, and omit the equivalent JSON output.

4.4.4 Example disk image chain

We will use the below disk image chain (and occasionally spelling it out where appropriate) when discussing various primitives:

```
[A] <-- [B] <-- [C] <-- [D]
```

Where [A] is the original base image; [B] and [C] are intermediate overlay images; image [D] is the active layer – i.e. live QEMU is writing to it. (The rule of thumb is: live QEMU will always be pointing to the rightmost image in a disk image chain.)

The above image chain can be created by invoking `blockdev-snapshot-sync` commands as following (which shows the creation of overlay image [B]) using the `qmp-shell` (our invocation also prints the raw JSON invocation of it):

```
(QEMU) blockdev-snapshot-sync node-name=node-A snapshot-file=b.qcow2 snapshot-node-
↪name=node-B format=qcow2
{
  "execute": "blockdev-snapshot-sync",
  "arguments": {
    "node-name": "node-A",
    "snapshot-file": "b.qcow2",
    "format": "qcow2",
    "snapshot-node-name": "node-B"
  }
}
```

Here, “node-A” is the name QEMU internally uses to refer to the base image [A] – it is the backing file, based on which the overlay image, [B], is created.

To create the rest of the overlay images, [C], and [D] (omitting the raw JSON output for brevity):

```
(QEMU) blockdev-snapshot-sync node-name=node-B snapshot-file=c.qcow2 snapshot-node-
↪name=node-C format=qcow2
(QEMU) blockdev-snapshot-sync node-name=node-C snapshot-file=d.qcow2 snapshot-node-
↪name=node-D format=qcow2
```


4.4.5 A note on points-in-time vs file names

In our disk image chain:

```
[A] <-- [B] <-- [C] <-- [D]
```

We have *three* points in time and an active layer:

- Point 1: Guest state when [B] was created is contained in file [A]
- Point 2: Guest state when [C] was created is contained in [A] + [B]
- Point 3: Guest state when [D] was created is contained in [A] + [B] + [C]
- Active layer: Current guest state is contained in [A] + [B] + [C] + [D]

Therefore, be aware with naming choices:

- Naming a file after the time it is created is misleading – the guest data for that point in time is *not* contained in that file (as explained earlier)
- Rather, think of files as a *delta* from the backing file

4.4.6 Live block streaming — `block-stream`

The `block-stream` command allows you to do live copy data from backing files into overlay images.

Given our original example disk image chain from earlier:

```
[A] <-- [B] <-- [C] <-- [D]
```

The disk image chain can be shortened in one of the following different ways (not an exhaustive list).

- (1) Merge everything into the active layer: I.e. copy all contents from the base image, [A], and overlay images, [B] and [C], into [D], *while* the guest is running. The resulting chain will be a standalone image, [D] – with contents from [A], [B] and [C] merged into it (where live QEMU writes go to):

```
[D]
```

- (2) Taking the same example disk image chain mentioned earlier, merge only images [B] and [C] into [D], the active layer. The result will be contents of images [B] and [C] will be copied into [D], and the backing file pointer of image [D] will be adjusted to point to image [A]. The resulting chain will be:

```
[A] <-- [D]
```

- (3) Intermediate streaming (available since QEMU 2.8): Starting afresh with the original example disk image chain, with a total of four images, it is possible to copy contents from image [B] into image [C]. Once the copy is finished, image [B] can now be (optionally) discarded; and the backing file pointer of image [C] will be adjusted to point to [A]. I.e. after performing “intermediate streaming” of [B] into [C], the resulting image chain will be (where live QEMU is writing to [D]):

```
[A] <-- [C] <-- [D]
```

QMP invocation for `block-stream`

For *Case-1*, to merge contents of all the backing files into the active layer, where ‘node-D’ is the current active image (by default `block-stream` will flatten the entire chain); `qmp-shell` (and its corresponding JSON output):

```
(QEMU) block-stream device=node-D job-id=job0
{
  "execute": "block-stream",
  "arguments": {
    "device": "node-D",
    "job-id": "job0"
  }
}
```

For *Case-2*, merge contents of the images [B] and [C] into [D], where image [D] ends up referring to image [A] as its backing file:

```
(QEMU) block-stream device=node-D base-node=node-A job-id=job0
```

And for *Case-3*, of “intermediate” streaming”, merge contents of images [B] into [C], where [C] ends up referring to [A] as its backing image:

```
(QEMU) block-stream device=node-C base-node=node-A job-id=job0
```

Progress of a `block-stream` operation can be monitored via the QMP command:

```
(QEMU) query-block-jobs
{
  "execute": "query-block-jobs",
  "arguments": {}
}
```

Once the `block-stream` operation has completed, QEMU will emit an event, `BLOCK_JOB_COMPLETED`. The intermediate overlays remain valid, and can now be (optionally) discarded, or retained to create further overlays based on them. Finally, the `block-stream` jobs can be restarted at anytime.

4.4.7 Live block commit — `block-commit`

The `block-commit` command lets you merge live data from overlay images into backing file(s). Since QEMU 2.0, this includes “live active commit” (i.e. it is possible to merge the “active layer”, the right-most image in a disk image chain where live QEMU will be writing to, into the base image). This is analogous to `block-stream`, but in the opposite direction.

Again, starting afresh with our example disk image chain, where live QEMU is writing to the right-most image in the chain, [D]:

```
[A] <-- [B] <-- [C] <-- [D]
```

The disk image chain can be shortened in one of the following ways:

- (1) Commit content from only image [B] into image [A]. The resulting chain is the following, where image [C] is adjusted to point at [A] as its new backing file:

```
[A] <-- [C] <-- [D]
```

- (2) Commit content from images [B] and [C] into image [A]. The resulting chain, where image [D] is adjusted to point to image [A] as its new backing file:

```
[A] <-- [D]
```

- (3) Commit content from images [B], [C], and the active layer [D] into image [A]. The resulting chain (in this case, a consolidated single image):

```
[A]
```

- (4) Commit content from image only image [C] into image [B]. The resulting chain:

```
[A] <-- [B] <-- [D]
```

- (5) Commit content from image [C] and the active layer [D] into image [B]. The resulting chain:

```
[A] <-- [B]
```

QMP invocation for `block-commit`

For *Case-1*, to merge contents only from image [B] into image [A], the invocation is as follows:

```
(QEMU) block-commit device=node-D base=a.qcow2 top=b.qcow2 job-id=job0
{
  "execute": "block-commit",
  "arguments": {
    "device": "node-D",
    "job-id": "job0",
    "top": "b.qcow2",
    "base": "a.qcow2"
  }
}
```

Once the above `block-commit` operation has completed, a `BLOCK_JOB_COMPLETED` event will be issued, and no further action is required. As the end result, the backing file of image [C] is adjusted to point to image [A], and the original 4-image chain will end up being transformed to:

```
[A] <-- [C] <-- [D]
```

Note: The intermediate image [B] is invalid (as in: no more further overlays based on it can be created).

Reasoning: An intermediate image after a ‘stream’ operation still represents that old point-in-time, and may be valid in that context. However, an intermediate image after a ‘commit’ operation no longer represents any point-in-time, and is invalid in any context.

However, *Case-3* (also called: “active `block-commit`”) is a *two-phase* operation: In the first phase, the content from the active overlay, along with the intermediate overlays, is copied into the backing file (also called the base image). In the second phase, adjust the said backing file as the current active image – possible via issuing the command `block-job-complete`. Optionally, the `block-commit` operation can be cancelled by issuing the command `block-job-cancel`, but be careful when doing this.

Once the `block-commit` operation has completed, the event `BLOCK_JOB_READY` will be emitted, signalling that the synchronization has finished. Now the job can be gracefully completed by issuing the command `block-job-complete` – until such a command is issued, the ‘commit’ operation remains active.

The following is the flow for *Case-3* to convert a disk image chain such as this:

```
[A] <-- [B] <-- [C] <-- [D]
```

Into:

[A]

Where content from all the subsequent overlays, [B], and [C], including the active layer, [D], is committed back to [A] – which is where live QEMU is performing all its current writes).

Start the “active block-commit” operation:

```
(QEMU) block-commit device=node-D base=a.qcow2 top=d.qcow2 job-id=job0
{
  "execute": "block-commit",
  "arguments": {
    "device": "node-D",
    "job-id": "job0",
    "top": "d.qcow2",
    "base": "a.qcow2"
  }
}
```

Once the synchronization has completed, the event `BLOCK_JOB_READY` will be emitted.

Then, optionally query for the status of the active block operations. We can see the ‘commit’ job is now ready to be completed, as indicated by the line “*ready*”: *true*:

```
(QEMU) query-block-jobs
{
  "execute": "query-block-jobs",
  "arguments": {}
}
{
  "return": [
    {
      "busy": false,
      "type": "commit",
      "len": 1376256,
      "paused": false,
      "ready": true,
      "io-status": "ok",
      "offset": 1376256,
      "device": "job0",
      "speed": 0
    }
  ]
}
```

Gracefully complete the ‘commit’ block device job:

```
(QEMU) block-job-complete device=job0
{
  "execute": "block-job-complete",
  "arguments": {
    "device": "job0"
  }
}
{
  "return": {}
}
```

Finally, once the above job is completed, an event `BLOCK_JOB_COMPLETED` will be emitted.

Note: The invocation for rest of the cases (2, 4, and 5), discussed in the previous section, is omitted for brevity.

4.4.8 Live disk synchronization — `drive-mirror` and `blockdev-mirror`

Synchronize a running disk image chain (all or part of it) to a target image.

Again, given our familiar disk image chain:

```
[A] <-- [B] <-- [C] <-- [D]
```

The `drive-mirror` (and its newer equivalent `blockdev-mirror`) allows you to copy data from the entire chain into a single target image (which can be located on a different host), [E].

Note: When you cancel an in-progress ‘mirror’ job *before* the source and target are synchronized, `block-job-cancel` will emit the event `BLOCK_JOB_CANCELLED`. However, note that if you cancel a ‘mirror’ job *after* it has indicated (via the event `BLOCK_JOB_READY`) that the source and target have reached synchronization, then the event emitted by `block-job-cancel` changes to `BLOCK_JOB_COMPLETED`.

Besides the ‘mirror’ job, the “active `block-commit`” is the only other block device job that emits the event `BLOCK_JOB_READY`. The rest of the block device jobs (‘stream’, “non-active `block-commit`”, and ‘backup’) end automatically.

So there are two possible actions to take, after a ‘mirror’ job has emitted the event `BLOCK_JOB_READY`, indicating that the source and target have reached synchronization:

- (1) Issuing the command `block-job-cancel` (after it emits the event `BLOCK_JOB_COMPLETED`) will create a point-in-time (which is at the time of *triggering* the cancel command) copy of the entire disk image chain (or only the top-most image, depending on the `sync` mode), contained in the target image [E]. One use case for this is live VM migration with non-shared storage.
- (2) Issuing the command `block-job-complete` (after it emits the event `BLOCK_JOB_COMPLETED`) will adjust the guest device (i.e. live QEMU) to point to the target image, [E], causing all the new writes from this point on to happen there.

About synchronization modes: The synchronization mode determines *which* part of the disk image chain will be copied to the target. Currently, there are four different kinds:

- (1) `full` – Synchronize the content of entire disk image chain to the target
- (2) `top` – Synchronize only the contents of the top-most disk image in the chain to the target
- (3) `none` – Synchronize only the new writes from this point on.

Note: In the case of `drive-backup` (or `blockdev-backup`), the behavior of `none` synchronization mode is different. Normally, a `backup` job consists of two parts: Anything that is overwritten by the guest is first copied out to the backup, and in the background the whole image is copied from start to end. With `sync=none`, it’s only the first part.

- (4) `incremental` – Synchronize content that is described by the dirty bitmap

Note: Refer to the [Dirty Bitmaps and Incremental Backup](#) document in the QEMU source tree to learn about the detailed workings of the `incremental` synchronization mode.

QMP invocation for `drive-mirror`

To copy the contents of the entire disk image chain, from [A] all the way to [D], to a new target (`drive-mirror` will create the destination file, if it doesn't already exist), call it [E]:

```
(QEMU) drive-mirror device=node-D target=e.qcow2 sync=full job-id=job0
{
  "execute": "drive-mirror",
  "arguments": {
    "device": "node-D",
    "job-id": "job0",
    "target": "e.qcow2",
    "sync": "full"
  }
}
```

The `"sync": "full"`, from the above, means: copy the *entire* chain to the destination.

Following the above, querying for active block jobs will show that a ‘mirror’ job is “ready” to be completed (and QEMU will also emit an event, `BLOCK_JOB_READY`):

```
(QEMU) query-block-jobs
{
  "execute": "query-block-jobs",
  "arguments": {}
}
{
  "return": [
    {
      "busy": false,
      "type": "mirror",
      "len": 21757952,
      "paused": false,
      "ready": true,
      "io-status": "ok",
      "offset": 21757952,
      "device": "job0",
      "speed": 0
    }
  ]
}
```

And, as noted in the previous section, there are two possible actions at this point:

- (a) Create a point-in-time snapshot by ending the synchronization. The point-in-time is at the time of *ending* the sync. (The result of the following being: the target image, [E], will be populated with content from the entire chain, [A] to [D]):

```
(QEMU) block-job-cancel device=job0
{
  "execute": "block-job-cancel",
  "arguments": {
    "device": "job0"
  }
}
```

- (b) Or, complete the operation and pivot the live QEMU to the target copy:

```
(QEMU) block-job-complete device=job0
```

In either of the above cases, if you once again run the *query-block-jobs* command, there should not be any active block operation.

Comparing ‘commit’ and ‘mirror’: In both then cases, the overlay images can be discarded. However, with ‘commit’, the *existing* base image will be modified (by updating it with contents from overlays); while in the case of ‘mirror’, a *new* target image is populated with the data from the disk image chain.

QMP invocation for live storage migration with `drive-mirror` + NBD

Live storage migration (without shared storage setup) is one of the most common use-cases that takes advantage of the `drive-mirror` primitive and QEMU’s built-in Network Block Device (NBD) server. Here’s a quick walk-through of this setup.

Given the disk image chain:

```
[A] <-- [B] <-- [C] <-- [D]
```

Instead of copying content from the entire chain, synchronize *only* the contents of the *top*-most disk image (i.e. the active layer), [D], to a target, say, [TargetDisk].

Important: The destination host must already have the contents of the backing chain, involving images [A], [B], and [C], visible via other means – whether by `cp`, `rsync`, or by some storage array-specific command.)

Sometimes, this is also referred to as “shallow copy” – because only the “active layer”, and not the rest of the image chain, is copied to the destination.

Note: In this example, for the sake of simplicity, we’ll be using the same `localhost` as both source and destination.

As noted earlier, on the destination host the contents of the backing chain – from images [A] to [C] – are already expected to exist in some form (e.g. in a file called, `Contents-of-A-B-C.qcow2`). Now, on the destination host, let’s create a target overlay image (with the image `Contents-of-A-B-C.qcow2` as its backing file), to which the contents of image [D] (from the source QEMU) will be mirrored to:

```
$ qemu-img create -f qcow2 -b ./Contents-of-A-B-C.qcow2 \
-F qcow2 ./target-disk.qcow2
```

And start the destination QEMU (we already have the source QEMU running – discussed in the section: *Interacting with a QEMU instance*) instance, with the following invocation. (As noted earlier, for simplicity’s sake, the destination QEMU is started on the same host, but it could be located elsewhere):

```
$ ./qemu-system-x86_64 -display none -no-user-config \
-M q35 -nodefaults -m 512 \
-blockdev node-name=node-TargetDisk,driver=qcow2,file.driver=file,file.node-
name=file,file.filename=./target-disk.qcow2 \
-device virtio-blk,drive=node-TargetDisk,id=virtio0 \
-S -monitor stdio -qmp unix:./qmp-sock2,server,nowait \
-incoming tcp:localhost:6666
```

Given the disk image chain on source QEMU:

```
[A] <-- [B] <-- [C] <-- [D]
```

On the destination host, it is expected that the contents of the chain [A] <-- [B] <-- [C] are *already* present, and therefore copy *only* the content of image [D].

- (1) [On *destination* QEMU] As part of the first step, start the built-in NBD server on a given host (local host, represented by ::) and port:

```
(QEMU) nbd-server-start addr={"type":"inet","data":{"host":"::","port":"49153"}}
{
  "execute": "nbd-server-start",
  "arguments": {
    "addr": {
      "data": {
        "host": "::",
        "port": "49153"
      },
      "type": "inet"
    }
  }
}
```

- (2) [On *destination* QEMU] And export the destination disk image using QEMU's built-in NBD server:

```
(QEMU) nbd-server-add device=node-TargetDisk writable=true
{
  "execute": "nbd-server-add",
  "arguments": {
    "device": "node-TargetDisk"
  }
}
```

- (3) [On *source* QEMU] Then, invoke drive-mirror (NB: since we're running drive-mirror with mode=existing (meaning: synchronize to a pre-created file, therefore 'existing', file on the target host), with the synchronization mode as 'top' ("sync: "top")):

```
(QEMU) drive-mirror device=node-D target=nbd:localhost:49153:exportname=node-
→TargetDisk sync=top mode=existing job-id=job0
{
  "execute": "drive-mirror",
  "arguments": {
    "device": "node-D",
    "mode": "existing",
    "job-id": "job0",
    "target": "nbd:localhost:49153:exportname=node-TargetDisk",
    "sync": "top"
  }
}
```

- (4) [On *source* QEMU] Once drive-mirror copies the entire data, and the event BLOCK_JOB_READY is emitted, issue block-job-cancel to gracefully end the synchronization, from source QEMU:

```
(QEMU) block-job-cancel device=job0
{
  "execute": "block-job-cancel",
  "arguments": {
    "device": "job0"
  }
}
```

(continues on next page)

(continued from previous page)

```
}
}
```

(5) [On *destination* QEMU] Then, stop the NBD server:

```
(QEMU) nbd-server-stop
{
  "execute": "nbd-server-stop",
  "arguments": {}
}
```

(6) [On *destination* QEMU] Finally, resume the guest vCPUs by issuing the QMP command *cont*:

```
(QEMU) cont
{
  "execute": "cont",
  "arguments": {}
}
```

Note: Higher-level libraries (e.g. libvirt) automate the entire above process (although note that libvirt does not allow same-host migrations to localhost for other reasons).

Notes on `blockdev-mirror`

The `blockdev-mirror` command is equivalent in core functionality to `drive-mirror`, except that it operates at node-level in a BDS graph.

Also: for `blockdev-mirror`, the ‘target’ image needs to be explicitly created (using `qemu-img`) and attach it to live QEMU via `blockdev-add`, which assigns a name to the to-be created target node.

E.g. the sequence of actions to create a point-in-time backup of an entire disk image chain, to a target, using `blockdev-mirror` would be:

- (0) Create the QCOW2 overlays, to arrive at a backing chain of desired depth
- (1) Create the target image (using `qemu-img`), say, `e.qcow2`
- (2) Attach the above created file (`e.qcow2`), run-time, using `blockdev-add` to QEMU
- (3) Perform `blockdev-mirror` (use `"sync": "full"` to copy the entire chain to the target). And notice the event `BLOCK_JOB_READY`
- (4) Optionally, query for active block jobs, there should be a ‘mirror’ job ready to be completed
- (5) Gracefully complete the ‘mirror’ block device job, and notice the the event `BLOCK_JOB_COMPLETED`
- (6) Shutdown the guest by issuing the QMP `quit` command so that caches are flushed
- (7) Then, finally, compare the contents of the disk image chain, and the target copy with `qemu-img compare`. You should notice: “Images are identical”

QMP invocation for `blockdev-mirror`

Given the disk image chain:

```
[A] <-- [B] <-- [C] <-- [D]
```

To copy the contents of the entire disk image chain, from [A] all the way to [D], to a new target, call it [E]. The following is the flow.

Create the overlay images, [B], [C], and [D]:

```
(QEMU) blockdev-snapshot-sync node-name=node-A snapshot-file=b.qcow2 snapshot-node-  
↪name=node-B format=qcow2  
(QEMU) blockdev-snapshot-sync node-name=node-B snapshot-file=c.qcow2 snapshot-node-  
↪name=node-C format=qcow2  
(QEMU) blockdev-snapshot-sync node-name=node-C snapshot-file=d.qcow2 snapshot-node-  
↪name=node-D format=qcow2
```

Create the target image, [E]:

```
$ qemu-img create -f qcow2 e.qcow2 39M
```

Add the above created target image to QEMU, via blockdev-add:

```
(QEMU) blockdev-add driver=qcow2 node-name=node-E file={"driver":"file","filename":"e.  
↪qcow2"}  
{  
  "execute": "blockdev-add",  
  "arguments": {  
    "node-name": "node-E",  
    "driver": "qcow2",  
    "file": {  
      "driver": "file",  
      "filename": "e.qcow2"  
    }  
  }  
}
```

Perform blockdev-mirror, and notice the event BLOCK_JOB_READY:

```
(QEMU) blockdev-mirror device=node-B target=node-E sync=full job-id=job0  
{  
  "execute": "blockdev-mirror",  
  "arguments": {  
    "device": "node-D",  
    "job-id": "job0",  
    "target": "node-E",  
    "sync": "full"  
  }  
}
```

Query for active block jobs, there should be a ‘mirror’ job ready:

```
(QEMU) query-block-jobs  
{  
  "execute": "query-block-jobs",  
  "arguments": {}  
}  
{  
  "return": [  
    {
```

(continues on next page)

(continued from previous page)

```

        "busy": false,
        "type": "mirror",
        "len": 21561344,
        "paused": false,
        "ready": true,
        "io-status": "ok",
        "offset": 21561344,
        "device": "job0",
        "speed": 0
    }
}

```

Gracefully complete the block device job operation, and notice the event `BLOCK_JOB_COMPLETED`:

```

(QEMU) block-job-complete device=job0
{
    "execute": "block-job-complete",
    "arguments": {
        "device": "job0"
    }
}
{
    "return": {}
}

```

Shutdown the guest, by issuing the `quit` QMP command:

```

(QEMU) quit
{
    "execute": "quit",
    "arguments": {}
}

```

4.4.9 Live disk backup — `drive-backup` and `blockdev-backup`

The `drive-backup` (and its newer equivalent `blockdev-backup`) allows you to create a point-in-time snapshot. In this case, the point-in-time is when you *start* the `drive-backup` (or its newer equivalent `blockdev-backup`) command.

QMP invocation for `drive-backup`

Yet again, starting afresh with our example disk image chain:

```
[A] <-- [B] <-- [C] <-- [D]
```

To create a target image [E], with content populated from image [A] to [D], from the above chain, the following is the syntax. (If the target image does not exist, `drive-backup` will create it):

```

(QEMU) drive-backup device=node-D sync=full target=e.qcow2 job-id=job0
{
    "execute": "drive-backup",
    "arguments": {

```

(continues on next page)

(continued from previous page)

```

        "device": "node-D",
        "job-id": "job0",
        "sync": "full",
        "target": "e.qcow2"
    }
}

```

Once the above drive-backup has completed, a `BLOCK_JOB_COMPLETED` event will be issued, indicating the live block device job operation has completed, and no further action is required.

Notes on blockdev-backup

The `blockdev-backup` command is equivalent in functionality to `drive-backup`, except that it operates at node-level in a Block Driver State (BDS) graph.

E.g. the sequence of actions to create a point-in-time backup of an entire disk image chain, to a target, using `blockdev-backup` would be:

- (0) Create the QCOW2 overlays, to arrive at a backing chain of desired depth
- (1) Create the target image (using `qemu-img`), say, `e.qcow2`
- (2) Attach the above created file (`e.qcow2`), run-time, using `blockdev-add` to QEMU
- (3) Perform `blockdev-backup` (use `"sync": "full"` to copy the entire chain to the target). And notice the event `BLOCK_JOB_COMPLETED`
- (4) Shutdown the guest, by issuing the QMP `quit` command, so that caches are flushed
- (5) Then, finally, compare the contents of the disk image chain, and the target copy with `qemu-img compare`. You should notice: "Images are identical"

The following section shows an example QMP invocation for `blockdev-backup`.

QMP invocation for blockdev-backup

Given a disk image chain of depth 1 where image [B] is the active overlay (live QEMU is writing to it):

```
[A] <-- [B]
```

The following is the procedure to copy the content from the entire chain to a target image (say, [E]), which has the full content from [A] and [B].

Create the overlay [B]:

```

(QEMU) blockdev-snapshot-sync node-name=node-A snapshot-file=b.qcow2 snapshot-node-
↪name=node-B format=qcow2
{
    "execute": "blockdev-snapshot-sync",
    "arguments": {
        "node-name": "node-A",
        "snapshot-file": "b.qcow2",
        "format": "qcow2",
        "snapshot-node-name": "node-B"
    }
}

```

Create a target image that will contain the copy:

```
$ qemu-img create -f qcow2 e.qcow2 39M
```

Then add it to QEMU via `blockdev-add`:

```
(QEMU) blockdev-add driver=qcow2 node-name=node-E file={"driver":"file","filename":"e.
↪qcow2"}
{
  "execute": "blockdev-add",
  "arguments": {
    "node-name": "node-E",
    "driver": "qcow2",
    "file": {
      "driver": "file",
      "filename": "e.qcow2"
    }
  }
}
```

Then invoke `blockdev-backup` to copy the contents from the entire image chain, consisting of images [A] and [B] to the target image ‘e.qcow2’:

```
(QEMU) blockdev-backup device=node-B target=node-E sync=full job-id=job0
{
  "execute": "blockdev-backup",
  "arguments": {
    "device": "node-B",
    "job-id": "job0",
    "target": "node-E",
    "sync": "full"
  }
}
```

Once the above ‘backup’ operation has completed, the event, `BLOCK_JOB_COMPLETED` will be emitted, signalling successful completion.

Next, query for any active block device jobs (there should be none):

```
(QEMU) query-block-jobs
{
  "execute": "query-block-jobs",
  "arguments": {}
}
```

Shutdown the guest:

```
(QEMU) quit
{
  "execute": "quit",
  "arguments": {}
}

"return": {}
}
```

Note: The above step is really important; if forgotten, an error, “Failed to get shared “write” lock on e.qcow2”, will be thrown when you do `qemu-img compare` to verify the integrity of the disk image with the backup content.

The end result will be the image ‘e.qcow2’ containing a point-in-time backup of the disk image chain – i.e. contents from images [A] and [B] at the time the `blockdev-backup` command was initiated.

One way to confirm the backup disk image contains the identical content with the disk image chain is to compare the backup and the contents of the chain, you should see “Images are identical”. (NB: this is assuming QEMU was launched with `-S` option, which will not start the CPUs at guest boot up):

```
$ qemu-img compare b.qcow2 e.qcow2
Warning: Image size mismatch!
Images are identical.
```

NOTE: The “Warning: Image size mismatch!” is expected, as we created the target image (e.qcow2) with 39M size.

4.5 Persistent reservation helper protocol

QEMU’s SCSI passthrough devices, `scsi-block` and `scsi-generic`, can delegate implementation of persistent reservations to an external (and typically privileged) program. Persistent Reservations allow restricting access to block devices to specific initiators in a shared storage setup.

For a more detailed reference please refer to the SCSI Primary Commands standard, specifically the section on Reservations and the “PERSISTENT RESERVE IN” and “PERSISTENT RESERVE OUT” commands.

This document describes the socket protocol used between QEMU’s `pr-manager-helper` object and the external program.

Contents

- *Persistent reservation helper protocol*
 - *Connection and initialization*
 - *Command format*

4.5.1 Connection and initialization

All data transmitted on the socket is big-endian.

After connecting to the helper program’s socket, the helper starts a simple feature negotiation process by writing four bytes corresponding to the features it exposes (`supported_features`). QEMU reads it, then writes four bytes corresponding to the desired features of the helper program (`requested_features`).

If a bit is 1 in `requested_features` and 0 in `supported_features`, the corresponding feature is not supported by the helper and the connection is closed. On the other hand, it is acceptable for a bit to be 0 in `requested_features` and 1 in `supported_features`; in this case, the helper will not enable the feature.

Right now no feature is defined, so the two parties always write four zero bytes.

4.5.2 Command format

It is invalid to send multiple commands concurrently on the same socket. It is however possible to connect multiple sockets to the helper and send multiple commands to the helper for one or more file descriptors.

A command consists of a request and a response. A request consists of a 16-byte SCSI CDB. A file descriptor must be passed to the helper together with the SCSI CDB using ancillary data.

The CDB has the following limitations:

- the command (stored in the first byte) must be one of 0x5E (PERSISTENT RESERVE IN) or 0x5F (PERSISTENT RESERVE OUT).
- the allocation length (stored in bytes 7-8 of the CDB for PERSISTENT RESERVE IN) or parameter list length (stored in bytes 5-8 of the CDB for PERSISTENT RESERVE OUT) is limited to 8 KiB.

For PERSISTENT RESERVE OUT, the parameter list is sent right after the CDB. The length of the parameter list is taken from the CDB itself.

The helper's reply has the following structure:

- 4 bytes for the SCSI status
- 4 bytes for the payload size (nonzero only for PERSISTENT RESERVE IN and only if the SCSI status is 0x00, i.e. GOOD)
- 96 bytes for the SCSI sense data
- if the size is nonzero, the payload follows

The sense data is always sent to keep the protocol simple, even though it is only valid if the SCSI status is CHECK CONDITION (0x02).

The payload size is always less than or equal to the allocation length specified in the CDB for the PERSISTENT RESERVE IN command.

If the protocol is violated, the helper closes the socket.

4.6 QEMU Guest Agent

4.6.1 Synopsis

qemu-ga [*OPTIONS*]

4.6.2 Description

The QEMU Guest Agent is a daemon intended to be run within virtual machines. It allows the hypervisor host to perform various operations in the guest, such as:

- get information from the guest
- set the guest's system time
- read/write a file
- sync and freeze the filesystems
- suspend the guest
- reconfigure guest local processors
- set user's password
- ...

qemu-ga will read a system configuration file on startup (located at `/etc/qemu/qemu-ga.conf` by default), then parse remaining configuration options on the command line. For the same key, the last option wins, but the lists accumulate (see below for configuration file format).

4.6.3 Options

- m, --method=METHOD**
Transport method: one of `unix-listen`, `virtio-serial`, or `isa-serial`, or `vsock-listen` (`virtio-serial` is the default).
- p, --path=PATH**
Device/socket path (the default for `virtio-serial` is `/dev/virtio-ports/org.qemu.guest_agent.0`, the default for `isa-serial` is `/dev/ttyS0`). Socket addresses for `vsock-listen` are written as `<cid>:<port>`.
- l, --logfile=PATH**
Set log file path (default is `stderr`).
- f, --pidfile=PATH**
Specify pid file (default is `/var/run/qemu-ga.pid`).
- F, --fsfreeze-hook=PATH**
Enable fsfreeze hook. Accepts an optional argument that specifies script to run on freeze/thaw. Script will be called with 'freeze'/'thaw' arguments accordingly (default is `/etc/qemu/fsfreeze-hook`). If using `-F` with an argument, do not follow `-F` with a space (for example: `-F/var/run/fsfreezehook.sh`).
- t, --statedir=PATH**
Specify the directory to store state information (absolute paths only, default is `/var/run`).
- v, --verbose**
Log extra debugging information.
- V, --version**
Print version information and exit.
- d, --daemon**
Daemonize after startup (detach from terminal).
- b, --blacklist=LIST**
Comma-separated list of RPCs to disable (no spaces, ? to list available RPCs).
- D, --dump-conf**
Dump the configuration in a format compatible with `qemu-ga.conf` and exit.
- h, --help**
Display this help and exit.

4.6.4 Files

The syntax of the `qemu-ga.conf` configuration file follows the Desktop Entry Specification, here is a quick summary: it consists of groups of key-value pairs, interspersed with comments.

```
# qemu-ga configuration sample
[general]
daemonize = 0
pidfile = /var/run/qemu-ga.pid
verbose = 0
method = virtio-serial
path = /dev/virtio-ports/org.qemu.guest_agent.0
statedir = /var/run
```

The list of keys follows the command line options:

Key	Key type
daemon	boolean
method	string
path	string
logfile	string
pidfile	string
fsfreeze-hook	string
statedir	string
verbose	boolean
blacklist	string list

4.6.5 See also

qemu (1)

4.7 QEMU Guest Agent Protocol Reference

4.7.1 General note concerning the use of guest agent interfaces

“unsupported” is a higher-level error than the errors that individual commands might document. The caller should always be prepared to receive QERR_UNSUPPORTED, even if the given command doesn’t specify it, or doesn’t document any failure mode at all.

4.7.2 QEMU guest agent protocol commands and structs

guest-sync-delimited (Command)

Echo back a unique integer value, and prepend to response a leading sentinel byte (0xFF) the client can check scan for.

This is used by clients talking to the guest agent over the wire to ensure the stream is in sync and doesn’t contain stale data from previous client. It must be issued upon initial connection, and after any client-side timeouts (including timeouts on receiving a response to this command).

After issuing this request, all guest agent responses should be ignored until the response containing the unique integer value the client passed in is returned. Receipt of the 0xFF sentinel byte must be handled as an indication that the client’s lexer/tokenizer/parser state should be flushed/reset in preparation for reliably receiving the subsequent response. As an optimization, clients may opt to ignore all data until a sentinel value is receiving to avoid unnecessary processing of stale data.

Similarly, clients should also precede this *request* with a 0xFF byte to make sure the guest agent flushes any partially read JSON data from a previous client connection.

Arguments

id: **int** randomly generated 64-bit integer

Returns

The unique integer id passed in by the client

Since

1.1

guest-sync (Command)

Echo back a unique integer value

This is used by clients talking to the guest agent over the wire to ensure the stream is in sync and doesn't contain stale data from previous client. All guest agent responses should be ignored until the provided unique integer value is returned, and it is up to the client to handle stale whole or partially-delivered JSON text in such a way that this response can be obtained.

In cases where a partial stale response was previously received by the client, this cannot always be done reliably. One particular scenario being if qemu-ga responses are fed character-by-character into a JSON parser. In these situations, using guest-sync-delimited may be optimal.

For clients that fetch responses line by line and convert them to JSON objects, guest-sync should be sufficient, but note that in cases where the channel is dirty some attempts at parsing the response may result in a parser error.

Such clients should also precede this command with a 0xFF byte to make sure the guest agent flushes any partially read JSON data from a previous session.

Arguments

id: int randomly generated 64-bit integer

Returns

The unique integer id passed in by the client

Since

0.15.0

guest-ping (Command)

Ping the guest agent, a non-error return implies success

Since

0.15.0

`guest-get-time` (Command)

Get the information about guest's System Time relative to the Epoch of 1970-01-01 in UTC.

Returns

Time in nanoseconds.

Since

1.5

`guest-set-time` (Command)

Set guest time.

When a guest is paused or migrated to a file then loaded from that file, the guest OS has no idea that there was a big gap in the time. Depending on how long the gap was, NTP might not be able to resynchronize the guest.

This command tries to set guest's System Time to the given value, then sets the Hardware Clock (RTC) to the current System Time. This will make it easier for a guest to resynchronize without waiting for NTP. If no `time` is specified, then the time to set is read from RTC. However, this may not be supported on all platforms (i.e. Windows). If that's the case users are advised to always pass a value.

Arguments

time: int (optional) time of nanoseconds, relative to the Epoch of 1970-01-01 in UTC.

Returns

Nothing on success.

Since

1.5

`GuestAgentCommandInfo` (Object)

Information about guest agent commands.

Members

name: string name of the command

enabled: boolean whether command is currently enabled by guest admin

success-response: boolean whether command returns a response on success (since 1.7)

Since

1.1.0

GuestAgentInfo (Object)

Information about guest agent.

Members

version: string guest agent version

supported_commands: array of GuestAgentCommandInfo Information about guest agent commands

Since

0.15.0

guest-info (Command)

Get some information about the guest agent.

Returns

GuestAgentInfo

Since

0.15.0

guest-shutdown (Command)

Initiate guest-activated shutdown. Note: this is an asynchronous shutdown request, with no guarantee of successful shutdown.

Arguments

mode: string (optional) “halt”, “powerdown” (default), or “reboot”

This command does NOT return a response on success. Success condition is indicated by the VM exiting with a zero exit status or, when running with `-no-shutdown`, by issuing the `query-status QMP` command to confirm the VM status is “shutdown”.

Since

0.15.0

guest-file-open (Command)

Open a file in the guest and retrieve a file handle for it

Arguments

path: string Full path to the file in the guest to open.

mode: string (optional) open mode, as per fopen(), “r” is the default.

Returns

Guest file handle on success.

Since

0.15.0

guest-file-close (Command)

Close an open file in the guest

Arguments

handle: int filehandle returned by guest-file-open

Returns

Nothing on success.

Since

0.15.0

GuestFileRead (Object)

Result of guest agent file-read operation

Members

count: int number of bytes read (note: count is *before* base64-encoding is applied)

buf-b64: string base64-encoded bytes read

eof: boolean whether EOF was encountered during read operation.

Since

0.15.0

`guest-file-read` (Command)

Read from an open file in the guest. Data will be base64-encoded. As this command is just for limited, ad-hoc debugging, such as log file access, the number of bytes to read is limited to 48 MB.

Arguments

handle: int filehandle returned by `guest-file-open`

count: int (optional) maximum number of bytes to read (default is 4KB, maximum is 48MB)

Returns

`GuestFileRead` on success.

Since

0.15.0

`GuestFileWrite` (Object)

Result of guest agent file-write operation

Members

count: int number of bytes written (note: count is actual bytes written, after base64-decoding of provided buffer)

eof: boolean whether EOF was encountered during write operation.

Since

0.15.0

`guest-file-write` (Command)

Write to an open file in the guest.

Arguments

handle: int filehandle returned by `guest-file-open`

buf-b64: string base64-encoded string representing data to be written

count: int (optional) bytes to write (actual bytes, after base64-decode), default is all content in `buf-b64` buffer after base64 decoding

Returns

`GuestFileWrite` on success.

Since

0.15.0

`GuestFileSeek` (Object)

Result of guest agent file-`seek` operation

Members

position: int current file position

eof: boolean whether EOF was encountered during file seek

Since

0.15.0

`QGASseek` (Enum)

Symbolic names for use in `guest-file-seek`

Values

set Set to the specified offset (same effect as `'whence':0`)

cur Add offset to the current location (same effect as `'whence':1`)

end Add offset to the end of the file (same effect as `'whence':2`)

Since

2.6

GuestFileWhence (Alternate)

Controls the meaning of offset to `guest-file-seek`.

Members

value: `int` Integral value (0 for set, 1 for cur, 2 for end), available for historical reasons, and might differ from the host's or guest's `SEEK_*` values (since: 0.15)

name: `QGASseek` Symbolic name, and preferred interface

Since

2.6

guest-file-seek (Command)

Seek to a position in the file, as with `fseek()`, and return the current file position afterward. Also encapsulates `ftell()`'s functionality, with `offset=0` and `whence=1`.

Arguments

handle: `int` filehandle returned by `guest-file-open`

offset: `int` bytes to skip over in the file stream

whence: `GuestFileWhence` Symbolic or numeric code for interpreting offset

Returns

`GuestFileSeek` on success.

Since

0.15.0

guest-file-flush (Command)

Write file changes buffered in userspace to disk/kernel buffers

Arguments

handle: `int` filehandle returned by `guest-file-open`

Returns

Nothing on success.

Since

0.15.0

GuestFsfreezeStatus (Enum)

An enumeration of filesystem freeze states

Values

thawed filesystems thawed/unfrozen

frozen all non-network guest filesystems frozen

Since

0.15.0

guest-fsfreeze-status (Command)

Get guest fsfreeze state. error state indicates

Returns

GuestFsfreezeStatus (“thawed”, “frozen”, etc., as defined below)

Note

This may fail to properly report the current state as a result of some other guest processes having issued an fs freeze/thaw.

Since

0.15.0

guest-fsfreeze-freeze (Command)

Sync and freeze all freezable, local guest filesystems. If this command succeeded, you may call `guest-fsfreeze-thaw` later to unfreeze.

Note

On Windows, the command is implemented with the help of a Volume Shadow-copy Service DLL helper. The frozen state is limited for up to 10 seconds by VSS.

Returns

Number of file systems currently frozen. On error, all filesystems will be thawed. If no filesystems are frozen as a result of this call, then `guest-fsfreeze-status` will remain “thawed” and calling `guest-fsfreeze-thaw` is not necessary.

Since

0.15.0

`guest-fsfreeze-freeze-list` (Command)

Sync and freeze specified guest filesystems. See also `guest-fsfreeze-freeze`.

Arguments

mountpoints: array of string (optional) an array of mountpoints of filesystems to be frozen. If omitted, every mounted filesystem is frozen. Invalid mount points are ignored.

Returns

Number of file systems currently frozen. On error, all filesystems will be thawed.

Since

2.2

`guest-fsfreeze-thaw` (Command)

Unfreeze all frozen guest filesystems

Returns

Number of file systems thawed by this call

Note

if return value does not match the previous call to `guest-fsfreeze-freeze`, this likely means some freezable filesystems were unfrozen before this call, and that the filesystem state may have changed before issuing this command.

Since

0.15.0

GuestFilesystemTrimResult (Object)

Members

path: `string` path that was trimmed

error: `string` (optional) an error message when trim failed

trimmed: `int` (optional) bytes trimmed for this path

minimum: `int` (optional) reported effective minimum for this path

Since

2.4

GuestFilesystemTrimResponse (Object)

Members

paths: `array of GuestFilesystemTrimResult` list of `GuestFilesystemTrimResult` per path that was trimmed

Since

2.4

guest-fstrim (Command)

Discard (or “trim”) blocks which are not in use by the filesystem.

Arguments

minimum: `int` (optional) Minimum contiguous free range to discard, in bytes. Free ranges smaller than this may be ignored (this is a hint and the guest may not respect it). By increasing this value, the fstrim operation will complete more quickly for filesystems with badly fragmented free space, although not all blocks will be discarded. The default value is zero, meaning “discard every free block”.

Returns

A `GuestFilesystemTrimResponse` which contains the status of all trimmed paths. (since 2.4)

Since

1.2

guest-suspend-disk (Command)

Suspend guest to disk.

This command attempts to suspend the guest using three strategies, in this order:

- systemd hibernate
- pm-utils (via pm-hibernate)
- manual write into sysfs

This command does NOT return a response on success. There is a high chance the command succeeded if the VM exits with a zero exit status or, when running with `--no-shutdown`, by issuing the `query-status QMP` command to confirm the VM status is “shutdown”. However, the VM could also exit (or set its status to “shutdown”) due to other reasons.

The following errors may be returned:

- If suspend to disk is not supported, `Unsupported`

Notes

It’s strongly recommended to issue the `guest-sync` command before sending commands when the guest resumes

Since

1.1

guest-suspend-ram (Command)

Suspend guest to ram.

This command attempts to suspend the guest using three strategies, in this order:

- systemd suspend
- pm-utils (via pm-suspend)
- manual write into sysfs

IMPORTANT: `guest-suspend-ram` requires working wakeup support in QEMU. You should check `QMP` command `query-current-machine` returns `wakeup-suspend-support: true` before issuing this command. Failure in doing so can result in a suspended guest that QEMU will not be able to awaken, forcing the user to power cycle the guest to bring it back.

This command does NOT return a response on success. There are two options to check for success:

1. Wait for the `SUSPEND QMP` event from QEMU
2. Issue the `query-status QMP` command to confirm the VM status is “suspended”

The following errors may be returned:

- If suspend to ram is not supported, `Unsupported`

Notes

It’s strongly recommended to issue the `guest-sync` command before sending commands when the guest resumes

Since

1.1

guest-suspend-hybrid (Command)

Save guest state to disk and suspend to ram.

This command attempts to suspend the guest by executing, in this order:

- systemd hybrid-sleep
- pm-utils (via pm-suspend-hybrid)

IMPORTANT: guest-suspend-hybrid requires working wakeup support in QEMU. You should check QMP command query-current-machine returns wakeup-suspend-support: true before issuing this command. Failure in doing so can result in a suspended guest that QEMU will not be able to awaken, forcing the user to power cycle the guest to bring it back.

This command does NOT return a response on success. There are two options to check for success:

1. Wait for the SUSPEND QMP event from QEMU
2. Issue the query-status QMP command to confirm the VM status is “suspended”

The following errors may be returned:

- If hybrid suspend is not supported, Unsupported

Notes

It’s strongly recommended to issue the guest-sync command before sending commands when the guest resumes

Since

1.1

GuestIpAddressType (Enum)

An enumeration of supported IP address types

Values

ipv4 IP version 4

ipv6 IP version 6

Since

1.1

GuestIpAddress (Object)

Members

ip-address: string IP address

ip-address-type: GuestIpAddressType Type of ip-address (e.g. ipv4, ipv6)

prefix: int Network prefix length of ip-address

Since

1.1

GuestNetworkInterfaceStat (Object)

Members

rx-bytes: int total bytes received

rx-packets: int total packets received

rx-errs: int bad packets received

rx-dropped: int receiver dropped packets

tx-bytes: int total bytes transmitted

tx-packets: int total packets transmitted

tx-errs: int packet transmit problems

tx-dropped: int dropped packets transmitted

Since

2.11

GuestNetworkInterface (Object)

Members

name: string The name of interface for which info are being delivered

hardware-address: string (optional) Hardware address of name

ip-addresses: array of GuestIpAddress (optional) List of addresses assigned to name

statistics: GuestNetworkInterfaceStat (optional) various statistic counters related to name (since 2.11)

Since

1.1

guest-network-get-interfaces (Command)

Get list of guest IP addresses, MAC addresses and netmasks.

Returns

List of GuestNetworkInfo on success.

Since

1.1

GuestLogicalProcessor (Object)**Members**

logical-id: int Arbitrary guest-specific unique identifier of the VCPU.

online: boolean Whether the VCPU is enabled.

can-offline: boolean (optional) Whether offlining the VCPU is possible. This member is always filled in by the guest agent when the structure is returned, and always ignored on input (hence it can be omitted then).

Since

1.5

guest-get-vcpus (Command)

Retrieve the list of the guest's logical processors.

This is a read-only operation.

Returns

The list of all VCPUs the guest knows about. Each VCPU is put on the list exactly once, but their order is unspecified.

Since

1.5

guest-set-vcpus (Command)

Attempt to reconfigure (currently: enable/disable) logical processors inside the guest.

The input list is processed node by node in order. In each node `logical-id` is used to look up the guest VCPU, for which `online` specifies the requested state. The set of distinct `logical-id`'s is only required to be a subset of the guest-supported identifiers. There's no restriction on list length or on repeating the same `logical-id` (with possibly

different `online` field). Preferably the input list should describe a modified subset of `guest-get-vcpus`' return value.

Arguments

vcpus: array of `GuestLogicalProcessor` Not documented

Returns

The length of the initial sublist that has been successfully processed. The guest agent maximizes this value. Possible cases:

- 0: if the `vcpus` list was empty on input. Guest state has not been changed. Otherwise,
- Error: processing the first node of `vcpus` failed for the reason returned. Guest state has not been changed. Otherwise,
- `< length(vcpus)`: more than zero initial nodes have been processed, but not the entire `vcpus` list. Guest state has changed accordingly. To retrieve the error (assuming it persists), repeat the call with the successfully processed initial sublist removed. Otherwise,
- `length(vcpus)`: call successful.

Since

1.5

GuestDiskBusType (Enum)

An enumeration of bus type of disks

Values

ide IDE disks

fdc floppy disks

scsi SCSI disks

virtio virtio disks

xen Xen disks

usb USB disks

uml UML disks

sata SATA disks

sd SD cards

unknown Unknown bus type

ieee1394 Win IEEE 1394 bus type

ssa Win SSA bus type

fibre Win fiber channel bus type
raid Win RAID bus type
iscsi Win iScsi bus type
sas Win serial-attaches SCSI bus type
mmc Win multimedia card (MMC) bus type
virtual Win virtual bus type
file-backed-virtual Win file-backed bus type

Since

2.2; ‘Unknown’ and all entries below since 2.4

GuestPCIAddress (Object)

Members

domain: int domain id
bus: int bus id
slot: int slot id
function: int function id

Since

2.2

GuestCCWAddress (Object)

Members

cssid: int channel subsystem image id
ssid: int subchannel set id
subchno: int subchannel number
devno: int device number

Since

6.0

GuestDiskAddress (Object)

Members

pci-controller: GuestPCIAddress controller's PCI address (fields are set to -1 if invalid)

bus-type: GuestDiskBusType bus type

bus: int bus id

target: int target id

unit: int unit id

serial: string (optional) serial number (since: 3.1)

dev: string (optional) device node (POSIX) or device UNC (Windows) (since: 3.1)

ccw-address: GuestCCWAddress (optional) CCW address on s390x (since: 6.0)

Since

2.2

GuestDiskInfo (Object)

Members

name: string device node (Linux) or device UNC (Windows)

partition: boolean whether this is a partition or disk

dependencies: array of string (optional) list of device dependencies; e.g. for LVs of the LVM this will hold the list of PVs, for LUKS encrypted volume this will contain the disk where the volume is placed. (Linux)

address: GuestDiskAddress (optional) disk address information (only for non-virtual devices)

alias: string (optional) optional alias assigned to the disk, on Linux this is a name assigned by device mapper

Since 5.2

guest-get-disks (Command)

Returns

The list of disks in the guest. For Windows these are only the physical disks. On Linux these are all root block devices of non-zero size including e.g. removable devices, loop devices, NBD, etc.

Since

5.2

GuestFilesystemInfo (Object)

Members

name: **string** disk name

mountpoint: **string** mount point path

type: **string** file system type string

used-bytes: **int** (optional) file system used bytes (since 3.0)

total-bytes: **int** (optional) non-root file system total bytes (since 3.0)

disk: **array of GuestDiskAddress** an array of disk hardware information that the volume lies on, which may be empty if the disk type is not supported

Since

2.2

guest-get-fsinfo (Command)

Returns

The list of filesystems information mounted in the guest. The returned mountpoints may be specified to `guest-fsfreeze-freeze-list`. Network filesystems (such as CIFS and NFS) are not listed.

Since

2.2

guest-set-user-password (Command)

Arguments

username: **string** the user account whose password to change

password: **string** the new password entry string, base64 encoded

encrypted: **boolean** true if password is already crypt()d, false if raw

If the `encrypted` flag is true, it is the caller's responsibility to ensure the correct `crypt()` encryption scheme is used. This command does not attempt to interpret or report on the encryption scheme. Refer to the documentation of the guest operating system in question to determine what is supported.

Not all guest operating systems will support use of the `encrypted` flag, as they may require the clear-text password

The `password` parameter must always be base64 encoded before transmission, even if already `crypt()`d, to ensure it is 8-bit safe when passed as JSON.

Returns

Nothing on success.

Since

2.3

GuestMemoryBlock (Object)

Members

phys-index: int Arbitrary guest-specific unique identifier of the MEMORY BLOCK.

online: boolean Whether the MEMORY BLOCK is enabled in guest.

can-offline: boolean (optional) Whether offlining the MEMORY BLOCK is possible. This member is always filled in by the guest agent when the structure is returned, and always ignored on input (hence it can be omitted then).

Since

2.3

guest-get-memory-blocks (Command)

Retrieve the list of the guest's memory blocks.

This is a read-only operation.

Returns

The list of all memory blocks the guest knows about. Each memory block is put on the list exactly once, but their order is unspecified.

Since

2.3

GuestMemoryBlockResponseType (Enum)

An enumeration of memory block operation result.

Values

success the operation of online/offline memory block is successful.

not-found can't find the corresponding memoryXXX directory in sysfs.

operation-not-supported for some old kernels, it does not support online or offline memory block.

operation-failed the operation of online/offline memory block fails, because of some errors happen.

Since

2.3

GuestMemoryBlockResponse (Object)

Members

phys-index: int same with the 'phys-index' member of GuestMemoryBlock.

response: GuestMemoryBlockResponseType the result of memory block operation.

error-code: int (optional) the error number. When memory block operation fails, we assign the value of 'errno' to this member, it indicates what goes wrong. When the operation succeeds, it will be omitted.

Since

2.3

guest-set-memory-blocks (Command)

Attempt to reconfigure (currently: enable/disable) state of memory blocks inside the guest.

The input list is processed node by node in order. In each node `phys-index` is used to look up the guest MEMORY BLOCK, for which `online` specifies the requested state. The set of distinct `phys-index`'s is only required to be a subset of the guest-supported identifiers. There's no restriction on list length or on repeating the same `phys-index` (with possibly different `online` field). Preferably the input list should describe a modified subset of `guest-get-memory-blocks`' return value.

Arguments

mem-blks: array of GuestMemoryBlock Not documented

Returns

The operation results, it is a list of GuestMemoryBlockResponse, which is corresponding to the input list.

Note: it will return NULL if the `mem-blks` list was empty on input, or there is an error, and in this case, guest state will not be changed.

Since

2.3

GuestMemoryBlockInfo (Object)

Members

size: int the size (in bytes) of the guest memory blocks, which are the minimal units of memory block on-line/offline operations (also called Logical Memory Hotplug).

Since

2.3

guest-get-memory-block-info (Command)

Get information relating to guest memory blocks.

Returns

GuestMemoryBlockInfo

Since

2.3

GuestExecStatus (Object)

Members

exited: boolean true if process has already terminated.

exitcode: int (optional) process exit code if it was normally terminated.

signal: int (optional) signal number (linux) or unhandled exception code (windows) if the process was abnormally terminated.

out-data: string (optional) base64-encoded stdout of the process

err-data: string (optional) base64-encoded stderr of the process Note: `out-data` and `err-data` are present only if 'capture-output' was specified for 'guest-exec'

out-truncated: boolean (optional) true if stdout was not fully captured due to size limitation.

err-truncated: boolean (optional) true if stderr was not fully captured due to size limitation.

Since

2.5

guest-exec-status (Command)

Check status of process associated with PID retrieved via `guest-exec`. Reap the process and associated metadata if it has exited.

Arguments

pid: int pid returned from guest-exec

Returns

GuestExecStatus on success.

Since

2.5

GuestExec (Object)

Members

pid: int pid of child process in guest OS

Since

2.5

guest-exec (Command)

Execute a command in the guest

Arguments

path: string path or executable name to execute

arg: array of string (optional) argument list to pass to executable

env: array of string (optional) environment variables to pass to executable

input-data: string (optional) data to be passed to process stdin (base64 encoded)

capture-output: boolean (optional) bool flag to enable capture of stdout/stderr of running process. defaults to false.

Returns

PID on success.

Since

2.5

GuestHostName (Object)

Members

host-name: string Fully qualified domain name of the guest OS

Since

2.10

guest-get-host-name (Command)

Return a name for the machine.

The returned name is not necessarily a fully-qualified domain name, or even present in DNS or some other name service at all. It need not even be unique on your local network or site, but usually it is.

Returns

the host name of the machine on success

Since

2.10

GuestUser (Object)

Members

user: string Username

domain: string (optional) Logon domain (windows only)

login-time: number Time of login of this user on the computer. If multiple instances of the user are logged in, the earliest login time is reported. The value is in fractional seconds since epoch time.

Since

2.10

guest-get-users (Command)

Retrieves a list of currently active users on the VM.

Returns

A unique list of users.

Since

2.10

GuestTimezone (Object)**Members**

zone: string (optional) Timezone name. These values may differ depending on guest/OS and should only be used for informational purposes.

offset: int Offset to UTC in seconds, negative numbers for time zones west of GMT, positive numbers for east

Since

2.10

guest-get-timezone (Command)

Retrieves the timezone information from the guest.

Returns

A GuestTimezone dictionary.

Since

2.10

GuestOSInfo (Object)**Members****kernel-release: string (optional)**

- POSIX: release field returned by uname(2)
- Windows: build number of the OS

kernel-version: string (optional)

- POSIX: version field returned by uname(2)
- Windows: version number of the OS

machine: string (optional)

- POSIX: machine field returned by uname(2)
- Windows: one of x86, x86_64, arm, ia64

id: string (optional)

- POSIX: as defined by `os-release(5)`
- Windows: contains string “mswindows”

name: string (optional)

- POSIX: as defined by `os-release(5)`
- Windows: contains string “Microsoft Windows”

pretty-name: string (optional)

- POSIX: as defined by `os-release(5)`
- Windows: product name, e.g. “Microsoft Windows 10 Enterprise”

version: string (optional)

- POSIX: as defined by `os-release(5)`
- Windows: long version string, e.g. “Microsoft Windows Server 2008”

version-id: string (optional)

- POSIX: as defined by `os-release(5)`
- Windows: short version identifier, e.g. “7” or “20012r2”

variant: string (optional)

- POSIX: as defined by `os-release(5)`
- Windows: contains string “server” or “client”

variant-id: string (optional)

- POSIX: as defined by `os-release(5)`
- Windows: contains string “server” or “client”

Notes

On POSIX systems the fields `id`, `name`, `pretty-name`, `version`, `version-id`, `variant` and `variant-id` follow the definition specified in `os-release(5)`. Refer to the manual page for exact description of the fields. Their values are taken from the `os-release` file. If the file is not present in the system, or the values are not present in the file, the fields are not included.

On Windows the values are filled from information gathered from the system.

Since

2.10

`guest-get-osinfo` (Command)

Retrieve guest operating system information

Returns

`GuestOSInfo`

Since

2.10

GuestDeviceType (Enum)**Values**

pci Not documented

GuestDeviceIdPCI (Object)**Members**

vendor-id: int vendor ID

device-id: int device ID

Since

5.2

GuestDeviceId (Object)

Id of the device - **pci**: PCI ID, since: 5.2

Members

type: GuestDeviceType Not documented

The members of **GuestDeviceIdPCI** when **type** is "pci"

Since

5.2

GuestDeviceInfo (Object)**Members**

driver-name: string name of the associated driver

driver-date: int (optional) driver release date, in nanoseconds since the epoch

driver-version: string (optional) driver version

id: GuestDeviceId (optional) device ID

Since

5.2

guest-get-devices (Command)

Retrieve information about device drivers in Windows guest

Returns

GuestDeviceInfo

Since

5.2

GuestAuthorizedKeys (Object)

Members

keys: **array of string** public keys (in OpenSSH/sshd(8) authorized_keys format)

Since

5.2

If

defined(CONFIG_POSIX)

guest-ssh-get-authorized-keys (Command)

Arguments

username: **string** the user account to add the authorized keys

Return the public keys from user .ssh/authorized_keys on Unix systems (not implemented for other systems).

Returns

GuestAuthorizedKeys

Since

5.2

If

`defined(CONFIG_POSIX)`

guest-ssh-add-authorized-keys (Command)

Arguments

username: string the user account to add the authorized keys

keys: array of string the public keys to add (in OpenSSH/sshd(8) authorized_keys format)

reset: boolean (optional) ignore the existing content, set it with the given keys only

Append public keys to user `.ssh/authorized_keys` on Unix systems (not implemented for other systems).

Returns

Nothing on success.

Since

5.2

If

`defined(CONFIG_POSIX)`

guest-ssh-remove-authorized-keys (Command)

Arguments

username: string the user account to remove the authorized keys

keys: array of string the public keys to remove (in OpenSSH/sshd(8) authorized_keys format)

Remove public keys from the user `.ssh/authorized_keys` on Unix systems (not implemented for other systems). It's not an error if the key is already missing.

Returns

Nothing on success.

Since

5.2

If

```
defined(CONFIG_POSIX)
```

4.8 QEMU QMP Reference Manual

4.8.1 Introduction

This document describes all commands currently supported by QMP.

Most of the time their usage is exactly the same as in the user Monitor, this means that any other document which also describe commands (the manpage, QEMU's manual, etc) can and should be consulted.

QMP has two types of commands: regular and query commands. Regular commands usually change the Virtual Machine's state somehow, while query commands just return information. The sections below are divided accordingly.

It's important to observe that all communication examples are formatted in a reader-friendly way, so that they're easier to understand. However, in real protocol usage, they're emitted as a single line.

Also, the following notation is used to denote data flow:

Example:

```
-> data issued by the Client
<- Server data response
```

Please, refer to the QMP specification (docs/interop/qmp-spec.txt) for detailed information on the Server command and response formats.

4.8.2 Stability Considerations

The current QMP command set (described in this file) may be useful for a number of use cases, however it's limited and several commands have bad defined semantics, specially with regard to command completion.

These problems are going to be solved incrementally in the next QEMU releases and we're going to establish a deprecation policy for badly defined commands.

If you're planning to adopt QMP, please observe the following:

1. The deprecation policy will take effect and be documented soon, please check the documentation of each used command as soon as a new release of QEMU is available
2. DO NOT rely on anything which is not explicit documented
3. Errors, in special, are not documented. Applications should NOT check for specific errors classes or data (it's strongly recommended to only check for the "error" key)

4.8.3 QMP errors

QapiErrorClass (Enum)

QEMU error classes

Values

GenericError this is used for errors that don't require a specific error class. This should be the default case for most errors

CommandNotFound the requested command has not been found

DeviceNotActive a device has failed to become active

DeviceNotFound the requested device has not been found

KVMMissingCap the requested operation can't be fulfilled because a required KVM capability is missing

Since

1.2

4.8.4 Common data types

IoOperationType (Enum)

An enumeration of the I/O operation types

Values

read read operation

write write operation

Since

2.1

OnOffAuto (Enum)

An enumeration of three options: on, off, and auto

Values

auto QEMU selects the value between on and off

on Enabled

off Disabled

Since

2.2

OnOffSplit (Enum)

An enumeration of three values: on, off, and split

Values

on Enabled

off Disabled

split Mixed

Since

2.6

String (Object)

A fat type wrapping 'str', to be embedded in lists.

Members

str: string Not documented

Since

1.2

StrOrNull (Alternate)

This is a string value or the explicit lack of a string (null pointer in C). Intended for cases when 'optional absent' already has a different meaning.

Members

s: string the string value

n: null no string value

Since

2.10

OffAutoPCIBAR (Enum)

An enumeration of options for specifying a PCI BAR

Values

- off** The specified feature is disabled
- auto** The PCI BAR for the feature is automatically selected
- bar0** PCI BAR0 is used for the feature
- bar1** PCI BAR1 is used for the feature
- bar2** PCI BAR2 is used for the feature
- bar3** PCI BAR3 is used for the feature
- bar4** PCI BAR4 is used for the feature
- bar5** PCI BAR5 is used for the feature

Since

2.12

PCIELinkSpeed (Enum)

An enumeration of PCIe link speeds in units of GT/s

Values

- 2_5** 2.5GT/s
- 5** 5.0GT/s
- 8** 8.0GT/s
- 16** 16.0GT/s

Since

4.0

PCIELinkWidth (Enum)

An enumeration of PCIe link width

Values

- 1** x1
- 2** x2
- 4** x4
- 8** x8
- 12** x12

16 x16

32 x32

Since

4.0

4.8.5 Socket data types

NetworkAddressFamily (Enum)

The network address family

Values

ipv4 IPV4 family

ipv6 IPV6 family

unix unix socket

vsock vsock family (since 2.8)

unknown otherwise

Since

2.1

InetSocketAddressBase (Object)

Members

host: string host part of the address

port: string port part of the address

InetSocketAddress (Object)

Captures a socket address or address range in the Internet namespace.

Members

numeric: boolean (optional) true if the host/port are guaranteed to be numeric, false if name resolution should be attempted. Defaults to false. (Since 2.9)

to: int (optional) If present, this is range of possible addresses, with port between `port` and `to`.

ipv4: boolean (optional) whether to accept IPv4 addresses, default try both IPv4 and IPv6

ipv6: boolean (optional) whether to accept IPv6 addresses, default try both IPv4 and IPv6

keep-alive: boolean (optional) enable keep-alive when connecting to this socket. Not supported for passive sockets. (Since 4.2)

The members of `InetSocketAddressBase`

Since

1.3

`UnixSocketAddress` (Object)

Captures a socket address in the local (“Unix socket”) namespace.

Members

path: string filesystem path to use

abstract: boolean (optional) (If: `defined(CONFIG_LINUX)`) if true, this is a Linux abstract socket address. `path` will be prefixed by a null byte, and optionally padded with null bytes. Defaults to false. (Since 5.1)

tight: boolean (optional) (If: `defined(CONFIG_LINUX)`) if false, pad an abstract socket address with enough null bytes to make it fill struct `sockaddr_un` member `sun_path`. Defaults to true. (Since 5.1)

Since

1.3

`VsockSocketAddress` (Object)

Captures a socket address in the vsock namespace.

Members

cid: string unique host identifier

port: string port

Note

string types are used to allow for possible future hostname or service resolution support.

Since

2.8

SocketAddressLegacy (Object)

Captures the address of a socket, which could also be a named file descriptor

Members

type One of `inet`, `unix`, `vsock`, `fd`

data: `InetSocketAddress` when **type** is `"inet"`

data: `UnixSocketAddress` when **type** is `"unix"`

data: `VsockSocketAddress` when **type** is `"vsock"`

data: `String` when **type** is `"fd"`

Note

This type is deprecated in favor of `SocketAddress`. The difference between `SocketAddressLegacy` and `SocketAddress` is that the latter is a flat union rather than a simple union. Flat is nicer because it avoids nesting on the wire, i.e. that form has fewer `{}`.

Since

1.3

SocketAddressType (Enum)

Available `SocketAddress` types

Values

inet Internet address

unix Unix domain socket

vsock VMCI address

fd decimal is for file descriptor number, otherwise a file descriptor name. Named file descriptors are permitted in monitor commands, in combination with the `'getfd'` command. Decimal file descriptors are permitted at startup or other contexts where no monitor context is active.

Since

2.9

SocketAddress (Object)

Captures the address of a socket, which could also be a named file descriptor

Members

type: SocketAddressType Transport type

The members of **InetSocketAddress** when **type** is "inet"

The members of **UnixSocketAddress** when **type** is "unix"

The members of **VsockSocketAddress** when **type** is "vsock"

The members of **String** when **type** is "fd"

Since

2.9

4.8.6 VM run state

RunState (Enum)

An enumeration of VM run states.

Values

debug QEMU is running on a debugger

finish-migrate guest is paused to finish the migration process

inmigrate guest is paused waiting for an incoming migration. Note that this state does not tell whether the machine will start at the end of the migration. This depends on the command-line -S option and any invocation of 'stop' or 'cont' that has happened since QEMU was started.

internal-error An internal error that prevents further guest execution has occurred

io-error the last IOP has failed and the device is configured to pause on I/O errors

paused guest has been paused via the 'stop' command

postmigrate guest is paused following a successful 'migrate'

prelaunch QEMU was started with -S and guest has not started

restore-vm guest is paused to restore VM state

running guest is actively running

save-vm guest is paused to save the VM state

shutdown guest is shut down (and -no-shutdown is in use)

suspended guest is suspended (ACPI S3)

watchdog the watchdog action is configured to pause and has been triggered

guest-panicked guest has been panicked as a result of guest OS panic

colo guest is paused to save/restore VM state under colo checkpoint, VM can not get into this state unless colo capability is enabled for migration. (since 2.8)

ShutdownCause (Enum)

An enumeration of reasons for a Shutdown.

Values

none No shutdown request pending

host-error An error prevents further use of guest

host-qmp-quit Reaction to the QMP command ‘quit’

host-qmp-system-reset Reaction to the QMP command ‘system_reset’

host-signal Reaction to a signal, such as SIGINT

host-ui Reaction to a UI event, like window close

guest-shutdown Guest shutdown/suspend request, via ACPI or other hardware-specific means

guest-reset Guest reset request, and command line turns that into a shutdown

guest-panic Guest panicked, and command line turns that into a shutdown

subsystem-reset Partial guest reset that does not trigger QMP events and ignores `–no-reboot`. This is useful for sanitizing hypercalls on s390 that are used during `kexec/kdump/boot`

StatusInfo (Object)

Information about VCPU run state

Members

running: boolean true if all VCPUs are runnable, false if not runnable

singlestep: boolean true if VCPUs are in single-step mode

status: RunState the virtual machine `RunState`

Since

0.14

Notes

`singlestep` is enabled through the GDB stub

query-status (Command)

Query the run status of all VCPUs

Returns

StatusInfo reflecting all VCPUs

Since

0.14

Example

```
-> { "execute": "query-status" }
<- { "return": { "running": true,
                  "singlestep": false,
                  "status": "running" } }
```

SHUTDOWN (Event)

Emitted when the virtual machine has shut down, indicating that qemu is about to exit.

Arguments

guest: boolean If true, the shutdown was triggered by a guest request (such as a guest-initiated ACPI shutdown request or other hardware-specific action) rather than a host request (such as sending qemu a SIGINT). (since 2.10)

reason: ShutdownCause The ShutdownCause which resulted in the SHUTDOWN. (since 4.0)

Note

If the command-line option “-no-shutdown” has been specified, qemu will not exit, and a STOP event will eventually follow the SHUTDOWN event

Since

0.12

Example

```
<- { "event": "SHUTDOWN", "data": { "guest": true },
      "timestamp": { "seconds": 1267040730, "microseconds": 682951 } }
```

POWERDOWN (Event)

Emitted when the virtual machine is powered down through the power control system, such as via ACPI.

Since

0.12

Example

```
<- { "event": "POWERDOWN",  
      "timestamp": { "seconds": 1267040730, "microseconds": 682951 } }
```

RESET (Event)

Emitted when the virtual machine is reset

Arguments

guest: boolean If true, the reset was triggered by a guest request (such as a guest-initiated ACPI reboot request or other hardware-specific action) rather than a host request (such as the QMP command `system_reset`). (since 2.10)

reason: ShutdownCause The ShutdownCause of the RESET. (since 4.0)

Since

0.12

Example

```
<- { "event": "RESET", "data": { "guest": false },  
      "timestamp": { "seconds": 1267041653, "microseconds": 9518 } }
```

STOP (Event)

Emitted when the virtual machine is stopped

Since

0.12

Example

```
<- { "event": "STOP",  
      "timestamp": { "seconds": 1267041730, "microseconds": 281295 } }
```


RESUME (Event)

Emitted when the virtual machine resumes execution

Since

0.12

Example

```
<- { "event": "RESUME",  
      "timestamp": { "seconds": 1271770767, "microseconds": 582542 } }
```

SUSPEND (Event)

Emitted when guest enters a hardware suspension state, for example, S3 state, which is sometimes called standby state

Since

1.1

Example

```
<- { "event": "SUSPEND",  
      "timestamp": { "seconds": 1344456160, "microseconds": 309119 } }
```

SUSPEND_DISK (Event)

Emitted when guest enters a hardware suspension state with data saved on disk, for example, S4 state, which is sometimes called hibernate state

Note

QEMU shuts down (similar to event SHUTDOWN) when entering this state

Since

1.2

Example

```
<- { "event": "SUSPEND_DISK",  
      "timestamp": { "seconds": 1344456160, "microseconds": 309119 } }
```

WAKEUP (Event)

Emitted when the guest has woken up from suspend state and is running

Since

1.1

Example

```
<- { "event": "WAKEUP",  
      "timestamp": { "seconds": 1344522075, "microseconds": 745528 } }
```

WATCHDOG (Event)

Emitted when the watchdog device's timer is expired

Arguments

action: WatchdogAction action that has been taken

Note

If action is “reset”, “shutdown”, or “pause” the WATCHDOG event is followed respectively by the RESET, SHUT-DOWN, or STOP events

Note

This event is rate-limited.

Since

0.13

Example

```
<- { "event": "WATCHDOG",  
      "data": { "action": "reset" },  
      "timestamp": { "seconds": 1267061043, "microseconds": 959568 } }
```

WatchdogAction (Enum)

An enumeration of the actions taken when the watchdog device's timer is expired

Values

reset system resets

shutdown system shutdown, note that it is similar to `powerdown`, which tries to set to system status and notify guest

poweroff system poweroff, the emulator program exits

pause system pauses, similar to `stop`

debug system enters debug state

none nothing is done

inject-nmi a non-maskable interrupt is injected into the first VCPU (all VCPUS on x86) (since 2.4)

Since

2.1

RebootAction (Enum)

Possible QEMU actions upon guest reboot

Values

reset Reset the VM

shutdown Shutdown the VM and exit, according to the shutdown action

Since

6.0

ShutdownAction (Enum)

Possible QEMU actions upon guest shutdown

Values

poweroff Shutdown the VM and exit

pause pause the VM#

Since

6.0

PanicAction (Enum)

Values

none Continue VM execution

pause Pause the VM

shutdown Shutdown the VM and exit, according to the shutdown action

Since

6.0

watchdog-set-action (Command)

Set watchdog action

Arguments

action: WatchdogAction Not documented

Since

2.11

set-action (Command)

Set the actions that will be taken by the emulator in response to guest events.

Arguments

reboot: RebootAction (optional) RebootAction action taken on guest reboot.

shutdown: ShutdownAction (optional) ShutdownAction action taken on guest shutdown.

panic: PanicAction (optional) PanicAction action taken on guest panic.

watchdog: WatchdogAction (optional) WatchdogAction action taken when watchdog timer expires .

Returns

Nothing on success.

Since

6.0

Example

```
-> { "execute": "set-action",
      "arguments": { "reboot": "shutdown",
                     "shutdown": "pause",
                     "panic": "pause",
                     "watchdog": "inject-nmi" } }
<- { "return": {} }
```

GUEST_PANICKED (Event)

Emitted when guest OS panic is detected

Arguments

action: `GuestPanicAction` action that has been taken, currently always “pause”

info: `GuestPanicInformation` (optional) information about a panic (since 2.9)

Since

1.5

Example

```
<- { "event": "GUEST_PANICKED",
      "data": { "action": "pause" } }
```

GUEST_CRASHLOADED (Event)

Emitted when guest OS crash loaded is detected

Arguments

action: `GuestPanicAction` action that has been taken, currently always “run”

info: `GuestPanicInformation` (optional) information about a panic

Since

5.0

Example

```
<- { "event": "GUEST_CRASHLOADED",
      "data": { "action": "run" } }
```

GuestPanicAction (Enum)

An enumeration of the actions taken when guest OS panic is detected

Values

pause system pauses

poweroff Not documented

run Not documented

Since

2.1 (poweroff since 2.8, run since 5.0)

GuestPanicInformationType (Enum)

An enumeration of the guest panic information types

Values

hyper-v hyper-v guest panic information type

s390 s390 guest panic information type (Since: 2.12)

Since

2.9

GuestPanicInformation (Object)

Information about a guest panic

Members

type: GuestPanicInformationType Crash type that defines the hypervisor specific information

The members of **GuestPanicInformationHyperV** when **type** is "hyper-v"

The members of **GuestPanicInformationS390** when **type** is "s390"

Since

2.9

GuestPanicInformationHyperV (Object)

Hyper-V specific guest panic information (HV crash MSRs)

Members

arg1: int Not documented

arg2: int Not documented

arg3: int Not documented

arg4: int Not documented

arg5: int Not documented

Since

2.9

S390CrashReason (Enum)

Reason why the CPU is in a crashed state.

Values

unknown no crash reason was set

disabled-wait the CPU has entered a disabled wait state

extint-loop clock comparator or cpu timer interrupt with new PSW enabled for external interrupts

pgmint-loop program interrupt with BAD new PSW

opint-loop operation exception interrupt with invalid code at the program interrupt new PSW

Since

2.12

GuestPanicInformationS390 (Object)

S390 specific guest panic information (PSW)

Members

core: int core id of the CPU that crashed

psw-mask: int control fields of guest PSW

psw-addr: int guest instruction address

reason: S390CrashReason guest crash reason

Since

2.12

MEMORY_FAILURE (Event)

Emitted when a memory failure occurs on host side.

Arguments

recipient: MemoryFailureRecipient recipient is defined as MemoryFailureRecipient.

action: MemoryFailureAction action that has been taken. action is defined as MemoryFailureAction.

flags: MemoryFailureFlags flags for MemoryFailureAction. action is defined as MemoryFailureFlags.

Since

5.2

Example

```
<- { "event": "MEMORY_FAILURE",  
      "data": { "recipient": "hypervisor",  
                 "action": "fatal",  
                 "flags": { 'action-required': false } } }
```

MemoryFailureRecipient (Enum)

Hardware memory failure occurs, handled by recipient.

Values

hypervisor memory failure at QEMU process address space. (none guest memory, but used by QEMU itself).

guest memory failure at guest memory,

Since

5.2

MemoryFailureAction (Enum)

Actions taken by QEMU in response to a hardware memory failure.

Values

ignore the memory failure could be ignored. This will only be the case for action-optional failures.

inject memory failure occurred in guest memory, the guest enabled MCE handling mechanism, and QEMU could inject the MCE into the guest successfully.

fatal the failure is unrecoverable. This occurs for action-required failures if the recipient is the hypervisor; QEMU will exit.

reset the failure is unrecoverable but confined to the guest. This occurs if the recipient is a guest which is not ready to handle memory failures.

Since

5.2

MemoryFailureFlags (Object)

Additional information on memory failures.

Members

action-required: boolean whether a memory failure event is action-required or action-optional (e.g. a failure during memory scrub).

recursive: boolean whether the failure occurred while the previous failure was still in progress.

Since

5.2

4.8.7 Cryptography

QCryptoTLSCredsEndpoint (Enum)

The type of network endpoint that will be using the credentials. Most types of credential require different setup / structures depending on whether they will be used in a server versus a client.

Values

client the network endpoint is acting as the client

server the network endpoint is acting as the server

Since

2.5

QCryptoSecretFormat (Enum)

The data format that the secret is provided in

Values

raw raw bytes. When encoded in JSON only valid UTF-8 sequences can be used

base64 arbitrary base64 encoded binary data

Since

2.6

QCryptoHashAlgorithm (Enum)

The supported algorithms for computing content digests

Values

md5 MD5. Should not be used in any new code, legacy compat only

sha1 SHA-1. Should not be used in any new code, legacy compat only

sha224 SHA-224. (since 2.7)

sha256 SHA-256. Current recommended strong hash.

sha384 SHA-384. (since 2.7)

sha512 SHA-512. (since 2.7)

ripemd160 RIPEMD-160. (since 2.7)

Since

2.6

QCryptoCipherAlgorithm (Enum)

The supported algorithms for content encryption ciphers

Values

aes-128 AES with 128 bit / 16 byte keys

aes-192 AES with 192 bit / 24 byte keys

aes-256 AES with 256 bit / 32 byte keys

des-rgb RFB specific variant of single DES. Do not use except in VNC.

3des 3DES(EDE) with 192 bit / 24 byte keys (since 2.9)

cast5-128 Cast5 with 128 bit / 16 byte keys

serpent-128 Serpent with 128 bit / 16 byte keys

serpent-192 Serpent with 192 bit / 24 byte keys

serpent-256 Serpent with 256 bit / 32 byte keys

twofish-128 Twofish with 128 bit / 16 byte keys

twofish-192 Twofish with 192 bit / 24 byte keys

twofish-256 Twofish with 256 bit / 32 byte keys

Since

2.6

QCryptoCipherMode (Enum)

The supported modes for content encryption ciphers

Values

ecb Electronic Code Book

cbc Cipher Block Chaining

xts XEX with tweaked code book and ciphertext stealing

ctr Counter (Since 2.8)

Since

2.6

QCryptoIVGenAlgorithm (Enum)

The supported algorithms for generating initialization vectors for full disk encryption. The 'plain' generator should not be used for disks with sector numbers larger than 2^{32} , except where compatibility with pre-existing Linux dm-crypt volumes is required.

Values

plain 64-bit sector number truncated to 32-bits

plain64 64-bit sector number

essiv 64-bit sector number encrypted with a hash of the encryption key

Since

2.6

QCryptoBlockFormat (Enum)

The supported full disk encryption formats

Values

qcow QCow/QCOW2 built-in AES-CBC encryption. Use only for liberating data from old images.

luks LUKS encryption format. Recommended for new images

Since

2.6

QCryptoBlockOptionsBase (Object)

The common options that apply to all full disk encryption formats

Members

format: QCryptoBlockFormat the encryption format

Since

2.6

QCryptoBlockOptionsQCOW (Object)

The options that apply to QCow/QCOW2 AES-CBC encryption format

Members

key-secret: string (optional) the ID of a QCryptoSecret object providing the decryption key. Mandatory except when probing image for metadata only.

Since

2.6

QCryptoBlockOptionsLUKS (Object)

The options that apply to LUKS encryption format

Members

key-secret: string (optional) the ID of a QCryptoSecret object providing the decryption key. Mandatory except when probing image for metadata only.

Since

2.6

`QCryptoBlockCreateOptionsLUKS` (Object)

The options that apply to LUKS encryption format initialization

Members

cipher-alg: `QCryptoCipherAlgorithm` (optional) the cipher algorithm for data encryption Currently defaults to 'aes-256'.

cipher-mode: `QCryptoCipherMode` (optional) the cipher mode for data encryption Currently defaults to 'xts'

ivgen-alg: `QCryptoIVGenAlgorithm` (optional) the initialization vector generator Currently defaults to 'plain64'

ivgen-hash-alg: `QCryptoHashAlgorithm` (optional) the initialization vector generator hash Currently defaults to 'sha256'

hash-alg: `QCryptoHashAlgorithm` (optional) the master key hash algorithm Currently defaults to 'sha256'

iter-time: `int` (optional) number of milliseconds to spend in PBKDF passphrase processing. Currently defaults to 2000. (since 2.8)

The members of `QCryptoBlockOptionsLUKS`

Since

2.6

`QCryptoBlockOpenOptions` (Object)

The options that are available for all encryption formats when opening an existing volume

Members

The members of `QCryptoBlockOptionsBase`

The members of `QCryptoBlockOptionsQcow` when format is "qcow"

The members of `QCryptoBlockOptionsLUKS` when format is "luks"

Since

2.6

`QCryptoBlockCreateOptions` (Object)

The options that are available for all encryption formats when initializing a new volume

Members

The members of `QCryptoBlockOptionsBase`

The members of `QCryptoBlockOptionsQCow` when format is "qcow"

The members of `QCryptoBlockCreateOptionsLUKS` when format is "luks"

Since

2.6

`QCryptoBlockInfoBase` (Object)

The common information that applies to all full disk encryption formats

Members

format: `QCryptoBlockFormat` the encryption format

Since

2.7

`QCryptoBlockInfoLUKSSlot` (Object)

Information about the LUKS block encryption key slot options

Members

active: `boolean` whether the key slot is currently in use

key-offset: `int` offset to the key material in bytes

iters: `int` (optional) number of PBKDF2 iterations for key material

stripes: `int` (optional) number of stripes for splitting key material

Since

2.7

`QCryptoBlockInfoLUKS` (Object)

Information about the LUKS block encryption options

Members

cipher-alg: `QCryptoCipherAlgorithm` the cipher algorithm for data encryption
cipher-mode: `QCryptoCipherMode` the cipher mode for data encryption
ivgen-alg: `QCryptoIVGenAlgorithm` the initialization vector generator
ivgen-hash-alg: `QCryptoHashAlgorithm` (optional) the initialization vector generator hash
hash-alg: `QCryptoHashAlgorithm` the master key hash algorithm
payload-offset: `int` offset to the payload data in bytes
master-key-iters: `int` number of PBKDF2 iterations for key material
uuid: `string` unique identifier for the volume
slots: `array of QCryptoBlockInfoLUKSSlot` information about each key slot

Since

2.7

`QCryptoBlockInfo` (Object)

Information about the block encryption options

Members

The members of `QCryptoBlockInfoBase`

The members of `QCryptoBlockInfoLUKS` when format is "luks"

Since

2.7

`QCryptoBlockLUKSKeyslotState` (Enum)

Defines state of keyslots that are affected by the update

Values

active The slots contain the given password and marked as active
inactive The slots are erased (contain garbage) and marked as inactive

Since

5.1

QCryptoBlockAmendOptionsLUKS (Object)

This struct defines the update parameters that activate/de-activate set of keyslots

Members

state: `QCryptoBlockLUKSKeyslotState` the desired state of the keyslots

new-secret: `string (optional)` The ID of a `QCryptoSecret` object providing the password to be written into added active keyslots

old-secret: `string (optional)` Optional (for deactivation only) If given will deactivate all keyslots that match password located in `QCryptoSecret` with this ID

iter-time: `int (optional)` Optional (for activation only) Number of milliseconds to spend in PBKDF passphrase processing for the newly activated keyslot. Currently defaults to 2000.

keyslot: `int (optional)` Optional. ID of the keyslot to activate/deactivate. For keyslot activation, keyslot should not be active already (this is unsafe to update an active keyslot), but possible if 'force' parameter is given. If keyslot is not given, first free keyslot will be written.

For keyslot deactivation, this parameter specifies the exact keyslot to deactivate

secret: `string (optional)` Optional. The ID of a `QCryptoSecret` object providing the password to use to retrieve current master key. Defaults to the same secret that was used to open the image

Since 5.1

QCryptoBlockAmendOptions (Object)

The options that are available for all encryption formats when amending encryption settings

Members

The members of `QCryptoBlockOptionsBase`

The members of `QCryptoBlockAmendOptionsLUKS` when `format` is "luks"

Since

5.1

4.8.8 Block devices

Block core (VM unrelated)

Background jobs

JobType (Enum)

Type of a background job.

Values

commit block commit job type, see “block-commit”
stream block stream job type, see “block-stream”
mirror drive mirror job type, see “drive-mirror”
backup drive backup job type, see “drive-backup”
create image creation job type, see “blockdev-create” (since 3.0)
amend image options amend job type, see “x-blockdev-amend” (since 5.1)
snapshot-load snapshot load job type, see “snapshot-load” (since 6.0)
snapshot-save snapshot save job type, see “snapshot-save” (since 6.0)
snapshot-delete snapshot delete job type, see “snapshot-delete” (since 6.0)

Since

1.7

JobStatus (Enum)

Indicates the present state of a given job in its lifetime.

Values

undefined Erroneous, default state. Should not ever be visible.
created The job has been created, but not yet started.
running The job is currently running.
paused The job is running, but paused. The pause may be requested by either the QMP user or by internal processes.
ready The job is running, but is ready for the user to signal completion. This is used for long-running jobs like mirror that are designed to run indefinitely.
standby The job is ready, but paused. This is nearly identical to `paused`. The job may return to `ready` or otherwise be canceled.
waiting The job is waiting for other jobs in the transaction to converge to the waiting state. This status will likely not be visible for the last job in a transaction.
pending The job has finished its work, but has finalization steps that it needs to make prior to completing. These changes will require manual intervention via `job-finalize` if `auto-finalize` was set to false. These pending changes may still fail.
aborting The job is in the process of being aborted, and will finish with an error. The job will afterwards report that it is concluded. This status may not be visible to the management process.
concluded The job has finished all work. If `auto-dismiss` was set to false, the job will remain in the query list until it is dismissed via `job-dismiss`.
null The job is in the process of being dismantled. This state should not ever be visible externally.

Since

2.12

JobVerb (Enum)

Represents command verbs that can be applied to a job.

Values

cancel see job-cancel

pause see job-pause

resume see job-resume

set-speed see block-job-set-speed

complete see job-complete

dismiss see job-dismiss

finalize see job-finalize

Since

2.12

JOB_STATUS_CHANGE (Event)

Emitted when a job transitions to a different status.

Arguments

id: string The job identifier

status: JobStatus The new job status

Since

3.0

job-pause (Command)

Pause an active job.

This command returns immediately after marking the active job for pausing. Pausing an already paused job is an error.

The job will pause as soon as possible, which means transitioning into the PAUSED state if it was RUNNING, or into STANDBY if it was READY. The corresponding JOB_STATUS_CHANGE event will be emitted.

Cancelling a paused job automatically resumes it.

Arguments

id: string The job identifier.

Since

3.0

job-resume (Command)

Resume a paused job.

This command returns immediately after resuming a paused job. Resuming an already running job is an error.

id : The job identifier.

Arguments

id: string Not documented

Since

3.0

job-cancel (Command)

Instruct an active background job to cancel at the next opportunity. This command returns immediately after marking the active job for cancellation.

The job will cancel as soon as possible and then emit a `JOB_STATUS_CHANGE` event. Usually, the status will change to `ABORTING`, but it is possible that a job successfully completes (e.g. because it was almost done and there was no opportunity to cancel earlier than completing the job) and transitions to `PENDING` instead.

Arguments

id: string The job identifier.

Since

3.0

job-complete (Command)

Manually trigger completion of an active job in the `READY` state.

Arguments

id: string The job identifier.

Since

3.0

job-dismiss (Command)

Deletes a job that is in the CONCLUDED state. This command only needs to be run explicitly for jobs that don't have automatic dismiss enabled.

This command will refuse to operate on any job that has not yet reached its terminal state, JOB_STATUS_CONCLUDED. For jobs that make use of JOB_READY event, job-cancel or job-complete will still need to be used as appropriate.

Arguments

id: string The job identifier.

Since

3.0

job-finalize (Command)

Instructs all jobs in a transaction (or a single job if it is not part of any transaction) to finalize any graph changes and do any necessary cleanup. This command requires that all involved jobs are in the PENDING state.

For jobs in a transaction, instructing one job to finalize will force ALL jobs in the transaction to finalize, so it is only necessary to instruct a single member job to finalize.

Arguments

id: string The identifier of any job in the transaction, or of a job that is not part of any transaction.

Since

3.0

JobInfo (Object)

Information about a job.

Members

id: string The job identifier

type: JobType The kind of job that is being performed

status: JobStatus Current job state/status

current-progress: int Progress made until now. The unit is arbitrary and the value can only meaningfully be used for the ratio of `current-progress` to `total-progress`. The value is monotonically increasing.

total-progress: int Estimated `current-progress` value at the completion of the job. This value can arbitrarily change while the job is running, in both directions.

error: string (optional) If this field is present, the job failed; if it is still missing in the `CONCLUDED` state, this indicates successful completion.

The value is a human-readable error message to describe the reason for the job failure. It should not be parsed by applications.

Since

3.0

query-jobs (Command)

Return information about jobs.

Returns

a list with a `JobInfo` for each active job

Since

3.0

SnapshotInfo (Object)

Members

id: string unique snapshot id

name: string user chosen name

vm-state-size: int size of the VM state

date-sec: int UTC date of the snapshot in seconds

date-nsec: int fractional part in nano seconds to be used with `date-sec`

vm-clock-sec: int VM clock relative to boot in seconds

vm-clock-nsec: int fractional part in nano seconds to be used with `vm-clock-sec`

icount: int (optional) Current instruction count. Appears when execution record/replay is enabled. Used for “time-traveling” to match the moment in the recorded execution with the snapshots. This counter may be obtained through `query-replay` command (since 5.2)

Since

1.3

ImageInfoSpecificQCow2EncryptionBase (Object)

Members

format: BlockdevQcow2EncryptionFormat The encryption format

Since

2.10

ImageInfoSpecificQCow2Encryption (Object)

Members

The members of `ImageInfoSpecificQCow2EncryptionBase`

The members of `QCryptoBlockInfoLUKS` when `format` is "luks"

Since

2.10

ImageInfoSpecificQCow2 (Object)

Members

compat: string compatibility level

data-file: string (optional) the filename of the external data file that is stored in the image and used as a default for opening the image (since: 4.0)

data-file-raw: boolean (optional) True if the external data file must stay valid as a standalone (read-only) raw image without looking at qcow2 metadata (since: 4.0)

extended-l2: boolean (optional) true if the image has extended L2 entries; only valid for `compat` \geq 1.1 (since 5.2)

lazy-refcounts: boolean (optional) on or off; only valid for `compat` \geq 1.1

corrupt: boolean (optional) true if the image has been marked corrupt; only valid for `compat` \geq 1.1 (since 2.2)

refcount-bits: int width of a refcount entry in bits (since 2.3)

encrypt: ImageInfoSpecificQCow2Encryption (optional) details about encryption parameters; only set if image is encrypted (since 2.10)

bitmaps: array of Qcow2BitmapInfo (optional) A list of qcow2 bitmap details (since 4.0)

compression-type: Qcow2CompressionType the image cluster compression method (since 5.1)

Since

1.7

ImageInfoSpecificVmdk (Object)

Members

create-type: string The create type of VMDK image

cid: int Content id of image

parent-cid: int Parent VMDK image's cid

extents: array of ImageInfo List of extent files

Since

1.7

ImageInfoSpecific (Object)

A discriminated record of image format specific information structures.

Members

type One of qcow2, vmdk, luks

data: ImageInfoSpecificQCow2 when type is "qcow2"

data: ImageInfoSpecificVmdk when type is "vmdk"

data: QCryptoBlockInfoLUKS when type is "luks"

Since

1.7

ImageInfo (Object)

Information about a QEMU image file

Members

filename: string name of the image file
format: string format of the image file
virtual-size: int maximum capacity in bytes of the image
actual-size: int (optional) actual size on disk in bytes of the image
dirty-flag: boolean (optional) true if image is not cleanly closed
cluster-size: int (optional) size of a cluster in bytes
encrypted: boolean (optional) true if the image is encrypted
compressed: boolean (optional) true if the image is compressed (Since 1.7)
backing-filename: string (optional) name of the backing file
full-backing-filename: string (optional) full path of the backing file
backing-filename-format: string (optional) the format of the backing file
snapshots: array of SnapshotInfo (optional) list of VM snapshots
backing-image: ImageInfo (optional) info of the backing image (since 1.6)
format-specific: ImageInfoSpecific (optional) structure supplying additional format-specific information (since 1.7)

Since

1.3

ImageCheck (Object)

Information about a QEMU image file check

Members

filename: string name of the image file checked
format: string format of the image file checked
check-errors: int number of unexpected errors occurred during check
image-end-offset: int (optional) offset (in bytes) where the image ends, this field is present if the driver for the image format supports it
corruptions: int (optional) number of corruptions found during the check if any
leaks: int (optional) number of leaks found during the check if any
corruptions-fixed: int (optional) number of corruptions fixed during the check if any
leaks-fixed: int (optional) number of leaks fixed during the check if any
total-clusters: int (optional) total number of clusters, this field is present if the driver for the image format supports it

allocated-clusters: int (optional) total number of allocated clusters, this field is present if the driver for the image format supports it

fragmented-clusters: int (optional) total number of fragmented clusters, this field is present if the driver for the image format supports it

compressed-clusters: int (optional) total number of compressed clusters, this field is present if the driver for the image format supports it

Since

1.4

MapEntry (Object)

Mapping information from a virtual block range to a host file range

Members

start: int virtual (guest) offset of the first byte described by this entry

length: int the number of bytes of the mapped virtual range

data: boolean reading the image will actually read data from a file (in particular, if `offset` is present this means that the sectors are not simply preallocated, but contain actual data in raw format)

zero: boolean whether the virtual blocks read as zeroes

depth: int number of layers (0 = top image, 1 = top image's backing file, ..., n - 1 = bottom image (where n is the number of images in the chain)) before reaching one for which the range is allocated

offset: int (optional) if present, the image file stores the data for this range in raw format at the given (host) offset

filename: string (optional) filename that is referred to by `offset`

Since

2.6

BlockdevCacheInfo (Object)

Cache mode information for a block device

Members

writeback: boolean true if writeback mode is enabled

direct: boolean true if the host page cache is bypassed (O_DIRECT)

no-flush: boolean true if flush requests are ignored for the device

Since

2.3

BlockDeviceInfo (Object)

Information about the backing device for a block device.

Members

file: string the filename of the backing device

node-name: string (optional) the name of the block driver node (Since 2.0)

ro: boolean true if the backing device was open read-only

drv: string the name of the block format used to open the backing device. As of 0.14 this can be: 'blkdebug', 'bochs', 'cloop', 'cow', 'dmg', 'file', 'file', 'ftp', 'ftps', 'host_cdrom', 'host_device', 'http', 'https', 'luks', 'nbd', 'parallels', 'qcow', 'qcow2', 'raw', 'vdi', 'vmdk', 'vpc', 'vvfat' 2.2: 'archipelago' added, 'cow' dropped 2.3: 'host_floppy' deprecated 2.5: 'host_floppy' dropped 2.6: 'luks' added 2.8: 'replication' added, 'tftp' dropped 2.9: 'archipelago' dropped

backing_file: string (optional) the name of the backing file (for copy-on-write)

backing_file_depth: int number of files in the backing file chain (since: 1.2)

encrypted: boolean true if the backing device is encrypted

encryption_key_missing: boolean always false

detect zeroes: BlockdevDetectZeroesOptions detect and optimize zero writes (Since 2.1)

bps: int total throughput limit in bytes per second is specified

bps_rd: int read throughput limit in bytes per second is specified

bps_wr: int write throughput limit in bytes per second is specified

iops: int total I/O operations per second is specified

iops_rd: int read I/O operations per second is specified

iops_wr: int write I/O operations per second is specified

image: ImageInfo the info of image used (since: 1.6)

bps_max: int (optional)

total throughput limit during bursts, in bytes (Since 1.7)

bps_rd_max: int (optional)

read throughput limit during bursts, in bytes (Since 1.7)

bps_wr_max: int (optional)

write throughput limit during bursts, in bytes (Since 1.7)

iops_max: int (optional)

total I/O operations per second during bursts, in bytes (Since 1.7)

iops_rd_max: int (optional)

read I/O operations per second during bursts, in bytes (Since 1.7)

iops_wr_max: int (optional)

write I/O operations per second during bursts, in bytes (Since 1.7)

bps_max_length: int (optional)

maximum length of the bps_max burst period, in seconds. (Since 2.6)

bps_rd_max_length: int (optional)

maximum length of the bps_rd_max burst period, in seconds. (Since 2.6)

bps_wr_max_length: int (optional)

maximum length of the bps_wr_max burst period, in seconds. (Since 2.6)

iops_max_length: int (optional)

maximum length of the iops burst period, in seconds. (Since 2.6)

iops_rd_max_length: int (optional)

maximum length of the iops_rd_max burst period, in seconds. (Since 2.6)

iops_wr_max_length: int (optional)

maximum length of the iops_wr_max burst period, in seconds. (Since 2.6)

iops_size: int (optional) an I/O size in bytes (Since 1.7)

group: string (optional) throttle group name (Since 2.4)

cache: BlockdevCacheInfo the cache mode used for the block device (since: 2.3)

write_threshold: int configured write threshold for the device. 0 if disabled. (Since 2.3)

dirty-bitmaps: array of BlockDirtyInfo (optional) dirty bitmaps information (only present if node has one or more dirty bitmaps) (Since 4.2)

Features

deprecated Member `encryption_key_missing` is deprecated. It is always false.

Since

0.14

BlockDeviceIoStatus (Enum)

An enumeration of block device I/O status.

Values

ok The last I/O operation has succeeded

failed The last I/O operation has failed

nospace The last I/O operation has failed due to a no-space condition

Since

1.0

DirtyBitmapStatus (Enum)

An enumeration of possible states that a dirty bitmap can report to the user.

Values

frozen The bitmap is currently in-use by some operation and is immutable. If the bitmap was `active` prior to the operation, new writes by the guest are being recorded in a temporary buffer, and will not be lost. Generally, bitmaps are cleared on successful use in an operation and the temporary buffer is committed into the bitmap. On failure, the temporary buffer is merged back into the bitmap without first clearing it. Please refer to the documentation for each bitmap-using operation, See also `blockdev-backup`, `drive-backup`.

disabled The bitmap is not currently recording new writes by the guest. This is requested explicitly via `block-dirty-bitmap-disable`. It can still be cleared, deleted, or used for backup operations.

active The bitmap is actively monitoring for new writes, and can be cleared, deleted, or used for backup operations.

locked The bitmap is currently in-use by some operation and is immutable. If the bitmap was `active` prior to the operation, it is still recording new writes. If the bitmap was `disabled`, it is not recording new writes. (Since 2.12)

inconsistent This is a persistent dirty bitmap that was marked in-use on disk, and is unusable by QEMU. It can only be deleted. Please rely on the `inconsistent` field in `BlockDirtyInfo` instead, as the `status` field is deprecated. (Since 4.0)

Since

2.4

BlockDirtyInfo (Object)

Block dirty bitmap information.

Members

name: string (optional) the name of the dirty bitmap (Since 2.4)

count: int number of dirty bytes according to the dirty bitmap

granularity: int granularity of the dirty bitmap in bytes (since 1.4)

status: DirtyBitmapStatus current status of the dirty bitmap (since 2.4)

recording: boolean true if the bitmap is recording new writes from the guest. Replaces *active* and *disabled* statuses. (since 4.0)

busy: boolean true if the bitmap is in-use by some operation (NBD or jobs) and cannot be modified via QMP or used by another operation. Replaces *locked* and *frozen* statuses. (since 4.0)

persistent: boolean true if the bitmap was stored on disk, is scheduled to be stored on disk, or both. (since 4.0)

inconsistent: boolean (optional) true if this is a persistent bitmap that was improperly stored. Implies `persistent` to be true; `recording` and `busy` to be false. This bitmap cannot be used. To remove it, use `block-dirty-bitmap-remove`. (Since 4.0)

Features

deprecated Member `status` is deprecated. Use `recording` and `locked` instead.

Since

1.3

Qcow2BitmapInfoFlags (Enum)

An enumeration of flags that a bitmap can report to the user.

Values

in-use This flag is set by any process actively modifying the qcow2 file, and cleared when the updated bitmap is flushed to the qcow2 image. The presence of this flag in an offline image means that the bitmap was not saved correctly after its last usage, and may contain inconsistent data.

auto The bitmap must reflect all changes of the virtual disk by any application that would write to this qcow2 file.

Since

4.0

Qcow2BitmapInfo (Object)

Qcow2 bitmap information.

Members

name: string the name of the bitmap

granularity: int granularity of the bitmap in bytes

flags: array of Qcow2BitmapInfoFlags flags of the bitmap

Since

4.0

BlockLatencyHistogramInfo (Object)

Block latency histogram.

Members

boundaries: array of int list of interval boundary values in nanoseconds, all greater than zero and in ascending order. For example, the list [10, 50, 100] produces the following histogram intervals: [0, 10), [10, 50), [50, 100), [100, +inf).

bins: array of int list of io request counts corresponding to histogram intervals. `len(bins) = len(boundaries) + 1` For the example above, bins may be something like [3, 1, 5, 2], and corresponding histogram looks like:



Since

4.0

BlockInfo (Object)

Block device information. This structure describes a virtual device and the backing device associated with it.

Members

device: string The device name associated with the virtual device.

qdev: string (optional) The qdev ID, or if no ID is assigned, the QOM path of the block device. (since 2.10)

type: string This field is returned only for compatibility reasons, it should not be used (always returns 'unknown')

removable: boolean True if the device supports removable media.

locked: boolean True if the guest has locked this device from having its media removed

tray_open: boolean (optional) True if the device's tray is open (only present if it has a tray)

dirty-bitmaps: array of BlockDirtyInfo (optional) dirty bitmaps information (only present if the driver has one or more dirty bitmaps) (Since 2.0)

io-status: BlockDeviceIoStatus (optional) BlockDeviceIoStatus. Only present if the device supports it and the VM is configured to stop on errors (supported device models: virtio-blk, IDE, SCSI except scsi-generic)

inserted: BlockDeviceInfo (optional) BlockDeviceInfo describing the device if media is present

Features

deprecated Member `dirty-bitmaps` is deprecated. Use inserted member `dirty-bitmaps` instead.

Since

0.14

BlockMeasureInfo (Object)

Image file size calculation information. This structure describes the size requirements for creating a new image file.

The size requirements depend on the new image file format. File size always equals virtual disk size for the ‘raw’ format, even for sparse POSIX files. Compact formats such as ‘qcow2’ represent unallocated and zero regions efficiently so file size may be smaller than virtual disk size.

The values are upper bounds that are guaranteed to fit the new image file. Subsequent modification, such as internal snapshot or further bitmap creation, may require additional space and is not covered here.

Members

required: int Size required for a new image file, in bytes, when copying just allocated guest-visible contents.

fully-allocated: int Image file size, in bytes, once data has been written to all sectors, when copying just guest-visible contents.

bitmaps: int (optional) Additional size required if all the top-level bitmap metadata in the source image were to be copied to the destination, present only when source and destination both support persistent bitmaps. (since 5.1)

Since

2.10

query-block (Command)

Get a list of `BlockInfo` for all virtual block devices.

Returns

a list of `BlockInfo` describing each virtual block device. Filter nodes that were created implicitly are skipped over.

Since

0.14

Example

```
-> { "execute": "query-block" }
<- {
  "return": [
    {
      "io-status": "ok",
      "device": "ide0-hd0",
      "locked": false,
      "removable": false,
      "inserted": {
        "ro": false,
        "drv": "qcow2",
        "encrypted": false,
        "file": "disks/test.qcow2",
        "backing_file_depth": 1,
        "bps": 1000000,
        "bps_rd": 0,
        "bps_wr": 0,
        "iops": 1000000,
        "iops_rd": 0,
        "iops_wr": 0,
        "bps_max": 8000000,
        "bps_rd_max": 0,
        "bps_wr_max": 0,
        "iops_max": 0,
        "iops_rd_max": 0,
        "iops_wr_max": 0,
        "iops_size": 0,
        "detect_zeroes": "on",
        "write_threshold": 0,
        "image": {
          "filename": "disks/test.qcow2",
          "format": "qcow2",
          "virtual-size": 2048000,
          "backing_file": "base.qcow2",
          "full-backing-filename": "disks/base.qcow2",
          "backing-filename-format": "qcow2",
          "snapshots": [
            {
              "id": "1",
              "name": "snapshot1",
              "vm-state-size": 0,
              "date-sec": 10000200,
              "date-nsec": 12,
              "vm-clock-sec": 206,
              "vm-clock-nsec": 30
            }
          ]
        },
        "backing-image": {
          "filename": "disks/base.qcow2",
          "format": "qcow2",
          "virtual-size": 2048000
        }
      }
    },
    { "qdev": "ide_disk",
```

(continues on next page)

(continued from previous page)

```

        "type": "unknown"
    },
    {
        "io-status": "ok",
        "device": "ide1-cd0",
        "locked": false,
        "removable": true,
        "qdev": "/machine/unattached/device[23]",
        "tray_open": false,
        "type": "unknown"
    },
    {
        "device": "floppy0",
        "locked": false,
        "removable": true,
        "qdev": "/machine/unattached/device[20]",
        "type": "unknown"
    },
    {
        "device": "sd0",
        "locked": false,
        "removable": true,
        "type": "unknown"
    }
]
}

```

BlockDeviceTimedStats (Object)

Statistics of a block device during a given interval of time.

Members

interval_length: int Interval used for calculating the statistics, in seconds.

min_rd_latency_ns: int Minimum latency of read operations in the defined interval, in nanoseconds.

min_wr_latency_ns: int Minimum latency of write operations in the defined interval, in nanoseconds.

min_flush_latency_ns: int Minimum latency of flush operations in the defined interval, in nanoseconds.

max_rd_latency_ns: int Maximum latency of read operations in the defined interval, in nanoseconds.

max_wr_latency_ns: int Maximum latency of write operations in the defined interval, in nanoseconds.

max_flush_latency_ns: int Maximum latency of flush operations in the defined interval, in nanoseconds.

avg_rd_latency_ns: int Average latency of read operations in the defined interval, in nanoseconds.

avg_wr_latency_ns: int Average latency of write operations in the defined interval, in nanoseconds.

avg_flush_latency_ns: int Average latency of flush operations in the defined interval, in nanoseconds.

avg_rd_queue_depth: number Average number of pending read operations in the defined interval.

avg_wr_queue_depth: number Average number of pending write operations in the defined interval.

Since

2.5

BlockDeviceStats (Object)

Statistics of a virtual block device or a block backing device.

Members

rd_bytes: int The number of bytes read by the device.

wr_bytes: int The number of bytes written by the device.

unmap_bytes: int The number of bytes unmapped by the device (Since 4.2)

rd_operations: int The number of read operations performed by the device.

wr_operations: int The number of write operations performed by the device.

flush_operations: int The number of cache flush operations performed by the device (since 0.15)

unmap_operations: int The number of unmap operations performed by the device (Since 4.2)

rd_total_time_ns: int Total time spent on reads in nanoseconds (since 0.15).

wr_total_time_ns: int Total time spent on writes in nanoseconds (since 0.15).

flush_total_time_ns: int Total time spent on cache flushes in nanoseconds (since 0.15).

unmap_total_time_ns: int Total time spent on unmap operations in nanoseconds (Since 4.2)

wr_highest_offset: int The offset after the greatest byte written to the device. The intended use of this information is for growable sparse files (like qcow2) that are used on top of a physical device.

rd_merged: int Number of read requests that have been merged into another request (Since 2.3).

wr_merged: int Number of write requests that have been merged into another request (Since 2.3).

unmap_merged: int Number of unmap requests that have been merged into another request (Since 4.2)

idle_time_ns: int (optional) Time since the last I/O operation, in nanoseconds. If the field is absent it means that there haven't been any operations yet (Since 2.5).

failed_rd_operations: int The number of failed read operations performed by the device (Since 2.5)

failed_wr_operations: int The number of failed write operations performed by the device (Since 2.5)

failed_flush_operations: int The number of failed flush operations performed by the device (Since 2.5)

failed_unmap_operations: int The number of failed unmap operations performed by the device (Since 4.2)

invalid_rd_operations: int

The number of invalid read operations performed by the device (Since 2.5)

invalid_wr_operations: int The number of invalid write operations performed by the device (Since 2.5)

invalid_flush_operations: int The number of invalid flush operations performed by the device (Since 2.5)

invalid_unmap_operations: int The number of invalid unmap operations performed by the device (Since 4.2)

account_invalid: boolean Whether invalid operations are included in the last access statistics (Since 2.5)

account_failed: boolean Whether failed operations are included in the latency and last access statistics (Since 2.5)

timed_stats: array of BlockDeviceTimedStats Statistics specific to the set of previously defined intervals of time (Since 2.5)

rd_latency_histogram: BlockLatencyHistogramInfo (optional) BlockLatencyHistogramInfo. (Since 4.0)

wr_latency_histogram: BlockLatencyHistogramInfo (optional) BlockLatencyHistogramInfo. (Since 4.0)

flush_latency_histogram: BlockLatencyHistogramInfo (optional) BlockLatencyHistogramInfo. (Since 4.0)

Since

0.14

BlockStatsSpecificFile (Object)

File driver statistics

Members

discard-nb-ok: int The number of successful discard operations performed by the driver.

discard-nb-failed: int The number of failed discard operations performed by the driver.

discard-bytes-ok: int The number of bytes discarded by the driver.

Since

4.2

BlockStatsSpecificNvme (Object)

NVMe driver statistics

Members

completion-errors: int The number of completion errors.

aligned-accesses: int The number of aligned accesses performed by the driver.

unaligned-accesses: int The number of unaligned accesses performed by the driver.

Since

5.2

BlockStatsSpecific (Object)

Block driver specific statistics

Members

driver: `BlockdevDriver` Not documented

The members of `BlockStatsSpecificFile` when `driver` is "file"

The members of `BlockStatsSpecificFile` when `driver` is "host_device"

The members of `BlockStatsSpecificNvme` when `driver` is "nvme"

Since

4.2

BlockStats (Object)

Statistics of a virtual block device or a block backing device.

Members

device: `string` (optional) If the stats are for a virtual block device, the name corresponding to the virtual block device.

node-name: `string` (optional) The node name of the device. (Since 2.3)

qdev: `string` (optional) The qdev ID, or if no ID is assigned, the QOM path of the block device. (since 3.0)

stats: `BlockDeviceStats` A `BlockDeviceStats` for the device.

driver-specific: `BlockStatsSpecific` (optional) Optional driver-specific stats. (Since 4.2)

parent: `BlockStats` (optional) This describes the file block device if it has one. Contains recursively the statistics of the underlying protocol (e.g. the host file for a qcow2 image). If there is no underlying protocol, this field is omitted

backing: `BlockStats` (optional) This describes the backing block device if it has one. (Since 2.0)

Since

0.14

query-blockstats (Command)

Query the `BlockStats` for all virtual block devices.

Arguments

query-nodes: boolean (optional) If true, the command will query all the block nodes that have a node name, in a list which will include “parent” information, but not “backing”. If false or omitted, the behavior is as before - query all the device backends, recursively including their “parent” and “backing”. Filter nodes that were created implicitly are skipped over in this mode. (Since 2.3)

Returns

A list of `BlockStats` for each virtual block devices.

Since

0.14

Example

```
-> { "execute": "query-blockstats" }
<- {
  "return":[
    {
      "device":"ide0-hd0",
      "parent":{
        "stats":{
          "wr_highest_offset":3686448128,
          "wr_bytes":9786368,
          "wr_operations":751,
          "rd_bytes":122567168,
          "rd_operations":36772
          "wr_total_times_ns":313253456
          "rd_total_times_ns":3465673657
          "flush_total_times_ns":49653
          "flush_operations":61,
          "rd_merged":0,
          "wr_merged":0,
          "idle_time_ns":2953431879,
          "account_invalid":true,
          "account_failed":false
        }
      },
      "stats":{
        "wr_highest_offset":2821110784,
        "wr_bytes":9786368,
        "wr_operations":692,
        "rd_bytes":122739200,
        "rd_operations":36604
        "flush_operations":51,
        "wr_total_times_ns":313253456
        "rd_total_times_ns":3465673657
        "flush_total_times_ns":49653,
        "rd_merged":0,
        "wr_merged":0,
```

(continues on next page)

(continued from previous page)

```

        "idle_time_ns":2953431879,
        "account_invalid":true,
        "account_failed":false
    },
    "qdev": "/machine/unattached/device[23]"
},
{
    "device":"ide1-cd0",
    "stats":{
        "wr_highest_offset":0,
        "wr_bytes":0,
        "wr_operations":0,
        "rd_bytes":0,
        "rd_operations":0,
        "flush_operations":0,
        "wr_total_times_ns":0,
        "rd_total_times_ns":0,
        "flush_total_times_ns":0,
        "rd_merged":0,
        "wr_merged":0,
        "account_invalid":false,
        "account_failed":false
    },
    "qdev": "/machine/unattached/device[24]"
},
{
    "device":"floppy0",
    "stats":{
        "wr_highest_offset":0,
        "wr_bytes":0,
        "wr_operations":0,
        "rd_bytes":0,
        "rd_operations":0,
        "flush_operations":0,
        "wr_total_times_ns":0,
        "rd_total_times_ns":0,
        "flush_total_times_ns":0,
        "rd_merged":0,
        "wr_merged":0,
        "account_invalid":false,
        "account_failed":false
    },
    "qdev": "/machine/unattached/device[16]"
},
{
    "device":"sd0",
    "stats":{
        "wr_highest_offset":0,
        "wr_bytes":0,
        "wr_operations":0,
        "rd_bytes":0,
        "rd_operations":0,
        "flush_operations":0,
        "wr_total_times_ns":0,
        "rd_total_times_ns":0,
        "flush_total_times_ns":0,
        "rd_merged":0,

```

(continues on next page)

(continued from previous page)

```

        "wr_merged":0,
        "account_invalid":false,
        "account_failed":false
    }
}
]
}

```

BlockdevOnError (Enum)

An enumeration of possible behaviors for errors on I/O operations. The exact meaning depends on whether the I/O was initiated by a guest or by a block job

Values

report for guest operations, report the error to the guest; for jobs, cancel the job

ignore ignore the error, only report a QMP event (BLOCK_IO_ERROR or BLOCK_JOB_ERROR). The backup, mirror and commit block jobs retry the failing request later and may still complete successfully. The stream block job continues to stream and will complete with an error.

enospc same as **stop** on ENOSPC, same as **report** otherwise.

stop for guest operations, stop the virtual machine; for jobs, pause the job

auto inherit the error handling policy of the backend (since: 2.7)

Since

1.3

MirrorSyncMode (Enum)

An enumeration of possible behaviors for the initial synchronization phase of storage mirroring.

Values

top copies data in the topmost image to the destination

full copies data from all images to the destination

none only copy data written from now on

incremental only copy data described by the dirty bitmap. (since: 2.4)

bitmap only copy data described by the dirty bitmap. (since: 4.2) Behavior on completion is determined by the BitmapSyncMode.

Since

1.3

BitmapSyncMode (Enum)

An enumeration of possible behaviors for the synchronization of a bitmap when used for data copy operations.

Values

on-success The bitmap is only synced when the operation is successful. This is the behavior always used for ‘INCREMENTAL’ backups.

never The bitmap is never synchronized with the operation, and is treated solely as a read-only manifest of blocks to copy.

always The bitmap is always synchronized with the operation, regardless of whether or not the operation was successful.

Since

4.2

MirrorCopyMode (Enum)

An enumeration whose values tell the mirror block job when to trigger writes to the target.

Values

background copy data in background only.

write-blocking when data is written to the source, write it (synchronously) to the target as well. In addition, data is copied in background just like in `background` mode.

Since

3.0

BlockJobInfo (Object)

Information about a long-running block device operation.

Members

type: string the job type (‘stream’ for image streaming)

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: int Estimated `offset` value at the completion of the job. This value can arbitrarily change while the job is running, in both directions.

offset: int Progress made until now. The unit is arbitrary and the value can only meaningfully be used for the ratio of `offset` to `len`. The value is monotonically increasing.

busy: boolean false if the job is known to be in a quiescent state, with no pending I/O. Since 1.3.

paused: boolean whether the job is paused or, if `busy` is true, will pause itself as soon as possible. Since 1.3.

speed: int the rate limit, bytes per second

io-status: BlockDeviceIoStatus the status of the job (since 1.3)

ready: boolean true if the job may be completed (since 2.2)

status: JobStatus Current job state/status (since 2.12)

auto-finalize: boolean Job will finalize itself when PENDING, moving to the CONCLUDED state. (since 2.12)

auto-dismiss: boolean Job will dismiss itself when CONCLUDED, moving to the NULL state and disappearing from the query list. (since 2.12)

error: string (optional) Error information if the job did not complete successfully. Not set if the job completed successfully. (since 2.12.1)

Since

1.1

query-block-jobs (Command)

Return information about long-running block device operations.

Returns

a list of `BlockJobInfo` for each active block job

Since

1.1

block_passwd (Command)

This command sets the password of a block device that has not been open with a password and requires one.

This command is now obsolete and will always return an error since 2.10

Arguments

device: string (optional) Not documented

node-name: string (optional) Not documented

password: string Not documented

block_resize (Command)

Resize a block image while a guest is running.

Either `device` or `node-name` must be set but not both.

Arguments

device: string (optional) the name of the device to get the image resized

node-name: string (optional) graph node name to get the image resized (Since 2.0)

size: int new image size in bytes

Returns

- nothing on success
- If `device` is not a valid block device, `DeviceNotFound`

Since

0.14

Example

```
-> { "execute": "block_resize",  
      "arguments": { "device": "scratch", "size": 1073741824 } }  
<- { "return": {} }
```

NewImageMode (Enum)

An enumeration that tells QEMU how to set the backing file path in a new image file.

Values

existing QEMU should look for an existing image file.

absolute-paths QEMU should create a new image with absolute paths for the backing file. If there is no backing file available, the new image will not be backed either.

Since

1.1

BlockdevSnapshotSync (Object)

Either `device` or `node-name` must be set but not both.

Members

device: string (optional) the name of the device to take a snapshot of.

node-name: string (optional) graph node name to generate the snapshot from (Since 2.0)

snapshot-file: string the target of the new overlay image. If the file exists, or if it is a device, the overlay will be created in the existing file/device. Otherwise, a new file will be created.

snapshot-node-name: string (optional) the graph node name of the new image (Since 2.0)

format: string (optional) the format of the overlay image, default is ‘qcow2’.

mode: NewImageMode (optional) whether and how QEMU should create a new image, default is ‘absolute-paths’.

BlockdevSnapshot (Object)

Members

node: string device or node name that will have a snapshot taken.

overlay: string reference to the existing block device that will become the overlay of `node`, as part of taking the snapshot. It must not have a current backing file (this can be achieved by passing “backing”: null to `blockdev-add`).

Since

2.5

BackupPerf (Object)

Optional parameters for backup. These parameters don’t affect functionality, but may significantly affect performance.

Members

use-copy-range: boolean (optional) Use copy offloading. Default false.

max-workers: int (optional) Maximum number of parallel requests for the sustained background copying process. Doesn’t influence copy-before-write operations. Default 64.

max-chunk: int (optional) Maximum request length for the sustained background copying process. Doesn’t influence copy-before-write operations. 0 means unlimited. If max-chunk is non-zero then it should not be less than job cluster size which is calculated as maximum of target image cluster size and 64k. Default 0.

Since

6.0

BackupCommon (Object)

Members

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device name or node-name of a root node which should be copied.

sync: MirrorSyncMode what parts of the disk image should be copied to the destination (all the disk, only the sectors allocated in the topmost image, from a dirty bitmap, or only new I/O).

speed: int (optional) the maximum speed, in bytes per second. The default is 0, for unlimited.

bitmap: string (optional) The name of a dirty bitmap to use. Must be present if sync is “bitmap” or “incremental”. Can be present if sync is “full” or “top”. Must not be present otherwise. (Since 2.4 (drive-backup), 3.1 (blockdev-backup))

bitmap-mode: BitmapSyncMode (optional) Specifies the type of data the bitmap should contain after the operation concludes. Must be present if a bitmap was provided, Must NOT be present otherwise. (Since 4.2)

compress: boolean (optional) true to compress data, if the target format supports it. (default: false) (since 2.8)

on-source-error: BlockdevOnError (optional) the action to take on an error on the source, default ‘report’. ‘stop’ and ‘enospc’ can only be used if the block device supports io-status (see BlockInfo).

on-target-error: BlockdevOnError (optional) the action to take on an error on the target, default ‘report’ (no limitations, since this applies to a different block device than `device`).

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 2.12)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 2.12)

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the backup job inserts into the graph above node specified by `drive`. If this option is not given, a node name is autogenerated. (Since: 4.2)

x-perf: BackupPerf (optional) Performance options. (Since 6.0)

Note

`on-source-error` and `on-target-error` only affect background I/O. If an error occurs during a guest write request, the device’s `error/werror` actions will be used.

Since

4.2

DriveBackup (Object)

Members

target: string the target of the new image. If the file exists, or if it is a device, the existing file/device will be used as the new destination. If it does not exist, a new file will be created.

format: string (optional) the format of the new destination, default is to probe if `mode` is ‘existing’, else the format of the source

mode: NewImageMode (optional) whether and how QEMU should create a new image, default is ‘absolute-paths’.

The members of `BackupCommon`

Since

1.6

BlockdevBackup (Object)

Members

target: string the device name or node-name of the backup target node.

The members of `BackupCommon`

Since

2.3

blockdev-snapshot-sync (Command)

Takes a synchronous snapshot of a block device.

For the arguments, see the documentation of `BlockdevSnapshotSync`.

Returns

- nothing on success
- If `device` is not a valid block device, `DeviceNotFound`

Since

0.14

Example

```
-> { "execute": "blockdev-snapshot-sync",
      "arguments": { "device": "ide-hd0",
                     "snapshot-file":
                       "/some/place/my-image",
                     "format": "qcow2" } }
<- { "return": {} }
```

blockdev-snapshot (Command)

Takes a snapshot of a block device.

Take a snapshot, by installing ‘node’ as the backing image of ‘overlay’. Additionally, if ‘node’ is associated with a block device, the block device changes to using ‘overlay’ as its new active image.

For the arguments, see the documentation of BlockdevSnapshot.

Features

allow-write-only-overlay If present, the check whether this operation is safe was relaxed so that it can be used to change backing file of a destination of a blockdev-mirror. (since 5.0)

Since

2.5

Example

```
-> { "execute": "blockdev-add",
      "arguments": { "driver": "qcow2",
                     "node-name": "node1534",
                     "file": { "driver": "file",
                               "filename": "hd1.qcow2" },
                     "backing": null } }
<- { "return": {} }

-> { "execute": "blockdev-snapshot",
      "arguments": { "node": "ide-hd0",
                     "overlay": "node1534" } }
<- { "return": {} }
```

change-backing-file (Command)

Change the backing file in the image file metadata. This does not cause QEMU to reopen the image file to reparse the backing filename (it may, however, perform a reopen to change permissions from r/o -> r/w -> r/o, if needed). The new backing file string is written into the image file metadata, and the QEMU internal strings are updated.

Arguments

image-node-name: string The name of the block driver state node of the image to modify. The “device” argument is used to verify “image-node-name” is in the chain described by “device”.

device: string The device name or node-name of the root node that owns image-node-name.

backing-file: string The string to write as the backing file. This string is not validated, so care should be taken when specifying the string or the image chain may not be able to be reopened again.

Returns

- Nothing on success
- If “device” does not exist or cannot be determined, DeviceNotFound

Since

2.1

block-commit (Command)

Live commit of data from overlay image nodes into backing nodes - i.e., writes data between ‘top’ and ‘base’ into ‘base’.

If top == base, that is an error. If top has no overlays on top of it, or if it is in use by a writer, the job will not be completed by itself. The user needs to complete the job with the block-job-complete command after getting the ready event. (Since 2.0)

If the base image is smaller than top, then the base image will be resized to be the same size as top. If top is smaller than the base image, the base will not be truncated. If you want the base image size to match the size of the smaller top, you can safely truncate it yourself once the commit operation successfully completes.

Arguments

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device name or node-name of a root node

base-node: string (optional) The node name of the backing image to write data into. If not specified, this is the deepest backing image. (since: 3.1)

base: string (optional) Same as base-node, except that it is a file name rather than a node name. This must be the exact filename string that was used to open the node; other strings, even if addressing the same file, are not accepted

top-node: string (optional) The node name of the backing image within the image chain which contains the topmost data to be committed down. If not specified, this is the active layer. (since: 3.1)

top: string (optional) Same as top-node, except that it is a file name rather than a node name. This must be the exact filename string that was used to open the node; other strings, even if addressing the same file, are not accepted

backing-file: string (optional) The backing file string to write into the overlay image of 'top'. If 'top' does not have an overlay image, or if 'top' is in use by a writer, specifying a backing file string is an error.

This filename is not validated. If a pathname string is such that it cannot be resolved by QEMU, that means that subsequent QMP or HMP commands must use node-names for the image in question, as filename lookup methods will fail.

If not specified, QEMU will automatically determine the backing file string to use, or error out if there is no obvious choice. Care should be taken when specifying the string, to specify a valid filename or protocol. (Since 2.1)

speed: int (optional) the maximum speed, in bytes per second

on-error: BlockdevOnError (optional) the action to take on an error. 'ignore' means that the request should be retried. (default: report; Since: 5.0)

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the commit job inserts into the graph above top. If this option is not given, a node name is autogenerated. (Since: 2.9)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for block-job-finalize before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits block-job-dismiss. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Features

deprecated Members `base` and `top` are deprecated. Use `base-node` and `top-node` instead.

Returns

- Nothing on success
- If `device` does not exist, `DeviceNotFound`
- Any other error returns a `GenericError`.

Since

1.3

Example

```
-> { "execute": "block-commit",  
    "arguments": { "device": "virtio0",  
                  "top": "/tmp/snap1.qcow2" } }  
<- { "return": {} }
```


drive-backup (Command)

Start a point-in-time copy of a block device to a new destination. The status of ongoing drive-backup operations can be checked with `query-block-jobs` where the `BlockJobInfo.type` field has the value `'backup'`. The operation can be stopped before it has completed using the `block-job-cancel` command.

Arguments

The members of DriveBackup

Returns

- nothing on success
- If `device` is not a valid block device, `GenericError`

Since

1.6

Example

```
-> { "execute": "drive-backup",
      "arguments": { "device": "drive0",
                     "sync": "full",
                     "target": "backup.img" } }
<- { "return": {} }
```

blockdev-backup (Command)

Start a point-in-time copy of a block device to a new destination. The status of ongoing blockdev-backup operations can be checked with `query-block-jobs` where the `BlockJobInfo.type` field has the value `'backup'`. The operation can be stopped before it has completed using the `block-job-cancel` command.

Arguments

The members of BlockdevBackup

Returns

- nothing on success
- If `device` is not a valid block device, `DeviceNotFound`

Since

2.3

Example

```
-> { "execute": "blockdev-backup",  
    "arguments": { "device": "src-id",  
                  "sync": "full",  
                  "target": "tgt-id" } }  
<- { "return": {} }
```

query-named-block-nodes (Command)

Get the named block driver list

Arguments

flat: boolean (optional) Omit the nested data about backing image (“backing-image” key) if true. Default is false (Since 5.0)

Returns

the list of BlockDeviceInfo

Since

2.0

Example

```
-> { "execute": "query-named-block-nodes" }  
<- { "return": [ { "ro":false,  
                  "drv":"qcow2",  
                  "encrypted":false,  
                  "file":"disks/test.qcow2",  
                  "node-name": "my-node",  
                  "backing_file_depth":1,  
                  "bps":1000000,  
                  "bps_rd":0,  
                  "bps_wr":0,  
                  "iops":1000000,  
                  "iops_rd":0,  
                  "iops_wr":0,  
                  "bps_max": 8000000,  
                  "bps_rd_max": 0,  
                  "bps_wr_max": 0,  
                  "iops_max": 0,  
                  "iops_rd_max": 0,  
                  "iops_wr_max": 0,  
                  "iops_size": 0,  
                  "write_threshold": 0,  
                  "image":{
```

(continues on next page)

(continued from previous page)

```

    "filename": "disks/test.qcow2",
    "format": "qcow2",
    "virtual-size": 2048000,
    "backing_file": "base.qcow2",
    "full-backing-filename": "disks/base.qcow2",
    "backing-filename-format": "qcow2",
    "snapshots": [
        {
            "id": "1",
            "name": "snapshot1",
            "vm-state-size": 0,
            "date-sec": 10000200,
            "date-nsec": 12,
            "vm-clock-sec": 206,
            "vm-clock-nsec": 30
        }
    ],
    "backing-image": {
        "filename": "disks/base.qcow2",
        "format": "qcow2",
        "virtual-size": 2048000
    }
} } ] }

```

XDbgBlockGraphNodeType (Enum)

Values

block-backend corresponds to BlockBackend

block-job corresponds to BlockJob

block-driver corresponds to BlockDriverState

Since

4.0

XDbgBlockGraphNode (Object)

Members

id: **int** Block graph node identifier. This id is generated only for x-debug-query-block-graph and does not relate to any other identifiers in Qemu.

type: **XDbgBlockGraphNodeType** Type of graph node. Can be one of block-backend, block-job or block-driver-state.

name: **string** Human readable name of the node. Corresponds to node-name for block-driver-state nodes; is not guaranteed to be unique in the whole graph (with block-jobs and block-backends).

Since

4.0

BlockPermission (Enum)

Enum of base block permissions.

Values

consistent-read A user that has the “permission” of consistent reads is guaranteed that their view of the contents of the block device is complete and self-consistent, representing the contents of a disk at a specific point. For most block devices (including their backing files) this is true, but the property cannot be maintained in a few situations like for intermediate nodes of a commit block job.

write This permission is required to change the visible disk contents.

write-unchanged This permission (which is weaker than `BLK_PERM_WRITE`) is both enough and required for writes to the block node when the caller promises that the visible disk content doesn’t change. As the `BLK_PERM_WRITE` permission is strictly stronger, either is sufficient to perform an unchanging write.

resize This permission is required to change the size of a block node.

graph-mod This permission is required to change the node that this `BdrvChild` points to.

Since

4.0

XDbgBlockGraphEdge (Object)

Block Graph edge description for `x-debug-query-block-graph`.

Members

parent: int parent id

child: int child id

name: string name of the relation (examples are ‘file’ and ‘backing’)

perm: array of BlockPermission granted permissions for the parent operating on the child

shared-perm: array of BlockPermission permissions that can still be granted to other users of the child while it is still attached to this parent

Since

4.0

XDbgBlockGraph (Object)

Block Graph - list of nodes and list of edges.

Members

nodes: array of **XDbgBlockGraphNode** Not documented

edges: array of **XDbgBlockGraphEdge** Not documented

Since

4.0

x-debug-query-block-graph (Command)

Get the block graph.

Since

4.0

drive-mirror (Command)

Start mirroring a block device's writes to a new destination. `target` specifies the target of the new image. If the file exists, or if it is a device, it will be used as the new destination for writes. If it does not exist, a new file will be created. `format` specifies the format of the mirror image, default is to probe if `mode='existing'`, else the format of the source.

Arguments

The members of **DriveMirror**

Returns

- nothing on success
- If `device` is not a valid block device, `GenericError`

Since

1.3

Example

```
-> { "execute": "drive-mirror",  
    "arguments": { "device": "ide-hd0",  
                  "target": "/some/place/my-image",  
                  "sync": "full",  
                  "format": "qcow2" } }  
<- { "return": {} }
```

DriveMirror (Object)

A set of parameters describing drive mirror setup.

Members

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device name or node-name of a root node whose writes should be mirrored.

target: string the target of the new image. If the file exists, or if it is a device, the existing file/device will be used as the new destination. If it does not exist, a new file will be created.

format: string (optional) the format of the new destination, default is to probe if mode is ‘existing’, else the format of the source

node-name: string (optional) the new block driver state node name in the graph (Since 2.1)

replaces: string (optional) with sync=full graph node name to be replaced by the new image when a whole image copy is done. This can be used to repair broken Quorum files. By default, device is replaced, although implicitly created filters on it are kept. (Since 2.1)

mode: NewImageMode (optional) whether and how QEMU should create a new image, default is ‘absolute-paths’.

speed: int (optional) the maximum speed, in bytes per second

sync: MirrorSyncMode what parts of the disk image should be copied to the destination (all the disk, only the sectors allocated in the topmost image, or only new I/O).

granularity: int (optional) granularity of the dirty bitmap, default is 64K if the image format doesn’t have clusters, 4K if the clusters are smaller than that, else the cluster size. Must be a power of 2 between 512 and 64M (since 1.4).

buf-size: int (optional) maximum amount of data in flight from source to target (since 1.4).

on-source-error: BlockdevOnError (optional) the action to take on an error on the source, default ‘report’. ‘stop’ and ‘enospc’ can only be used if the block device supports io-status (see BlockInfo).

on-target-error: BlockdevOnError (optional) the action to take on an error on the target, default ‘report’ (no limitations, since this applies to a different block device than device).

unmap: boolean (optional) Whether to try to unmap target sectors where source has only zero. If true, and target unallocated sectors will read as zero, target image sectors will be unmapped; otherwise, zeroes will be written. Both will result in identical contents. Default is true. (Since 2.4)

copy-mode: MirrorCopyMode (optional) when to copy data to the destination; defaults to ‘background’ (Since: 3.0)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for block-job-finalize before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Since

1.3

BlockDirtyBitmap (Object)

Members

node: string name of device/node which the bitmap is tracking

name: string name of the dirty bitmap

Since

2.4

BlockDirtyBitmapAdd (Object)

Members

node: string name of device/node which the bitmap is tracking

name: string name of the dirty bitmap (must be less than 1024 bytes)

granularity: int (optional) the bitmap granularity, default is 64k for block-dirty-bitmap-add

persistent: boolean (optional) the bitmap is persistent, i.e. it will be saved to the corresponding block device image file on its close. For now only Qcow2 disks support persistent bitmaps. Default is false for block-dirty-bitmap-add. (Since: 2.10)

disabled: boolean (optional) the bitmap is created in the disabled state, which means that it will not track drive changes. The bitmap may be enabled with block-dirty-bitmap-enable. Default is false. (Since: 4.0)

Since

2.4

BlockDirtyBitmapMergeSource (Alternate)

Members

local: string name of the bitmap, attached to the same node as target bitmap.

external: BlockDirtyBitmap bitmap with specified node

Since

4.1

BlockDirtyBitmapMerge (Object)

Members

node: **string** name of device/node which the `target` bitmap is tracking

target: **string** name of the destination dirty bitmap

bitmaps: **array of BlockDirtyBitmapMergeSource** name(s) of the source dirty bitmap(s) at `node` and/or fully specified `BlockDirtyBitmap` elements. The latter are supported since 4.1.

Since

4.0

block-dirty-bitmap-add (Command)

Create a dirty bitmap with a name on the node, and start tracking the writes.

Returns

- nothing on success
- If `node` is not a valid block device or node, `DeviceNotFound`
- If `name` is already taken, `GenericError` with an explanation

Since

2.4

Example

```
-> { "execute": "block-dirty-bitmap-add",  
    "arguments": { "node": "drive0", "name": "bitmap0" } }  
<- { "return": {} }
```

block-dirty-bitmap-remove (Command)

Stop write tracking and remove the dirty bitmap that was created with `block-dirty-bitmap-add`. If the bitmap is persistent, remove it from its storage too.

Returns

- nothing on success
- If `node` is not a valid block device or node, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation
- if `name` is frozen by an operation, `GenericError`

Since

2.4

Example

```
-> { "execute": "block-dirty-bitmap-remove",
      "arguments": { "node": "drive0", "name": "bitmap0" } }
<- { "return": {} }
```

block-dirty-bitmap-clear (Command)

Clear (reset) a dirty bitmap on the device, so that an incremental backup from this point in time forward will only backup clusters modified after this clear operation.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation

Since

2.4

Example

```
-> { "execute": "block-dirty-bitmap-clear",
      "arguments": { "node": "drive0", "name": "bitmap0" } }
<- { "return": {} }
```

block-dirty-bitmap-enable (Command)

Enables a dirty bitmap so that it will begin tracking disk changes.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation

Since

4.0

Example

```
-> { "execute": "block-dirty-bitmap-enable",  
      "arguments": { "node": "drive0", "name": "bitmap0" } }  
<- { "return": {} }
```

block-dirty-bitmap-disable (Command)

Disables a dirty bitmap so that it will stop tracking disk changes.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation

Since

4.0

Example

```
-> { "execute": "block-dirty-bitmap-disable",  
      "arguments": { "node": "drive0", "name": "bitmap0" } }  
<- { "return": {} }
```

block-dirty-bitmap-merge (Command)

Merge dirty bitmaps listed in `bitmaps` to the `target` dirty bitmap. Dirty bitmaps in `bitmaps` will be unchanged, except if it also appears as the `target` bitmap. Any bits already set in `target` will still be set after the merge, i.e., this operation does not clear the `target`. On error, `target` is unchanged.

The resulting bitmap will count as dirty any clusters that were dirty in any of the source bitmaps. This can be used to achieve backup checkpoints, or in simpler usages, to copy bitmaps.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If any bitmap in `bitmaps` or `target` is not found, `GenericError`
- If any of the bitmaps have different sizes or granularities, `GenericError`

Since

4.0

Example

```
-> { "execute": "block-dirty-bitmap-merge",
      "arguments": { "node": "drive0", "target": "bitmap0",
                     "bitmaps": ["bitmap1"] } }
<- { "return": {} }
```

BlockDirtyBitmapSha256 (Object)

SHA256 hash of dirty bitmap data

Members

sha256: `string` ASCII representation of SHA256 bitmap hash

Since

2.10

x-debug-block-dirty-bitmap-sha256 (Command)

Get bitmap SHA256.

Returns

- `BlockDirtyBitmapSha256` on success
- If `node` is not a valid block device, `DeviceNotFound`
- If `name` is not found or if hashing has failed, `GenericError` with an explanation

Since

2.10

blockdev-mirror (Command)

Start mirroring a block device's writes to a new destination.

Arguments

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string The device name or node-name of a root node whose writes should be mirrored.

target: string the id or node-name of the block device to mirror to. This mustn't be attached to guest.

replaces: string (optional) with sync=full graph node name to be replaced by the new image when a whole image copy is done. This can be used to repair broken Quorum files. By default, `device` is replaced, although implicitly created filters on it are kept.

speed: int (optional) the maximum speed, in bytes per second

sync: MirrorSyncMode what parts of the disk image should be copied to the destination (all the disk, only the sectors allocated in the topmost image, or only new I/O).

granularity: int (optional) granularity of the dirty bitmap, default is 64K if the image format doesn't have clusters, 4K if the clusters are smaller than that, else the cluster size. Must be a power of 2 between 512 and 64M

buf-size: int (optional) maximum amount of data in flight from source to target

on-source-error: BlockdevOnError (optional) the action to take on an error on the source, default 'report'. 'stop' and 'enospc' can only be used if the block device supports io-status (see BlockInfo).

on-target-error: BlockdevOnError (optional) the action to take on an error on the target, default 'report' (no limitations, since this applies to a different block device than `device`).

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the mirror job inserts into the graph above `device`. If this option is not given, a node name is autogenerated. (Since: 2.9)

copy-mode: MirrorCopyMode (optional) when to copy data to the destination; defaults to 'background' (Since: 3.0)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Returns

nothing on success.

Since

2.6

Example

```
-> { "execute": "blockdev-mirror",
    "arguments": { "device": "ide-hd0",
                  "target": "target0",
                  "sync": "full" } }
<- { "return": {} }
```

BlockIOThrottle (Object)

A set of parameters describing block throttling.

Members

device: string (optional) Block device name

id: string (optional) The name or QOM path of the guest device (since: 2.8)

bps: int total throughput limit in bytes per second

bps_rd: int read throughput limit in bytes per second

bps_wr: int write throughput limit in bytes per second

iops: int total I/O operations per second

iops_rd: int read I/O operations per second

iops_wr: int write I/O operations per second

bps_max: int (optional) total throughput limit during bursts, in bytes (Since 1.7)

bps_rd_max: int (optional) read throughput limit during bursts, in bytes (Since 1.7)

bps_wr_max: int (optional) write throughput limit during bursts, in bytes (Since 1.7)

iops_max: int (optional) total I/O operations per second during bursts, in bytes (Since 1.7)

iops_rd_max: int (optional) read I/O operations per second during bursts, in bytes (Since 1.7)

iops_wr_max: int (optional) write I/O operations per second during bursts, in bytes (Since 1.7)

bps_max_length: int (optional) maximum length of the `bps_max` burst period, in seconds. It must only be set if `bps_max` is set as well. Defaults to 1. (Since 2.6)

bps_rd_max_length: int (optional) maximum length of the `bps_rd_max` burst period, in seconds. It must only be set if `bps_rd_max` is set as well. Defaults to 1. (Since 2.6)

bps_wr_max_length: int (optional) maximum length of the `bps_wr_max` burst period, in seconds. It must only be set if `bps_wr_max` is set as well. Defaults to 1. (Since 2.6)

iops_max_length: int (optional) maximum length of the `iops` burst period, in seconds. It must only be set if `iops_max` is set as well. Defaults to 1. (Since 2.6)

iops_rd_max_length: int (optional) maximum length of the `iops_rd_max` burst period, in seconds. It must only be set if `iops_rd_max` is set as well. Defaults to 1. (Since 2.6)

iops_wr_max_length: int (optional) maximum length of the `iops_wr_max` burst period, in seconds. It must only be set if `iops_wr_max` is set as well. Defaults to 1. (Since 2.6)

iops_size: int (optional) an I/O size in bytes (Since 1.7)

group: **string** (optional) throttle group name (Since 2.4)

Features

deprecated Member `device` is deprecated. Use `id` instead.

Since

1.1

ThrottleLimits (Object)

Limit parameters for throttling. Since some limit combinations are illegal, limits should always be set in one transaction. All fields are optional. When setting limits, if a field is missing the current value is not changed.

Members

iops-total: **int** (optional) limit total I/O operations per second

iops-total-max: **int** (optional) I/O operations burst

iops-total-max-length: **int** (optional) length of the iops-total-max burst period, in seconds It must only be set if iops-total-max is set as well.

iops-read: **int** (optional) limit read operations per second

iops-read-max: **int** (optional) I/O operations read burst

iops-read-max-length: **int** (optional) length of the iops-read-max burst period, in seconds It must only be set if iops-read-max is set as well.

iops-write: **int** (optional) limit write operations per second

iops-write-max: **int** (optional) I/O operations write burst

iops-write-max-length: **int** (optional) length of the iops-write-max burst period, in seconds It must only be set if iops-write-max is set as well.

bps-total: **int** (optional) limit total bytes per second

bps-total-max: **int** (optional) total bytes burst

bps-total-max-length: **int** (optional) length of the bps-total-max burst period, in seconds. It must only be set if bps-total-max is set as well.

bps-read: **int** (optional) limit read bytes per second

bps-read-max: **int** (optional) total bytes read burst

bps-read-max-length: **int** (optional) length of the bps-read-max burst period, in seconds It must only be set if bps-read-max is set as well.

bps-write: **int** (optional) limit write bytes per second

bps-write-max: **int** (optional) total bytes write burst

bps-write-max-length: **int** (optional) length of the bps-write-max burst period, in seconds It must only be set if bps-write-max is set as well.

iops-size: int (optional) when limiting by iops max size of an I/O in bytes

Since

2.11

block-stream (Command)

Copy data from a backing file into a block device.

The block streaming operation is performed in the background until the entire backing file has been copied. This command returns immediately once streaming has started. The status of ongoing block streaming operations can be checked with `query-block-jobs`. The operation can be stopped before it has completed using the `block-job-cancel` command.

The node that receives the data is called the top image, can be located in any part of the chain (but always above the base image; see below) and can be specified using its device or node name. Earlier qemu versions only allowed ‘device’ to name the top level node; presence of the ‘base-node’ parameter during introspection can be used as a witness of the enhanced semantics of ‘device’.

If a base file is specified then sectors are not copied from that base file and its backing chain. This can be used to stream a subset of the backing file chain instead of flattening the entire image. When streaming completes the image file will have the base file as its backing file, unless that node was changed while the job was running. In that case, base’s parent’s backing (or filtered, whichever exists) child (i.e., base at the beginning of the job) will be the new backing file.

On successful completion the image file is updated to drop the backing file and the `BLOCK_JOB_COMPLETED` event is emitted.

In case `device` is a filter node, `block-stream` modifies the first non-filter overlay node below it to point to the new backing node instead of modifying `device` itself.

Arguments

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device or node name of the top image

base: string (optional) the common backing file name. It cannot be set if `base-node` or `bottom` is also set.

base-node: string (optional) the node name of the backing file. It cannot be set if `base` or `bottom` is also set. (Since 2.8)

bottom: string (optional) the last node in the chain that should be streamed into top. It cannot be set if `base` or `base-node` is also set. It cannot be filter node. (Since 6.0)

backing-file: string (optional) The backing file string to write into the top image. This filename is not validated.

If a pathname string is such that it cannot be resolved by QEMU, that means that subsequent QMP or HMP commands must use node-names for the image in question, as filename lookup methods will fail.

If not specified, QEMU will automatically determine the backing file string to use, or error out if there is no obvious choice. Care should be taken when specifying the string, to specify a valid filename or protocol. (Since 2.1)

speed: int (optional) the maximum speed, in bytes per second

on-error: BlockdevOnError (optional) the action to take on an error (default report). ‘stop’ and ‘enospc’ can only be used if the block device supports io-status (see BlockInfo). Since 1.3.

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the stream job inserts into the graph above `device`. If this option is not given, a node name is autogenerated. (Since: 6.0)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Returns

- Nothing on success.
- If device does not exist, DeviceNotFound.

Since

1.1

Example

```
-> { "execute": "block-stream",  
    "arguments": { "device": "virtio0",  
                  "base": "/tmp/master.qcow2" } }  
<- { "return": {} }
```

block-job-set-speed (Command)

Set maximum speed for a background block operation.

This command can only be issued when there is an active block job.

Throttling can be disabled by setting the speed to 0.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

speed: int the maximum speed, in bytes per second, or 0 for unlimited. Defaults to 0.

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.1

`block-job-cancel` (Command)

Stop an active background block operation.

This command returns immediately after marking the active background block operation for cancellation. It is an error to call this command if no operation is in progress.

The operation will cancel as soon as possible and then emit the `BLOCK_JOB_CANCELLED` event. Before that happens the job is still visible when enumerated using `query-block-jobs`.

Note that if you issue ‘`block-job-cancel`’ after ‘`drive-mirror`’ has indicated (via the event `BLOCK_JOB_READY`) that the source and destination are synchronized, then the event triggered by this command changes to `BLOCK_JOB_COMPLETED`, to indicate that the mirroring has ended and the destination now has a point-in-time copy tied to the time of the cancellation.

For streaming, the image file retains its backing file unless the streaming operation happens to complete just as it is being cancelled. A new streaming operation can be started at a later time to finish copying all data from the backing file.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

force: boolean (optional) If true, and the job has already emitted the event `BLOCK_JOB_READY`, abandon the job immediately (even if it is paused) instead of waiting for the destination to complete its final synchronization (since 1.3)

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.1

block-job-pause (Command)

Pause an active background block operation.

This command returns immediately after marking the active background block operation for pausing. It is an error to call this command if no operation is in progress or if the job is already paused.

The operation will pause as soon as possible. No event is emitted when the operation is actually paused. Cancelling a paused job automatically resumes it.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.3

block-job-resume (Command)

Resume an active background block operation.

This command returns immediately after resuming a paused background block operation. It is an error to call this command if no operation is in progress or if the job is not paused.

This command also clears the error status of the job.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.3

block-job-complete (Command)

Manually trigger completion of an active background block operation. This is supported for drive mirroring, where it also switches the device to write to the target path only. The ability to complete is signaled with a BLOCK_JOB_READY event.

This command completes an active background block operation synchronously. The ordering of this command's return with the BLOCK_JOB_COMPLETED event is not defined. Note that if an I/O error occurs during the processing of this command: 1) the command itself will fail; 2) the error will be processed according to the `error`/`werror` arguments that were specified when starting the operation.

A cancelled or paused job cannot be completed.

Arguments

device: **string** The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

Returns

- Nothing on success
- If no background operation is active on this device, `DeviceNotActive`

Since

1.3

block-job-dismiss (Command)

For jobs that have already concluded, remove them from the block-job-query list. This command only needs to be run for jobs which were started with QEMU 2.12+ job lifetime management semantics.

This command will refuse to operate on any job that has not yet reached its terminal state, `JOB_STATUS_CONCLUDED`. For jobs that make use of the `BLOCK_JOB_READY` event, `block-job-cancel` or `block-job-complete` will still need to be used as appropriate.

Arguments

id: **string** The job identifier.

Returns

Nothing on success

Since

2.12

block-job-finalize (Command)

Once a job that has `manual=true` reaches the pending state, it can be instructed to finalize any graph changes and do any necessary cleanup via this command. For jobs in a transaction, instructing one job to finalize will force ALL jobs in the transaction to finalize, so it is only necessary to instruct a single member job to finalize.

Arguments

id: `string` The job identifier.

Returns

Nothing on success

Since

2.12

BlockdevDiscardOptions (Enum)

Determines how to handle discard requests.

Values

ignore Ignore the request

unmap Forward as an unmap request

Since

2.9

BlockdevDetectZeroesOptions (Enum)

Describes the operation mode for the automatic conversion of plain zero writes by the OS to driver specific optimized zero write commands.

Values

off Disabled (default)

on Enabled

unmap Enabled and even try to unmap blocks if possible. This requires also that `BlockdevDiscardOptions` is set to unmap for this device.

Since

2.1

BlockdevAioOptions (Enum)

Selects the AIO backend to handle I/O requests

Values

threads Use qemu's thread pool

native Use native AIO backend (only Linux and Windows)

io_uring (If: `defined(CONFIG_LINUX_IO_URING)`) Use linux io_uring (since 5.0)

Since

2.9

BlockdevCacheOptions (Object)

Includes cache-related options for block devices

Members

direct: boolean (optional) enables use of O_DIRECT (bypass the host page cache; default: false)

no-flush: boolean (optional) ignore any flush requests for the device (default: false)

Since

2.9

BlockdevDriver (Enum)

Drivers that are supported in block device operations.

Values

throttle Since 2.11

nvme Since 2.12

copy-on-read Since 3.0

blklogwrites Since 3.0

blkreplay Since 4.2

compress Since 5.0
blkdebug Not documented
blkverify Not documented
bochs Not documented
cloop Not documented
dmg Not documented
file Not documented
ftp Not documented
ftps Not documented
gluster Not documented
host_cdrom Not documented
host_device Not documented
http Not documented
https Not documented
iscsi Not documented
luks Not documented
nbd Not documented
nfs Not documented
null-aio Not documented
null-co Not documented
parallels Not documented
preallocate Not documented
qcow Not documented
qcow2 Not documented
qed Not documented
quorum Not documented
raw Not documented
rbd Not documented
replication (If: defined(CONFIG_REPLICATION)) Not documented
sheepdog Not documented
ssh Not documented
vdi Not documented
vhdx Not documented
vmdk Not documented
vpc Not documented
vvfat Not documented

Since

2.9

BlockdevOptionsFile (Object)

Driver specific block device options for the file backend.

Members

filename: string path to the image file

pr-manager: string (optional) the id for the object that will handle persistent reservations for this device (default: none, forward the commands via SG_IO; since 2.11)

aio: BlockdevAioOptions (optional) AIO backend (default: threads) (since: 2.8)

locking: OnOffAuto (optional) whether to enable file locking. If set to 'auto', only enable when Open File Descriptor (OFD) locking API is available (default: auto, since 2.10)

drop-cache: boolean (optional) (If: defined(CONFIG_LINUX)) invalidate page cache during live migration. This prevents stale data on the migration destination with cache.direct=off. Currently only supported on Linux hosts. (default: on, since: 4.0)

x-check-cache-dropped: boolean (optional) whether to check that page cache was dropped on live migration. May cause noticeable delays if the image file is large, do not use in production. (default: off) (since: 3.0)

Features

dynamic-auto-read-only If present, enabled auto-read-only means that the driver will open the image read-only at first, dynamically reopen the image file read-write when the first writer is attached to the node and reopen read-only when the last writer is detached. This allows giving QEMU write permissions only on demand when an operation actually needs write access.

Since

2.9

BlockdevOptionsNull (Object)

Driver specific block device options for the null backend.

Members

size: int (optional) size of the device in bytes.

latency-ns: int (optional) emulated latency (in nanoseconds) in processing requests. Default to zero which completes requests immediately. (Since 2.4)

read-zeroes: boolean (optional) if true, reads from the device produce zeroes; if false, the buffer is left unchanged. (default: false; since: 4.1)

Since

2.9

BlockdevOptionsNVMe (Object)

Driver specific block device options for the NVMe backend.

Members

device: string PCI controller address of the NVMe device in format hhhh:bb:ss.f (host:bus:slot.function)

namespace: int namespace number of the device, starting from 1.

Note that the PCI `device` must have been unbound from any host kernel driver before instructing QEMU to add the `blockdev`.

Since

2.12

BlockdevOptionsVVFAT (Object)

Driver specific block device options for the vvfat protocol.

Members

dir: string directory to be exported as FAT image

fat-type: int (optional) FAT type: 12, 16 or 32

floppy: boolean (optional) whether to export a floppy image (true) or partitioned hard disk (false; default)

label: string (optional) set the volume label, limited to 11 bytes. FAT16 and FAT32 traditionally have some restrictions on labels, which are ignored by most operating systems. Defaults to “QEMU VVFAT”. (since 2.4)

rw: boolean (optional) whether to allow write operations (default: false)

Since

2.9

BlockdevOptionsGenericFormat (Object)

Driver specific block device options for image format that have no option besides their data source.

Members

file: BlockdevRef reference to or definition of the data source block device

Since

2.9

BlockdevOptionsLUKS (Object)

Driver specific block device options for LUKS.

Members

key-secret: string (optional) the ID of a QCryptoSecret object providing the decryption key (since 2.6). Mandatory except when doing a metadata-only probe of the image.

The members of BlockdevOptionsGenericFormat

Since

2.9

BlockdevOptionsGenericCOWFormat (Object)

Driver specific block device options for image format that have no option besides their data source and an optional backing file.

Members

backing: BlockdevRefOrNull (optional) reference to or definition of the backing file block device, null disables the backing file entirely. Defaults to the backing file stored the image file.

The members of BlockdevOptionsGenericFormat

Since

2.9

Qcow2OverlapCheckMode (Enum)

General overlap check modes.

Values

none Do not perform any checks

constant Perform only checks which can be done in constant time and without reading anything from disk

cached Perform only checks which can be done without reading anything from disk

all Perform all available overlap checks

Since

2.9

Qcow2OverlapCheckFlags (Object)

Structure of flags for each metadata structure. Setting a field to ‘true’ makes qemu guard that structure against unintended overwriting. The default value is chosen according to the template given.

Members

template: Qcow2OverlapCheckMode (optional) Specifies a template mode which can be adjusted using the other flags, defaults to ‘cached’

bitmap-directory: boolean (optional) since 3.0

main-header: boolean (optional) Not documented

active-l1: boolean (optional) Not documented

active-l2: boolean (optional) Not documented

refcount-table: boolean (optional) Not documented

refcount-block: boolean (optional) Not documented

snapshot-table: boolean (optional) Not documented

inactive-l1: boolean (optional) Not documented

inactive-l2: boolean (optional) Not documented

Since

2.9

Qcow2OverlapChecks (Alternate)

Specifies which metadata structures should be guarded against unintended overwriting.

Members

flags: `Qcow2OverlapCheckFlags` set of flags for separate specification of each metadata structure type

mode: `Qcow2OverlapCheckMode` named mode which chooses a specific set of flags

Since

2.9

`BlockdevQcowEncryptionFormat` (Enum)

Values

aes AES-CBC with plain64 initialization vectors

Since

2.10

`BlockdevQcowEncryption` (Object)

Members

format: `BlockdevQcowEncryptionFormat` Not documented

The members of `QCryptoBlockOptionsQCow` when **format** is "aes"

Since

2.10

`BlockdevOptionsQcow` (Object)

Driver specific block device options for qcow.

Members

encrypt: `BlockdevQcowEncryption` (optional) Image decryption options. Mandatory for encrypted images, except when doing a metadata-only probe of the image.

The members of `BlockdevOptionsGenericCOWFormat`

Since

2.10

BlockdevQcow2EncryptionFormat (Enum)

Values

aes AES-CBC with plain64 initialization vectors

luks Not documented

Since

2.10

BlockdevQcow2Encryption (Object)

Members

format: BlockdevQcow2EncryptionFormat Not documented

The members of QCryptoBlockOptionsQcow when format is "aes"

The members of QCryptoBlockOptionsLUKS when format is "luks"

Since

2.10

BlockdevOptionsPreallocate (Object)

Filter driver intended to be inserted between format and protocol node and do preallocation in protocol node on write.

Members

prealloc-align: int (optional) on preallocation, align file length to this number, default 1048576 (1M)

prealloc-size: int (optional) how much to preallocate, default 134217728 (128M)

The members of BlockdevOptionsGenericFormat

Since

6.0

BlockdevOptionsQcow2 (Object)

Driver specific block device options for qcow2.

Members

- lazy-refcounts: boolean (optional)** whether to enable the lazy refcounts feature (default is taken from the image file)
- pass-discard-request: boolean (optional)** whether discard requests to the qcow2 device should be forwarded to the data source
- pass-discard-snapshot: boolean (optional)** whether discard requests for the data source should be issued when a snapshot operation (e.g. deleting a snapshot) frees clusters in the qcow2 file
- pass-discard-other: boolean (optional)** whether discard requests for the data source should be issued on other occasions where a cluster gets freed
- overlap-check: Qcow2OverlapChecks (optional)** which overlap checks to perform for writes to the image, defaults to ‘cached’ (since 2.2)
- cache-size: int (optional)** the maximum total size of the L2 table and refcount block caches in bytes (since 2.2)
- l2-cache-size: int (optional)** the maximum size of the L2 table cache in bytes (since 2.2)
- l2-cache-entry-size: int (optional)** the size of each entry in the L2 cache in bytes. It must be a power of two between 512 and the cluster size. The default value is the cluster size (since 2.12)
- refcount-cache-size: int (optional)** the maximum size of the refcount block cache in bytes (since 2.2)
- cache-clean-interval: int (optional)** clean unused entries in the L2 and refcount caches. The interval is in seconds. The default value is 600 on supporting platforms, and 0 on other platforms. 0 disables this feature. (since 2.5)
- encrypt: BlockdevQcow2Encryption (optional)** Image decryption options. Mandatory for encrypted images, except when doing a metadata-only probe of the image. (since 2.10)
- data-file: BlockdevRef (optional)** reference to or definition of the external data file. This may only be specified for images that require an external data file. If it is not specified for such an image, the data file name is loaded from the image file. (since 4.0)

The members of `BlockdevOptionsGenericCOWFormat`

Since

2.9

SshHostKeyCheckMode (Enum)

Values

- none** Don’t check the host key at all
- hash** Compare the host key with a given hash
- known_hosts** Check the host key against the known_hosts file

Since

2.12

SshHostKeyCheckHashType (Enum)

Values

md5 The given hash is an md5 hash

sha1 The given hash is an sha1 hash

Since

2.12

SshHostKeyHash (Object)

Members

type: SshHostKeyCheckHashType The hash algorithm used for the hash

hash: string The expected hash value

Since

2.12

SshHostKeyCheck (Object)

Members

mode: SshHostKeyCheckMode Not documented

The members of SshHostKeyHash when mode is "hash"

Since

2.12

BlockdevOptionsSsh (Object)

Members

server: InetSocketAddress host address

path: string path to the image on the host

user: string (optional) user as which to connect, defaults to current local user name

host-key-check: SshHostKeyCheck (optional) Defines how and what to check the host key against (default: known_hosts)

Since

2.9

BlkdebugEvent (Enum)

Trigger events supported by blkdebug.

Values

l1_shrink_write_table write zeros to the l1 table to shrink image. (since 2.11)

l1_shrink_free_l2_clusters discard the l2 tables. (since 2.11)

cor_write a write due to copy-on-read (since 2.11)

cluster_alloc_space an allocation of file space for a cluster (since 4.1)

none triggers once at creation of the blkdebug node (since 4.1)

l1_update Not documented

l1_grow_alloc_table Not documented

l1_grow_write_table Not documented

l1_grow_activate_table Not documented

l2_load Not documented

l2_update Not documented

l2_update_compressed Not documented

l2_alloc_cow_read Not documented

l2_alloc_write Not documented

read_aio Not documented

read_backing_aio Not documented

read_compressed Not documented

write_aio Not documented

write_compressed Not documented

vmstate_load Not documented

vmstate_save Not documented

cow_read Not documented

cow_write Not documented

reftable_load Not documented

reftable_grow Not documented

reftable_update Not documented

refblock_load Not documented

refblock_update Not documented

refblock_update_part Not documented
refblock_alloc Not documented
refblock_alloc_hookup Not documented
refblock_alloc_write Not documented
refblock_alloc_write_blocks Not documented
refblock_alloc_write_table Not documented
refblock_alloc_switch_table Not documented
cluster_alloc Not documented
cluster_alloc_bytes Not documented
cluster_free Not documented
flush_to_os Not documented
flush_to_disk Not documented
pwritev_rmw_head Not documented
pwritev_rmw_after_head Not documented
pwritev_rmw_tail Not documented
pwritev_rmw_after_tail Not documented
pwritev Not documented
pwritev_zero Not documented
pwritev_done Not documented
empty_image_prepare Not documented

Since

2.9

BlkdebugIOType (Enum)

Kinds of I/O that blkdebug can inject errors in.

Values

read .bdrv_co_preadv()
write .bdrv_co_pwritev()
write-zeroes .bdrv_co_pwrite_zeroes()
discard .bdrv_co_pdiscard()
flush .bdrv_co_flush_to_disk()
block-status .bdrv_co_block_status()

Since

4.1

BlkdebugInjectErrorOptions (Object)

Describes a single error injection for blkdebug.

Members

event: **BlkdebugEvent** trigger event

state: **int (optional)** the state identifier blkdebug needs to be in to actually trigger the event; defaults to “any”

iotype: **BlkdebugIOType (optional)** the type of I/O operations on which this error should be injected; defaults to “all read, write, write-zeroes, discard, and flush operations” (since: 4.1)

errno: **int (optional)** error identifier (errno) to be returned; defaults to EIO

sector: **int (optional)** specifies the sector index which has to be affected in order to actually trigger the event; defaults to “any sector”

once: **boolean (optional)** disables further events after this one has been triggered; defaults to false

immediately: **boolean (optional)** fail immediately; defaults to false

Since

2.9

BlkdebugSetStateOptions (Object)

Describes a single state-change event for blkdebug.

Members

event: **BlkdebugEvent** trigger event

state: **int (optional)** the current state identifier blkdebug needs to be in; defaults to “any”

new_state: **int** the state identifier blkdebug is supposed to assume if this event is triggered

Since

2.9

BlockdevOptionsBlkdebug (Object)

Driver specific block device options for blkdebug.

Members

image: BlockdevRef underlying raw block device (or image file)

config: string (optional) filename of the configuration file

align: int (optional) required alignment for requests in bytes, must be positive power of 2, or 0 for default

max-transfer: int (optional) maximum size for I/O transfers in bytes, must be positive multiple of `align` and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

opt-write-zero: int (optional) preferred alignment for write zero requests in bytes, must be positive multiple of `align` and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

max-write-zero: int (optional) maximum size for write zero requests in bytes, must be positive multiple of `align`, of `opt-write-zero`, and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

opt-discard: int (optional) preferred alignment for discard requests in bytes, must be positive multiple of `align` and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

max-discard: int (optional) maximum size for discard requests in bytes, must be positive multiple of `align`, of `opt-discard`, and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

inject-error: array of BlkdebugInjectErrorOptions (optional) array of error injection descriptions

set-state: array of BlkdebugSetStateOptions (optional) array of state-change descriptions

take-child-perms: array of BlockPermission (optional) Permissions to take on `image` in addition to what is necessary anyway (which depends on how the blkdebug node is used). Defaults to none. (since 5.0)

unshare-child-perms: array of BlockPermission (optional) Permissions not to share on `image` in addition to what cannot be shared anyway (which depends on how the blkdebug node is used). Defaults to none. (since 5.0)

Since

2.9

BlockdevOptionsBlklogwrites (Object)

Driver specific block device options for blklogwrites.

Members

file: BlockdevRef block device

log: BlockdevRef block device used to log writes to `file`

log-sector-size: int (optional) sector size used in logging writes to `file`, determines granularity of offsets and sizes of writes (default: 512)

log-append: boolean (optional) append to an existing log (default: false)

log-super-update-interval: int (optional) interval of write requests after which the log super block is updated to disk (default: 4096)

Since

3.0

BlockdevOptionsBlkverify (Object)

Driver specific block device options for blkverify.

Members

test: BlockdevRef block device to be tested

raw: BlockdevRef raw image used for verification

Since

2.9

BlockdevOptionsBlkreplay (Object)

Driver specific block device options for blkreplay.

Members

image: BlockdevRef disk image which should be controlled with blkreplay

Since

4.2

QuorumReadPattern (Enum)

An enumeration of quorum read patterns.

Values

quorum read all the children and do a quorum vote on reads

fifo read only from the first child that has not failed

Since

2.9

BlockdevOptionsQuorum (Object)

Driver specific block device options for Quorum

Members

blkverify: boolean (optional)

true if the driver must print content mismatch set to false by default

children: array of BlockdevRef the children block devices to use

vote-threshold: int the vote limit under which a read will fail

rewrite-corrupted: boolean (optional) rewrite corrupted data when quorum is reached (Since 2.1)

read-pattern: QuorumReadPattern (optional) choose read pattern and set to quorum by default (Since 2.2)

Since

2.9

BlockdevOptionsGluster (Object)

Driver specific block device options for Gluster

Members

volume: string name of gluster volume where VM image resides

path: string absolute path to image file in gluster volume

server: array of SocketAddress gluster servers description

debug: int (optional) libgfs2 log level (default '4' which is Error) (Since 2.8)

logfile: string (optional) libgfs2 log file (default /dev/stderr) (Since 2.8)

Since

2.9

IscsiTransport (Enum)

An enumeration of libiscsi transport types

Values

tcp Not documented

iser Not documented

Since

2.9

IscsiHeaderDigest (Enum)

An enumeration of header digests supported by libiscsi

Values

crc32c Not documented

none Not documented

crc32c-none Not documented

none-crc32c Not documented

Since

2.9

BlockdevOptionsIscsi (Object)

Members

transport: IscsiTransport The iscsi transport type

portal: string The address of the iscsi portal

target: string The target iqname

lun: int (optional) LUN to connect to. Defaults to 0.

user: string (optional) User name to log in with. If omitted, no CHAP authentication is performed.

password-secret: string (optional) The ID of a QCryptoSecret object providing the password for the login.
This option is required if `user` is specified.

initiator-name: string (optional) The iqname we want to identify to the target as. If this option is not specified, an initiator name is generated automatically.

header-digest: IscsiHeaderDigest (optional) The desired header digest. Defaults to none-crc32c.

timeout: int (optional) Timeout in seconds after which a request will timeout. 0 means no timeout and is the default.

Driver specific block device options for iscsi

Since

2.9

RbdAuthMode (Enum)

Values

cephx Not documented

none Not documented

Since

3.0

BlockdevOptionsRbd (Object)

Members

pool: string Ceph pool name.

namespace: string (optional) Rados namespace name in the Ceph pool. (Since 5.0)

image: string Image name in the Ceph pool.

conf: string (optional) path to Ceph configuration file. Values in the configuration file will be overridden by options specified via QAPI.

snapshot: string (optional) Ceph snapshot name.

user: string (optional) Ceph id name.

auth-client-required: array of RbdAuthMode (optional) Acceptable authentication modes. This maps to Ceph configuration option “auth_client_required”. (Since 3.0)

key-secret: string (optional) ID of a QCryptoSecret object providing a key for cephx authentication. This maps to Ceph configuration option “key”. (Since 3.0)

server: array of InetSocketAddressBase (optional) Monitor host address and port. This maps to the “mon_host” Ceph option.

Since

2.9

BlockdevOptionsSheepdog (Object)

Driver specific block device options for sheepdog

Members

vdi: string Virtual disk image name

server: SocketAddress The Sheepdog server to connect to

snap-id: int (optional) Snapshot ID

tag: string (optional) Snapshot tag name

Only one of `snap-id` and `tag` may be present.

Since

2.9

ReplicationMode (Enum)

An enumeration of replication modes.

Values

primary Primary mode, the vm's state will be sent to secondary QEMU.

secondary Secondary mode, receive the vm's state from primary QEMU.

Since

2.9

If

`defined(CONFIG_REPLICATION)`

BlockdevOptionsReplication (Object)

Driver specific block device options for replication

Members

mode: ReplicationMode the replication mode

top-id: string (optional) In secondary mode, node name or device ID of the root node who owns the replication node chain. Must not be given in primary mode.

The members of `BlockdevOptionsGenericFormat`

Since

2.9

If

`defined(CONFIG_REPLICATION)`

NFSTransport (Enum)

An enumeration of NFS transport types

Values

inet TCP transport

Since

2.9

NFSServer (Object)

Captures the address of the socket

Members

type: NFSTransport transport type used for NFS (only TCP supported)

host: string host address for NFS server

Since

2.9

BlockdevOptionsNfs (Object)

Driver specific block device option for NFS

Members

server: NFSServer host address

path: string path of the image on the host

user: int (optional) UID value to use when talking to the server (defaults to 65534 on Windows and `getuid()` on unix)

group: int (optional) GID value to use when talking to the server (defaults to 65534 on Windows and getgid() in unix)

tcp-syn-count: int (optional) number of SYNs during the session establishment (defaults to libnfs default)

readahead-size: int (optional) set the readahead size in bytes (defaults to libnfs default)

page-cache-size: int (optional) set the pagecache size in bytes (defaults to libnfs default)

debug: int (optional) set the NFS debug level (max 2) (defaults to libnfs default)

Since

2.9

BlockdevOptionsCurlBase (Object)

Driver specific block device options shared by all protocols supported by the curl backend.

Members

url: string URL of the image file

readahead: int (optional) Size of the read-ahead cache; must be a multiple of 512 (defaults to 256 kB)

timeout: int (optional) Timeout for connections, in seconds (defaults to 5)

username: string (optional) Username for authentication (defaults to none)

password-secret: string (optional) ID of a QCryptoSecret object providing a password for authentication (defaults to no password)

proxy-username: string (optional) Username for proxy authentication (defaults to none)

proxy-password-secret: string (optional) ID of a QCryptoSecret object providing a password for proxy authentication (defaults to no password)

Since

2.9

BlockdevOptionsCurlHttp (Object)

Driver specific block device options for HTTP connections over the curl backend. URLs must start with “[http://](#)”.

Members

cookie: string (optional) List of cookies to set; format is “name1=content1; name2=content2;” as explained by CURLOPT_COOKIE(3). Defaults to no cookies.

cookie-secret: string (optional) ID of a QCryptoSecret object providing the cookie data in a secure way. See `cookie` for the format. (since 2.10)

The members of `BlockdevOptionsCurlBase`

Since

2.9

BlockdevOptionsCurlHttps (Object)

Driver specific block device options for HTTPS connections over the curl backend. URLs must start with “[https://](#)”.

Members

cookie: string (optional) List of cookies to set; format is “name1=content1; name2=content2;” as explained by CURLOPT_COOKIE(3). Defaults to no cookies.

sslverify: boolean (optional) Whether to verify the SSL certificate’s validity (defaults to true)

cookie-secret: string (optional) ID of a QCryptoSecret object providing the cookie data in a secure way. See `cookie` for the format. (since 2.10)

The members of `BlockdevOptionsCurlBase`

Since

2.9

BlockdevOptionsCurlFtp (Object)

Driver specific block device options for FTP connections over the curl backend. URLs must start with “[ftp://](#)”.

Members

The members of `BlockdevOptionsCurlBase`

Since

2.9

BlockdevOptionsCurlFtps (Object)

Driver specific block device options for FTPS connections over the curl backend. URLs must start with “[ftps://](#)”.

Members

sslverify: boolean (optional) Whether to verify the SSL certificate’s validity (defaults to true)

The members of `BlockdevOptionsCurlBase`

Since

2.9

BlockdevOptionsNbd (Object)

Driver specific block device options for NBD.

Members

server: SocketAddress NBD server address

export: string (optional) export name

tls-creds: string (optional) TLS credentials ID

x-dirty-bitmap: string (optional) A metadata context name such as “qemu:dirty-bitmap:NAME” or “qemu:allocation-depth” to query in place of the traditional “base:allocation” block status (see NBD_OPT_LIST_META_CONTEXT in the NBD protocol; and yes, naming this option x-context would have made more sense) (since 3.0)

reconnect-delay: int (optional) On an unexpected disconnect, the nbd client tries to connect again until succeeding or encountering a serious error. During the first `reconnect-delay` seconds, all requests are paused and will be rerun on a successful reconnect. After that time, any delayed requests and all future requests before a successful reconnect will immediately fail. Default 0 (Since 4.2)

Since

2.9

BlockdevOptionsRaw (Object)

Driver specific block device options for the raw driver.

Members

offset: int (optional) position where the block device starts

size: int (optional) the assumed size of the device

The members of **BlockdevOptionsGenericFormat**

Since

2.9

BlockdevOptionsThrottle (Object)

Driver specific block device options for the throttle driver

Members

throttle-group: string the name of the throttle-group object to use. It must already exist.

file: BlockdevRef reference to or definition of the data source block device

Since

2.11

BlockdevOptionsCor (Object)

Driver specific block device options for the copy-on-read driver.

Members

bottom: string (optional) The name of a non-filter node (allocation-bearing layer) that limits the COR operations in the backing chain (inclusive), so that no data below this node will be copied by this filter. If option is absent, the limit is not applied, so that data from all backing layers may be copied.

The members of **BlockdevOptionsGenericFormat**

Since

6.0

BlockdevOptions (Object)

Options for creating a block device. Many options are available for all block devices, independent of the block driver:

Members

driver: BlockdevDriver block driver name

node-name: string (optional) the node name of the new node (Since 2.0). This option is required on the top level of blockdev-add. Valid node names start with an alphabetic character and may contain only alphanumeric characters, '-', '.', and '_'. Their maximum length is 31 characters.

discard: BlockdevDiscardOptions (optional) discard-related options (default: ignore)

cache: BlockdevCacheOptions (optional) cache-related options

read-only: boolean (optional) whether the block device should be read-only (default: false). Note that some block drivers support only read-only access, either generally or in certain configurations. In this case, the default value does not work and the option must be specified explicitly.

auto-read-only: boolean (optional) if true and `read-only` is false, QEMU may automatically decide not to open the image read-write as requested, but fall back to read-only instead (and switch between the modes later), e.g. depending on whether the image file is writable or whether a writing user is attached to the node (default: false, since 3.1)

detect-zeroes: `BlockdevDetectZeroesOptions` (optional) detect and optimize zero writes (Since 2.1) (default: off)

force-share: `boolean` (optional) force share all permission on added nodes. Requires read-only=true. (Since 2.10)

The members of `BlockdevOptionsBlkdebug` when driver is "blkdebug"

The members of `BlockdevOptionsBlklogwrites` when driver is "blklogwrites"

The members of `BlockdevOptionsBlkverify` when driver is "blkverify"

The members of `BlockdevOptionsBlkreplay` when driver is "blkreplay"

The members of `BlockdevOptionsGenericFormat` when driver is "bochs"

The members of `BlockdevOptionsGenericFormat` when driver is "cloop"

The members of `BlockdevOptionsGenericFormat` when driver is "compress"

The members of `BlockdevOptionsCor` when driver is "copy-on-read"

The members of `BlockdevOptionsGenericFormat` when driver is "dmg"

The members of `BlockdevOptionsFile` when driver is "file"

The members of `BlockdevOptionsCurlFtp` when driver is "ftp"

The members of `BlockdevOptionsCurlFtps` when driver is "ftps"

The members of `BlockdevOptionsGluster` when driver is "gluster"

The members of `BlockdevOptionsFile` when driver is "host_cdrom"

The members of `BlockdevOptionsFile` when driver is "host_device"

The members of `BlockdevOptionsCurlHttp` when driver is "http"

The members of `BlockdevOptionsCurlHttps` when driver is "https"

The members of `BlockdevOptionsIscsi` when driver is "iscsi"

The members of `BlockdevOptionsLUKS` when driver is "luks"

The members of `BlockdevOptionsNbd` when driver is "nbd"

The members of `BlockdevOptionsNfs` when driver is "nfs"

The members of `BlockdevOptionsNull` when driver is "null-aio"

The members of `BlockdevOptionsNull` when driver is "null-co"

The members of `BlockdevOptionsNVMe` when driver is "nvme"

The members of `BlockdevOptionsGenericFormat` when driver is "parallels"

The members of `BlockdevOptionsPreallocate` when driver is "preallocate"

The members of `BlockdevOptionsQcow2` when driver is "qcow2"

The members of `BlockdevOptionsQcow` when driver is "qcow"

The members of `BlockdevOptionsGenericCOWFormat` when driver is "qed"

The members of `BlockdevOptionsQuorum` when driver is "quorum"

The members of `BlockdevOptionsRaw` when driver is "raw"

The members of `BlockdevOptionsRbd` when driver is "rbd"

The members of `BlockdevOptionsReplication` when driver is "replication" (If: defined(`CONFIG_REPLICATION`))

The members of `BlockdevOptionsSheepdog` when driver is "sheepdog"

The members of `BlockdevOptionsSsh` when driver is "ssh"

The members of `BlockdevOptionsThrottle` when driver is "throttle"

The members of `BlockdevOptionsGenericFormat` when driver is "vdi"

The members of `BlockdevOptionsGenericFormat` when driver is "vhdx"

The members of `BlockdevOptionsGenericCOWFormat` when driver is "vmdk"

The members of `BlockdevOptionsGenericFormat` when driver is "vpc"

The members of `BlockdevOptionsVVFAT` when driver is "vvfat"

Remaining options are determined by the block driver.

Since

2.9

BlockdevRef (Alternate)

Reference to a block device.

Members

definition: `BlockdevOptions` defines a new block device inline

reference: `string` references the ID of an existing block device

Since

2.9

BlockdevRefOrNull (Alternate)

Reference to a block device.

Members

definition: `BlockdevOptions` defines a new block device inline

reference: `string` references the ID of an existing block device. An empty string means that no block device should be referenced. Deprecated; use null instead.

null: `null` No block device should be referenced (since 2.10)

Since

2.9

blockdev-add (Command)

Creates a new block device.

Arguments**The members of BlockdevOptions****Since**

2.9

Example

```

1.
-> { "execute": "blockdev-add",
    "arguments": {
        "driver": "qcow2",
        "node-name": "test1",
        "file": {
            "driver": "file",
            "filename": "test.qcow2"
        }
    }
}
<- { "return": {} }

2.
-> { "execute": "blockdev-add",
    "arguments": {
        "driver": "qcow2",
        "node-name": "node0",
        "discard": "unmap",
        "cache": {
            "direct": true
        },
        "file": {
            "driver": "file",
            "filename": "/tmp/test.qcow2"
        },
        "backing": {
            "driver": "raw",
            "file": {
                "driver": "file",
                "filename": "/dev/fdset/4"
            }
        }
    }
}

```

(continues on next page)

(continued from previous page)

```
    }  
  }  
  
<- { "return": {} }
```

x-blockdev-reopen (Command)

Reopens a block device using the given set of options. Any option not specified will be reset to its default value regardless of its previous status. If an option cannot be changed or a particular driver does not support reopening then the command will return an error.

The top-level `node-name` option (from `BlockdevOptions`) must be specified and is used to select the block device to be reopened. Other `node-name` options must be either omitted or set to the current name of the appropriate node. This command won't change any node name and any attempt to do it will result in an error.

In the case of options that refer to child nodes, the behavior of this command depends on the value:

- 1) A set of options (`BlockdevOptions`): the child is reopened with the specified set of options.
- 2) A reference to the current child: the child is reopened using its existing set of options.
- 3) A reference to a different node: the current child is replaced with the specified one.
- 4) `NULL`: the current child (if any) is detached.

Options (1) and (2) are supported in all cases, but at the moment only `backing` allows replacing or detaching an existing child.

Unlike with `blockdev-add`, the `backing` option must always be present unless the node being reopened does not have a backing file and its image does not have a default backing file name as part of its metadata.

Arguments

The members of `BlockdevOptions`

Since

4.0

blockdev-del (Command)

Deletes a block device that has been added using `blockdev-add`. The command will fail if the node is attached to a device or is otherwise being used.

Arguments

node-name: **string** Name of the graph node to delete.

Since

2.9

Example

```

-> { "execute": "blockdev-add",
    "arguments": {
        "driver": "qcow2",
        "node-name": "node0",
        "file": {
            "driver": "file",
            "filename": "test.qcow2"
        }
    }
}
<- { "return": {} }

-> { "execute": "blockdev-del",
    "arguments": { "node-name": "node0" }
}
<- { "return": {} }

```

BlockdevCreateOptionsFile (Object)

Driver specific image creation options for file.

Members

filename: string Filename for the new image file

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, falloc (if defined CONFIG_POSIX_FALLOCATE), full (if defined CONFIG_POSIX))

nocow: boolean (optional) Turn off copy-on-write (valid only on btrfs; default: off)

extent-size-hint: int (optional) Extent size hint to add to the image file; 0 for not adding an extent size hint (default: 1 MB, since 5.1)

Since

2.12

BlockdevCreateOptionsGluster (Object)

Driver specific image creation options for gluster.

Members

location: BlockdevOptionsGluster Where to store the new image file

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, falloc (if defined CONFIG_GLUSTERFS_FALLOCATE), full (if defined CONFIG_GLUSTERFS_ZEROFILL))

Since

2.12

BlockdevCreateOptionsLUKS (Object)

Driver specific image creation options for LUKS.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (since: 4.2) (default: off; allowed values: off, metadata, falloc, full)

The members of `QCryptoBlockCreateOptionsLUKS`

Since

2.12

BlockdevCreateOptionsNfs (Object)

Driver specific image creation options for NFS.

Members

location: BlockdevOptionsNfs Where to store the new image file

size: int Size of the virtual disk in bytes

Since

2.12

BlockdevCreateOptionsParallels (Object)

Driver specific image creation options for parallels.

Members

file: BlockdevRef Node to create the image format on
size: int Size of the virtual disk in bytes
cluster-size: int (optional) Cluster size in bytes (default: 1 MB)

Since

2.12

BlockdevCreateOptionsQcow (Object)

Driver specific image creation options for qcow.

Members

file: BlockdevRef Node to create the image format on
size: int Size of the virtual disk in bytes
backing-file: string (optional) File name of the backing file if a backing file should be used
encrypt: QCryptoBlockCreateOptions (optional) Encryption options if the image should be encrypted

Since

2.12

BlockdevQcow2Version (Enum)

Values

v2 The original QCOW2 format as introduced in qemu 0.10 (version 2)
v3 The extended QCOW2 format as introduced in qemu 1.1 (version 3)

Since

2.12

Qcow2CompressionType (Enum)

Compression type used in qcow2 image file

Values

zlib zlib compression, see <<http://zlib.net/>>

zstd (If: defined(CONFIG_ZSTD)) zstd compression, see <<http://github.com/facebook/zstd>>

Since

5.1

BlockdevCreateOptionsQcow2 (Object)

Driver specific image creation options for qcow2.

Members

file: BlockdevRef Node to create the image format on

data-file: BlockdevRef (optional) Node to use as an external data file in which all guest data is stored so that only metadata remains in the qcow2 file (since: 4.0)

data-file-raw: boolean (optional) True if the external data file must stay valid as a standalone (read-only) raw image without looking at qcow2 metadata (default: false; since: 4.0)

extended-l2: boolean (optional) True to make the image have extended L2 entries (default: false; since 5.2)

size: int Size of the virtual disk in bytes

version: BlockdevQcow2Version (optional) Compatibility level (default: v3)

backing-file: string (optional) File name of the backing file if a backing file should be used

backing-fmt: BlockdevDriver (optional) Name of the block driver to use for the backing file

encrypt: QCryptoBlockCreateOptions (optional) Encryption options if the image should be encrypted

cluster-size: int (optional) qcow2 cluster size in bytes (default: 65536)

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, falloc, full, metadata)

lazy-refcounts: boolean (optional) True if refcounts may be updated lazily (default: off)

refcount-bits: int (optional) Width of reference counts in bits (default: 16)

compression-type: Qcow2CompressionType (optional) The image cluster compression method (default: zlib, since 5.1)

Since

2.12

BlockdevCreateOptionsQed (Object)

Driver specific image creation options for qed.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

backing-file: string (optional) File name of the backing file if a backing file should be used

backing-fmt: BlockdevDriver (optional) Name of the block driver to use for the backing file

cluster-size: int (optional) Cluster size in bytes (default: 65536)

table-size: int (optional) L1/L2 table size (in clusters)

Since

2.12

BlockdevCreateOptionsRbd (Object)

Driver specific image creation options for rbd/Ceph.

Members

location: BlockdevOptionsRbd Where to store the new image file. This location cannot point to a snapshot.

size: int Size of the virtual disk in bytes

cluster-size: int (optional) RBD object size

Since

2.12

BlockdevVmdkSubformat (Enum)

Subformat options for VMDK images

Values

monolithicSparse Single file image with sparse cluster allocation

monolithicFlat Single flat data image and a descriptor file

twoGbMaxExtentSparse Data is split into 2GB (per virtual LBA) sparse extent files, in addition to a descriptor file

twoGbMaxExtentFlat Data is split into 2GB (per virtual LBA) flat extent files, in addition to a descriptor file

streamOptimized Single file image sparse cluster allocation, optimized for streaming over network.

Since

4.0

BlockdevVmdkAdapterType (Enum)

Adapter type info for VMDK images

Values

ide Not documented

buslogic Not documented

lsilogic Not documented

legacyESX Not documented

Since

4.0

BlockdevCreateOptionsVmdk (Object)

Driver specific image creation options for VMDK.

Members

file: BlockdevRef Where to store the new image file. This refers to the image file for monolithicSparse and streamOptimized format, or the descriptor file for other formats.

size: int Size of the virtual disk in bytes

extents: array of BlockdevRef (optional) Where to store the data extents. Required for monolithicFlat, twoGbMaxExtentSparse and twoGbMaxExtentFlat formats. For monolithicFlat, only one entry is required; for twoGbMaxExtent* formats, the number of entries required is calculated as $\text{extent_number} = \text{virtual_size} / 2\text{GB}$. Providing more extents than will be used is an error.

subformat: BlockdevVmdkSubformat (optional) The subformat of the VMDK image. Default: “monolithic-Sparse”.

backing-file: string (optional) The path of backing file. Default: no backing file is used.

adapter-type: BlockdevVmdkAdapterType (optional) The adapter type used to fill in the descriptor. Default: ide.

hwversion: string (optional) Hardware version. The meaningful options are “4” or “6”. Default: “4”.

zeroed-grain: boolean (optional) Whether to enable zeroed-grain feature for sparse subformats. Default: false.

Since

4.0

SheepdogRedundancyType (Enum)

Values

full Create a fully replicated vdi with x copies

erasure-coded Create an erasure coded vdi with x data strips and y parity strips

Since

2.12

SheepdogRedundancyFull (Object)

Members

copies: int Number of copies to use (between 1 and 31)

Since

2.12

SheepdogRedundancyErasureCoded (Object)

Members

data-strips: int Number of data strips to use (one of {2,4,8,16})

parity-strips: int Number of parity strips to use (between 1 and 15)

Since

2.12

SheepdogRedundancy (Object)

Members

type: SheepdogRedundancyType Not documented

The members of SheepdogRedundancyFull when type is "full"

The members of SheepdogRedundancyErasureCoded when type is "erasure-coded"

Since

2.12

BlockdevCreateOptionsSheepdog (Object)

Driver specific image creation options for Sheepdog.

Members

location: BlockdevOptionsSheepdog Where to store the new image file

size: int Size of the virtual disk in bytes

backing-file: string (optional) File name of a base image

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, full)

redundancy: SheepdogRedundancy (optional) Redundancy of the image

object-size: int (optional) Object size of the image

Since

2.12

BlockdevCreateOptionsSsh (Object)

Driver specific image creation options for SSH.

Members

location: BlockdevOptionsSsh Where to store the new image file

size: int Size of the virtual disk in bytes

Since

2.12

BlockdevCreateOptionsVdi (Object)

Driver specific image creation options for VDI.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, metadata)

Since

2.12

BlockdevVhdxSubformat (Enum)

Values

dynamic Growing image file

fixed Preallocated fixed-size image file

Since

2.12

BlockdevCreateOptionsVhdx (Object)

Driver specific image creation options for vhdx.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

log-size: int (optional) Log size in bytes, must be a multiple of 1 MB (default: 1 MB)

block-size: int (optional) Block size in bytes, must be a multiple of 1 MB and not larger than 256 MB (default: automatically choose a block size depending on the image size)

subformat: BlockdevVhdxSubformat (optional) vhdx subformat (default: dynamic)

block-state-zero: boolean (optional) Force use of payload blocks of type 'ZERO'. Non-standard, but default. Do not set to 'off' when using 'qemu-img convert' with subformat=dynamic.

Since

2.12

BlockdevVpcSubformat (Enum)

Values

dynamic Growing image file

fixed Preallocated fixed-size image file

Since

2.12

BlockdevCreateOptionsVpc (Object)

Driver specific image creation options for vpc (VHD).

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

subformat: BlockdevVpcSubformat (optional) vhd subformat (default: dynamic)

force-size: boolean (optional) Force use of the exact byte size instead of rounding to the next size that can be represented in CHS geometry (default: false)

Since

2.12

BlockdevCreateOptions (Object)

Options for creating an image format on a given node.

Members

driver: BlockdevDriver block driver to create the image format

The members of **BlockdevCreateOptionsFile** when **driver** is "file"

The members of **BlockdevCreateOptionsGluster** when **driver** is "gluster"

The members of **BlockdevCreateOptionsLUKS** when **driver** is "luks"

The members of **BlockdevCreateOptionsNfs** when **driver** is "nfs"

The members of **BlockdevCreateOptionsParallels** when **driver** is "parallels"

The members of **BlockdevCreateOptionsQcow** when **driver** is "qcow"

The members of **BlockdevCreateOptionsQcow2** when **driver** is "qcow2"

The members of `BlockdevCreateOptionsQed` when driver is "qed"

The members of `BlockdevCreateOptionsRbd` when driver is "rbd"

The members of `BlockdevCreateOptionsSheepdog` when driver is "sheepdog"

The members of `BlockdevCreateOptionsSsh` when driver is "ssh"

The members of `BlockdevCreateOptionsVdi` when driver is "vdi"

The members of `BlockdevCreateOptionsVhdx` when driver is "vhdx"

The members of `BlockdevCreateOptionsVmdk` when driver is "vmdk"

The members of `BlockdevCreateOptionsVpc` when driver is "vpc"

Since

2.12

`blockdev-create` (Command)

Starts a job to create an image format on a given node. The job is automatically finalized, but a manual job-dismiss is required.

Arguments

job-id: `string` Identifier for the newly created job.

options: `BlockdevCreateOptions` Options for the image creation.

Since

3.0

`BlockdevAmendOptionsLUKS` (Object)

Driver specific image amend options for LUKS.

Members

The members of `QCryptoBlockAmendOptionsLUKS`

Since

5.1

BlockdevAmendOptionsQcow2 (Object)

Driver specific image amend options for qcow2. For now, only encryption options can be amended

`encrypt` Encryption options to be amended

Members

encrypt: `QCryptoBlockAmendOptions` (optional) Not documented

Since

5.1

BlockdevAmendOptions (Object)

Options for amending an image format

Members

driver: `BlockdevDriver` Block driver of the node to amend.

The members of `BlockdevAmendOptionsLUKS` when `driver` is "luks"

The members of `BlockdevAmendOptionsQcow2` when `driver` is "qcow2"

Since

5.1

x-blockdev-amend (Command)

Starts a job to amend format specific options of an existing open block device The job is automatically finalized, but a manual job-dismiss is required.

Arguments

job-id: `string` Identifier for the newly created job.

node-name: `string` Name of the block node to work on

options: `BlockdevAmendOptions` Options (driver specific)

force: `boolean` (optional) Allow unsafe operations, format specific For luks that allows erase of the last active keyslot (permanent loss of data), and replacement of an active keyslot (possible loss of data if IO error happens)

Since

5.1

BlockErrorAction (Enum)

An enumeration of action that has been taken when a DISK I/O occurs

Values

ignore error has been ignored

report error has been reported to the device

stop error caused VM to be stopped

Since

2.1

BLOCK_IMAGE_CORRUPTED (Event)

Emitted when a disk image is being marked corrupt. The image can be identified by its device or node name. The 'device' field is always present for compatibility reasons, but it can be empty ("") if the image does not have a device name associated.

Arguments

device: string device name. This is always present for compatibility reasons, but it can be empty ("") if the image does not have a device name associated.

node-name: string (optional) node name (Since: 2.4)

msg: string informative message for human consumption, such as the kind of corruption being detected. It should not be parsed by machine as it is not guaranteed to be stable

offset: int (optional) if the corruption resulted from an image access, this is the host's access offset into the image

size: int (optional) if the corruption resulted from an image access, this is the access size

fatal: boolean if set, the image is marked corrupt and therefore unusable after this event and must be repaired (Since 2.2; before, every BLOCK_IMAGE_CORRUPTED event was fatal)

Note

If action is "stop", a STOP event will eventually follow the BLOCK_IO_ERROR event.

Example

```

<- { "event": "BLOCK_IMAGE_CORRUPTED",
      "data": { "device": "ide0-hd0", "node-name": "node0",
                 "msg": "Prevented active L1 table overwrite", "offset": 196608,
                 "size": 65536 },
      "timestamp": { "seconds": 1378126126, "microseconds": 966463 } }

```

Since

1.7

BLOCK_IO_ERROR (Event)

Emitted when a disk I/O error occurs

Arguments

device: string device name. This is always present for compatibility reasons, but it can be empty ("") if the image does not have a device name associated.

node-name: string (optional) node name. Note that errors may be reported for the root node that is directly attached to a guest device rather than for the node where the error occurred. The node name is not present if the drive is empty. (Since: 2.8)

operation: IoOperationType I/O operation

action: BlockErrorAction action that has been taken

nospace: boolean (optional) true if I/O error was caused due to a no-space condition. This key is only present if query-block's io-status is present, please see query-block documentation for more information (since: 2.2)

reason: string human readable string describing the error cause. (This field is a debugging aid for humans, it should not be parsed by applications) (since: 2.2)

Note

If action is "stop", a STOP event will eventually follow the BLOCK_IO_ERROR event

Since

0.13

Example

```
<- { "event": "BLOCK_IO_ERROR",
      "data": { "device": "ide0-hd1",
                 "node-name": "#block212",
                 "operation": "write",
                 "action": "stop" },
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

BLOCK_JOB_COMPLETED (Event)

Emitted when a block job has completed

Arguments

type: JobType job type

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: int maximum progress value

offset: int current progress value. On success this is equal to len. On failure this is less than len

speed: int rate limit, bytes per second

error: string (optional) error message. Only present on failure. This field contains a human-readable error message. There are no semantics other than that streaming has failed and clients should not try to interpret the error string

Since

1.1

Example

```
<- { "event": "BLOCK_JOB_COMPLETED",
      "data": { "type": "stream", "device": "virtio-disk0",
                  "len": 10737418240, "offset": 10737418240,
                  "speed": 0 },
      "timestamp": { "seconds": 1267061043, "microseconds": 959568 } }
```

BLOCK_JOB_CANCELLED (Event)

Emitted when a block job has been cancelled

Arguments

type: JobType job type

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: int maximum progress value

offset: int current progress value. On success this is equal to len. On failure this is less than len

speed: int rate limit, bytes per second

Since

1.1

Example

```
<- { "event": "BLOCK_JOB_CANCELLED",  
      "data": { "type": "stream", "device": "virtio-disk0",  
                 "len": 10737418240, "offset": 134217728,  
                 "speed": 0 },  
      "timestamp": { "seconds": 1267061043, "microseconds": 959568 } }
```

BLOCK_JOB_ERROR (Event)

Emitted when a block job encounters an error

Arguments

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

operation: IoOperationType I/O operation

action: BlockErrorAction action that has been taken

Since

1.3

Example

```
<- { "event": "BLOCK_JOB_ERROR",  
      "data": { "device": "ide0-hd1",  
                 "operation": "write",  
                 "action": "stop" },  
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

BLOCK_JOB_READY (Event)

Emitted when a block job is ready to complete

Arguments

type: JobType job type

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: int maximum progress value

offset: int current progress value. On success this is equal to len. On failure this is less than len

speed: int rate limit, bytes per second

Note

The “ready to complete” status is always reset by a BLOCK_JOB_ERROR event

Since

1.3

Example

```
<- { "event": "BLOCK_JOB_READY",
      "data": { "device": "drive0", "type": "mirror", "speed": 0,
                "len": 2097152, "offset": 2097152 }
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

BLOCK_JOB_PENDING (Event)

Emitted when a block job is awaiting explicit authorization to finalize graph changes via `block-job-finalize`. If this job is part of a transaction, it will not emit this event until the transaction has converged first.

Arguments

type: `JobType` job type

id: `string` The job identifier.

Since

2.12

Example

```
<- { "event": "BLOCK_JOB_WAITING",
      "data": { "device": "drive0", "type": "mirror" },
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

PreallocMode (Enum)

Preallocation mode of QEMU image file

Values

off no preallocation

metadata preallocate only for metadata

fallloc like `full` preallocation but allocate disk space by `posix_fallocate()` rather than writing data.

full preallocate all data by writing it to the device to ensure disk space is really available. This data may or may not be zero, depending on the image format and storage. `full` preallocation also sets up metadata correctly.

Since

2.2

`BLOCK_WRITE_THRESHOLD` (Event)

Emitted when writes on block device reaches or exceeds the configured write threshold. For thin-provisioned devices, this means the device should be extended to avoid pausing for disk exhaustion. The event is one shot. Once triggered, it needs to be re-registered with another `block-set-write-threshold` command.

Arguments

node-name: `string` graph node name on which the threshold was exceeded.

amount-exceeded: `int` amount of data which exceeded the threshold, in bytes.

write-threshold: `int` last configured threshold, in bytes.

Since

2.3

`block-set-write-threshold` (Command)

Change the write threshold for a block drive. An event will be delivered if a write to this block drive crosses the configured threshold. The threshold is an offset, thus must be non-negative. Default is no write threshold. Setting the threshold to zero disables it.

This is useful to transparently resize thin-provisioned drives without the guest OS noticing.

Arguments

node-name: `string` graph node name on which the threshold must be set.

write-threshold: `int` configured threshold for the block device, bytes. Use 0 to disable the threshold.

Since

2.3

Example

```
-> { "execute": "block-set-write-threshold",
      "arguments": { "node-name": "mydev",
                     "write-threshold": 17179869184 } }
<- { "return": {} }
```

x-blockdev-change (Command)

Dynamically reconfigure the block driver state graph. It can be used to add, remove, insert or replace a graph node. Currently only the Quorum driver implements this feature to add or remove its child. This is useful to fix a broken quorum child.

If `node` is specified, it will be inserted under `parent`. `child` may not be specified in this case. If both `parent` and `child` are specified but `node` is not, `child` will be detached from `parent`.

Arguments

parent: string the id or name of the parent node.

child: string (optional) the name of a child under the given parent node.

node: string (optional) the name of the node that will be added.

Note

this command is experimental, and its API is not stable. It does not support all kinds of operations, all kinds of children, nor all block drivers.

FIXME Removing children from a quorum node means introducing gaps in the child indices. This cannot be represented in the 'children' list of `BlockdevOptionsQuorum`, as returned by `.bdrv_refresh_filename()`.

Warning: The data in a new quorum child **MUST** be consistent with that of the rest of the array.

Since

2.7

Example

```
1. Add a new node to a quorum
-> { "execute": "blockdev-add",
    "arguments": {
        "driver": "raw",
        "node-name": "new_node",
        "file": { "driver": "file",
                  "filename": "test.raw" } } }
<- { "return": {} }

-> { "execute": "x-blockdev-change",
    "arguments": { "parent": "disk1",
                  "node": "new_node" } }
<- { "return": {} }

2. Delete a quorum's node
-> { "execute": "x-blockdev-change",
    "arguments": { "parent": "disk1",
                  "child": "children.1" } }
<- { "return": {} }
```

x-blockdev-set-iothread (Command)

Move `node` and its children into the `iothread`. If `iothread` is null then move `node` and its children into the main loop.

The node must not be attached to a BlockBackend.

Arguments

node-name: string the name of the block driver node

iothread: StrOrNull the name of the IOThread object or null for the main loop

force: boolean (optional) true if the node and its children should be moved when a BlockBackend is already attached

Note

this command is experimental and intended for test cases that need control over IOThreads only.

Since

2.12

Example

```
1. Move a node into an IOThread
-> { "execute": "x-blockdev-set-iothread",
    "arguments": { "node-name": "disk1",
                  "iothread": "iothread0" } }
<- { "return": {} }

2. Move a node into the main loop
-> { "execute": "x-blockdev-set-iothread",
    "arguments": { "node-name": "disk1",
                  "iothread": null } }
<- { "return": {} }
```

QuorumOpType (Enum)

An enumeration of the quorum operation types

Values

read read operation

write write operation

flush flush operation

Since

2.6

QUORUM_FAILURE (Event)

Emitted by the Quorum block driver if it fails to establish a quorum

Arguments

reference: string device name if defined else node name

sector-num: int number of the first sector of the failed read operation

sectors-count: int failed read operation sector count

Note

This event is rate-limited.

Since

2.0

Example

```
<- { "event": "QUORUM_FAILURE",
      "data": { "reference": "usr1", "sector-num": 345435, "sectors-count": 5 },
      "timestamp": { "seconds": 1344522075, "microseconds": 745528 } }
```

QUORUM_REPORT_BAD (Event)

Emitted to report a corruption of a Quorum file

Arguments

type: QuorumOpType quorum operation type (Since 2.6)

error: string (optional) error message. Only present on failure. This field contains a human-readable error message. There are no semantics other than that the block layer reported an error and clients should not try to interpret the error string.

node-name: string the graph node name of the block driver state

sector-num: int number of the first sector of the failed read operation

sectors-count: int failed read operation sector count

Note

This event is rate-limited.

Since

2.0

Example

```
1. Read operation

{ "event": "QUORUM_REPORT_BAD",
  "data": { "node-name": "node0", "sector-num": 345435, "sectors-count": 5,
            "type": "read" },
  "timestamp": { "seconds": 1344522075, "microseconds": 745528 } }

2. Flush operation

{ "event": "QUORUM_REPORT_BAD",
  "data": { "node-name": "node0", "sector-num": 0, "sectors-count": 2097120,
            "type": "flush", "error": "Broken pipe" },
  "timestamp": { "seconds": 1456406829, "microseconds": 291763 } }
```

BlockdevSnapshotInternal (Object)

Members

device: **string** the device name or node-name of a root node to generate the snapshot from

name: **string** the name of the internal snapshot to be created

Notes

In transaction, if name is empty, or any snapshot matching name exists, the operation will fail. Only some image formats support it, for example, qcow2, rbd, and sheepdog.

Since

1.7

blockdev-snapshot-internal-sync (Command)

Synchronously take an internal snapshot of a block device, when the format of the image used supports it. If the name is an empty string, or a snapshot with name already exists, the operation will fail.

For the arguments, see the documentation of BlockdevSnapshotInternal.

Returns

- nothing on success
- If `device` is not a valid block device, `GenericError`
- If any snapshot matching `name` exists, or `name` is empty, `GenericError`
- If the format of the image used does not support it, `BlockFormatFeatureNotSupported`

Since

1.7

Example

```
-> { "execute": "blockdev-snapshot-internal-sync",
      "arguments": { "device": "ide-hd0",
                     "name": "snapshot0" }
    }
<- { "return": {} }
```

blockdev-snapshot-delete-internal-sync (Command)

Synchronously delete an internal snapshot of a block device, when the format of the image used support it. The snapshot is identified by name or id or both. One of the name or id is required. Return `SnapshotInfo` for the successfully deleted snapshot.

Arguments

device: string the device name or node-name of a root node to delete the snapshot from

id: string (optional) optional the snapshot's ID to be deleted

name: string (optional) optional the snapshot's name to be deleted

Returns

- `SnapshotInfo` on success
- If `device` is not a valid block device, `GenericError`
- If snapshot not found, `GenericError`
- If the format of the image used does not support it, `BlockFormatFeatureNotSupported`
- If `id` and `name` are both not specified, `GenericError`

Since

1.7

Example

```
-> { "execute": "blockdev-snapshot-delete-internal-sync",
    "arguments": { "device": "ide-hd0",
                   "name": "snapshot0" }
  }
<- { "return": {
    "id": "1",
    "name": "snapshot0",
    "vm-state-size": 0,
    "date-sec": 1000012,
    "date-nsec": 10,
    "vm-clock-sec": 100,
    "vm-clock-nsec": 20,
    "icount": 220414
  }
}
```

Additional block stuff (VM related)

BiosAtaTranslation (Enum)

Policy that BIOS should use to interpret cylinder/head/sector addresses. Note that Bochs BIOS and SeaBIOS will not actually translate logical CHS to physical; instead, they will use logical block addressing.

Values

auto If cylinder/heads/sizes are passed, choose between none and LBA depending on the size of the disk. If they are not passed, choose none if QEMU can guess that the disk had 16 or fewer heads, large if QEMU can guess that the disk had 131072 or fewer tracks across all heads (i.e. cylinders*heads<131072), otherwise LBA.

none The physical disk geometry is equal to the logical geometry.

lba Assume 63 sectors per track and one of 16, 32, 64, 128 or 255 heads (if fewer than 255 are enough to cover the whole disk with 1024 cylinders/head). The number of cylinders/head is then computed based on the number of sectors and heads.

large The number of cylinders per head is scaled down to 1024 by correspondingly scaling up the number of heads.

rechs Same as large, but first convert a 16-head geometry to 15-head, by proportionally scaling up the number of cylinders/head.

Since

2.0

FloppyDriveType (Enum)

Type of Floppy drive to be emulated by the Floppy Disk Controller.

Values

144 1.44MB 3.5" drive

288 2.88MB 3.5" drive

120 1.2MB 5.25" drive

none No drive connected

auto Automatically determined by inserted media at boot

Since

2.6

PRManagerInfo (Object)

Information about a persistent reservation manager

Members

id: string the identifier of the persistent reservation manager

connected: boolean true if the persistent reservation manager is connected to the underlying storage or helper

Since

3.0

query-pr-managers (Command)

Returns a list of information about each persistent reservation manager.

Returns

a list of `PRManagerInfo` for each persistent reservation manager

Since

3.0

eject (Command)

Ejects the medium from a removable drive.

Arguments

device: string (optional) Block device name

id: string (optional) The name or QOM path of the guest device (since: 2.8)

force: boolean (optional) If true, eject regardless of whether the drive is locked. If not specified, the default value is false.

Features

deprecated Member `device` is deprecated. Use `id` instead.

Returns

- Nothing on success
- If `device` is not a valid block device, `DeviceNotFound`

Notes

Ejecting a device with no media results in success

Since

0.14

Example

```
-> { "execute": "eject", "arguments": { "id": "ide1-0-1" } }  
<- { "return": {} }
```

blockdev-open-tray (Command)

Opens a block device's tray. If there is a block driver state tree inserted as a medium, it will become inaccessible to the guest (but it will remain associated to the block device, so closing the tray will make it accessible again).

If the tray was already open before, this will be a no-op.

Once the tray opens, a `DEVICE_TRAY_MOVED` event is emitted. There are cases in which no such event will be generated, these include:

- if the guest has locked the tray, `force` is false and the guest does not respond to the eject request
- if the `BlockBackend` denoted by `device` does not have a guest device attached to it
- if the guest device does not have an actual tray

Arguments

device: string (optional) Block device name

id: string (optional) The name or QOM path of the guest device (since: 2.8)

force: boolean (optional) if false (the default), an eject request will be sent to the guest if it has locked the tray (and the tray will not be opened immediately); if true, the tray will be opened regardless of whether it is locked

Features

deprecated Member `device` is deprecated. Use `id` instead.

Since

2.5

Example

```
-> { "execute": "blockdev-open-tray",
      "arguments": { "id": "ide0-1-0" } }

<- { "timestamp": { "seconds": 1418751016,
                    "microseconds": 716996 },
      "event": "DEVICE_TRAY_MOVED",
      "data": { "device": "ide1-cd0",
                 "id": "ide0-1-0",
                 "tray-open": true } }

<- { "return": {} }
```

blockdev-close-tray (Command)

Closes a block device's tray. If there is a block driver state tree associated with the block device (which is currently ejected), that tree will be loaded as the medium.

If the tray was already closed before, this will be a no-op.

Arguments

device: string (optional) Block device name

id: string (optional) The name or QOM path of the guest device (since: 2.8)

Features

deprecated Member `device` is deprecated. Use `id` instead.

Since

2.5

Example

```
-> { "execute": "blockdev-close-tray",
      "arguments": { "id": "ide0-1-0" } }

<- { "timestamp": { "seconds": 1418751345,
                    "microseconds": 272147 },
      "event": "DEVICE_TRAY_MOVED",
      "data": { "device": "ide1-cd0",
                 "id": "ide0-1-0",
                 "tray-open": false } }

<- { "return": {} }
```

blockdev-remove-medium (Command)

Removes a medium (a block driver state tree) from a block device. That block device's tray must currently be open (unless there is no attached guest device).

If the tray is open and there is no medium inserted, this will be a no-op.

Arguments

id: string The name or QOM path of the guest device

Since

2.12

Example

```
-> { "execute": "blockdev-remove-medium",
      "arguments": { "id": "ide0-1-0" } }

<- { "error": { "class": "GenericError",
                "desc": "Tray of device 'ide0-1-0' is not open" } }

-> { "execute": "blockdev-open-tray",
      "arguments": { "id": "ide0-1-0" } }

<- { "timestamp": { "seconds": 1418751627,
                    "microseconds": 549958 },
      "event": "DEVICE_TRAY_MOVED",
      "data": { "device": "ide1-cd0",
                 "id": "ide0-1-0",
```

(continues on next page)

(continued from previous page)

```

        "tray-open": true } }

<- { "return": {} }

-> { "execute": "blockdev-remove-medium",
      "arguments": { "id": "ide0-1-0" } }

<- { "return": {} }

```

blockdev-insert-medium (Command)

Inserts a medium (a block driver state tree) into a block device. That block device's tray must currently be open (unless there is no attached guest device) and there must be no medium inserted already.

Arguments

id: string The name or QOM path of the guest device

node-name: string name of a node in the block driver state graph

Since

2.12

Example

```

-> { "execute": "blockdev-add",
      "arguments": {
        "node-name": "node0",
        "driver": "raw",
        "file": { "driver": "file",
                  "filename": "fedora.iso" } } }
<- { "return": {} }

-> { "execute": "blockdev-insert-medium",
      "arguments": { "id": "ide0-1-0",
                    "node-name": "node0" } }

<- { "return": {} }

```

BlockdevChangeReadOnlyMode (Enum)

Specifies the new read-only mode of a block device subject to the `blockdev-change-medium` command.

Values

retain Retains the current read-only mode

read-only Makes the device read-only

read-write Makes the device writable

Since

2.3

blockdev-change-medium (Command)

Changes the medium inserted into a block device by ejecting the current medium and loading a new image file which is inserted as the new medium (this command combines `blockdev-open-tray`, `blockdev-remove-medium`, `blockdev-insert-medium` and `blockdev-close-tray`).

Arguments

device: string (optional) Block device name

id: string (optional) The name or QOM path of the guest device (since: 2.8)

filename: string filename of the new image to be loaded

format: string (optional) format to open the new image with (defaults to the probed format)

read-only-mode: BlockdevChangeReadOnlyMode (optional) change the read-only mode of the device; defaults to 'retain'

Features

deprecated Member `device` is deprecated. Use `id` instead.

Since

2.5

Examples

```
1. Change a removable medium
-> { "execute": "blockdev-change-medium",
    "arguments": { "id": "ide0-1-0",
                  "filename": "/srv/images/Fedora-12-x86_64-DVD.iso",
                  "format": "raw" } }
<- { "return": {} }

2. Load a read-only medium into a writable drive
-> { "execute": "blockdev-change-medium",
    "arguments": { "id": "floppyA",
                  "filename": "/srv/images/ro.img",
```

(continues on next page)

(continued from previous page)

```

        "format": "raw",
        "read-only-mode": "retain" } }

<- { "error":
    { "class": "GenericError",
      "desc": "Could not open '/srv/images/ro.img': Permission denied" } }

-> { "execute": "blockdev-change-medium",
    "arguments": { "id": "floppyA",
                  "filename": "/srv/images/ro.img",
                  "format": "raw",
                  "read-only-mode": "read-only" } }

<- { "return": {} }

```

DEVICE_TRAY_MOVED (Event)

Emitted whenever the tray of a removable device is moved by the guest or by HMP/QMP commands

Arguments

device: string Block device name. This is always present for compatibility reasons, but it can be empty ("") if the image does not have a device name associated.

id: string The name or QOM path of the guest device (since 2.8)

tray-open: boolean true if the tray has been opened or false if it has been closed

Since

1.1

Example

```

<- { "event": "DEVICE_TRAY_MOVED",
    "data": { "device": "ide1-cd0",
              "id": "/machine/unattached/device[22]",
              "tray-open": true
            },
    "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }

```

PR_MANAGER_STATUS_CHANGED (Event)

Emitted whenever the connected status of a persistent reservation manager changes.

Arguments

id: string The id of the PR manager object

connected: **boolean** true if the PR manager is connected to a backend

Since

3.0

Example

```
<- { "event": "PR_MANAGER_STATUS_CHANGED",  
      "data": { "id": "pr-helper0",  
                 "connected": true  
            },  
      "timestamp": { "seconds": 1519840375, "microseconds": 450486 } }
```

block_set_io_throttle (Command)

Change I/O throttle limits for a block drive.

Since QEMU 2.4, each device with I/O limits is member of a throttle group.

If two or more devices are members of the same group, the limits will apply to the combined I/O of the whole group in a round-robin fashion. Therefore, setting new I/O limits to a device will affect the whole group.

The name of the group can be specified using the ‘group’ parameter. If the parameter is unset, it is assumed to be the current group of that device. If it’s not in any group yet, the name of the device will be used as the name for its group.

The ‘group’ parameter can also be used to move a device to a different group. In this case the limits specified in the parameters will be applied to the new group only.

I/O limits can be disabled by setting all of them to 0. In this case the device will be removed from its group and the rest of its members will not be affected. The ‘group’ parameter is ignored.

Arguments

The members of BlockIOThrottle

Returns

- Nothing on success
- If `device` is not a valid block device, `DeviceNotFound`

Since

1.1

Example


```

-> { "execute": "block_set_io_throttle",
    "arguments": { "id": "virtio-blk-pci0/virtio-backend",
        "bps": 0,
        "bps_rd": 0,
        "bps_wr": 0,
        "iops": 512,
        "iops_rd": 0,
        "iops_wr": 0,
        "bps_max": 0,
        "bps_rd_max": 0,
        "bps_wr_max": 0,
        "iops_max": 0,
        "iops_rd_max": 0,
        "iops_wr_max": 0,
        "bps_max_length": 0,
        "iops_size": 0 } }
<- { "return": {} }

-> { "execute": "block_set_io_throttle",
    "arguments": { "id": "ide0-l-0",
        "bps": 1000000,
        "bps_rd": 0,
        "bps_wr": 0,
        "iops": 0,
        "iops_rd": 0,
        "iops_wr": 0,
        "bps_max": 8000000,
        "bps_rd_max": 0,
        "bps_wr_max": 0,
        "iops_max": 0,
        "iops_rd_max": 0,
        "iops_wr_max": 0,
        "bps_max_length": 60,
        "iops_size": 0 } }
<- { "return": {} }

```

block-latency-histogram-set (Command)

Manage read, write and flush latency histograms for the device.

If only `id` parameter is specified, remove all present latency histograms for the device. Otherwise, add/reset some of (or all) latency histograms.

Arguments

id: `string` The name or QOM path of the guest device.

boundaries: `array of int (optional)` list of interval boundary values (see description in `BlockLatencyHistogramInfo` definition). If specified, all latency histograms are removed, and empty ones created for all io types with intervals corresponding to `boundaries` (except for io types, for which specific boundaries are set through the following parameters).

boundaries-read: `array of int (optional)` list of interval boundary values for read latency histogram. If specified, old read latency histogram is removed, and empty one created with intervals corresponding to `boundaries-read`. The parameter has higher priority then `boundaries`.

boundaries-write: array of int (optional) list of interval boundary values for write latency histogram.

boundaries-flush: array of int (optional) list of interval boundary values for flush latency histogram.

Returns

error if device is not found or any boundary arrays are invalid.

Since

4.0

Example

```
set new histograms for all io types with intervals
[0, 10), [10, 50), [50, 100), [100, +inf):

-> { "execute": "block-latency-histogram-set",
    "arguments": { "id": "drive0",
                  "boundaries": [10, 50, 100] } }
<- { "return": {} }
```

Example

```
set new histogram only for write, other histograms will remain
not changed (or not created):

-> { "execute": "block-latency-histogram-set",
    "arguments": { "id": "drive0",
                  "boundaries-write": [10, 50, 100] } }
<- { "return": {} }
```

Example

```
set new histograms with the following intervals:
  read, flush: [0, 10), [10, 50), [50, 100), [100, +inf)
  write: [0, 1000), [1000, 5000), [5000, +inf)

-> { "execute": "block-latency-histogram-set",
    "arguments": { "id": "drive0",
                  "boundaries": [10, 50, 100],
                  "boundaries-write": [1000, 5000] } }
<- { "return": {} }
```

Example

```
remove all latency histograms:

-> { "execute": "block-latency-histogram-set",
    "arguments": { "id": "drive0" } }
<- { "return": {} }
```

Block device exports

NbdServerOptions (Object)

Keep this type consistent with the nbd-server-start arguments. The only intended difference is using SocketAddress instead of SocketAddressLegacy.

Members

addr: SocketAddress Address on which to listen.

tls-creds: string (optional) ID of the TLS credentials object (since 2.6).

tls-authz: string (optional) ID of the QAuthZ authorization object used to validate the client's x509 distinguished name. This object is only resolved at time of use, so can be deleted and recreated on the fly while the NBD server is active. If missing, it will default to denying access (since 4.0).

max-connections: int (optional) The maximum number of connections to allow at the same time, 0 for unlimited. (since 5.2; default: 0)

Since

4.2

nbd-server-start (Command)

Start an NBD server listening on the given host and port. Block devices can then be exported using nbd-server-add. The NBD server will present them as named exports; for example, another QEMU instance could refer to them as "nbd:HOST:PORT:exportname=NAME".

Keep this type consistent with the NbdServerOptions type. The only intended difference is using SocketAddressLegacy instead of SocketAddress.

Arguments

addr: SocketAddressLegacy Address on which to listen.

tls-creds: string (optional) ID of the TLS credentials object (since 2.6).

tls-authz: string (optional) ID of the QAuthZ authorization object used to validate the client's x509 distinguished name. This object is only resolved at time of use, so can be deleted and recreated on the fly while the NBD server is active. If missing, it will default to denying access (since 4.0).

max-connections: int (optional) The maximum number of connections to allow at the same time, 0 for unlimited. (since 5.2; default: 0)

Returns

error if the server is already running.

Since

1.3

BlockExportOptionsNbdBase (Object)

An NBD block export (common options shared between nbd-server-add and the NBD branch of block-export-add).

Members

name: string (optional) Export name. If unspecified, the `device` parameter is used as the export name. (Since 2.12)

description: string (optional) Free-form description of the export, up to 4096 bytes. (Since 5.0)

Since

5.0

BlockExportOptionsNbd (Object)

An NBD block export (distinct options used in the NBD branch of block-export-add).

Members

bitmaps: array of string (optional) Also export each of the named dirty bitmaps reachable from `device`, so the NBD client can use `NBD_OPT_SET_META_CONTEXT` with the metadata context name “qemu:dirty-bitmap:BITMAP” to inspect each bitmap.

allocation-depth: boolean (optional) Also export the allocation depth map for `device`, so the NBD client can use `NBD_OPT_SET_META_CONTEXT` with the metadata context name “qemu:allocation-depth” to inspect allocation details. (since 5.2)

The members of `BlockExportOptionsNbdBase`

Since

5.2

BlockExportOptionsVhostUserBlk (Object)

A vhost-user-blk block export.

Members

addr: SocketAddress The vhost-user socket on which to listen. Both ‘unix’ and ‘fd’ SocketAddress types are supported. Passed fds must be UNIX domain sockets.

logical-block-size: int (optional) Logical block size in bytes. Defaults to 512 bytes.

num-queues: int (optional) Number of request virtqueues. Must be greater than 0. Defaults to 1.

Since

5.2

BlockExportOptionsFuse (Object)

Options for exporting a block graph node on some (file) mountpoint as a raw image.

Members

mountpoint: string Path on which to export the block device via FUSE. This must point to an existing regular file.

growable: boolean (optional) Whether writes beyond the EOF should grow the block node accordingly. (default: false)

Since

6.0

If

`defined(CONFIG_FUSE)`

NbdServerAddOptions (Object)

An NBD block export, per legacy nbd-server-add command.

Members

device: string The device name or node name of the node to be exported

writable: boolean (optional) Whether clients should be able to write to the device via the NBD connection (default false).

bitmap: string (optional) Also export a single dirty bitmap reachable from `device`, so the NBD client can use `NBD_OPT_SET_META_CONTEXT` with the metadata context name “qemu:dirty-bitmap:BITMAP” to inspect the bitmap (since 4.0).

The members of `BlockExportOptionsNbdBase`

Since

5.0

`nbd-server-add` (Command)

Export a block node to QEMU's embedded NBD server.

The export name will be used as the id for the resulting block export.

Arguments

The members of `NbdServerAddOptions`

Features

deprecated This command is deprecated. Use `block-export-add` instead.

Returns

error if the server is not running, or export with the same name already exists.

Since

1.3

`BlockExportRemoveMode` (Enum)

Mode for removing a block export.

Values

safe Remove export if there are no existing connections, fail otherwise.

hard Drop all connections immediately and remove export.

Potential additional modes to be added in the future:

hide: Just hide export from new clients, leave existing connections as is. Remove export after all clients are disconnected.

soft: Hide export from new clients, answer with ESHUTDOWN for all further requests from existing clients.

Since

2.12

nbd-server-remove (Command)

Remove NBD export by name.

Arguments

name: string Block export id.

mode: BlockExportRemoveMode (optional) Mode of command operation. See `BlockExportRemoveMode` description. Default is 'safe'.

Features

deprecated This command is deprecated. Use `block-export-del` instead.

Returns

error if

- the server is not running
- export is not found
- mode is 'safe' and there are existing connections

Since

2.12

nbd-server-stop (Command)

Stop QEMU's embedded NBD server, and unregister all devices previously added via `nbd-server-add`.

Since

1.3

BlockExportType (Enum)

An enumeration of block export types

Values

nbd NBD export

vhost-user-blk vhost-user-blk export (since 5.2)

fuse (If: defined(CONFIG_FUSE)) FUSE export (since: 6.0)

Since

4.2

BlockExportOptions (Object)

Describes a block export, i.e. how single node should be exported on an external interface.

Members

id: string A unique identifier for the block export (across all export types)

node-name: string The node name of the block node to be exported (since: 5.2)

writable: boolean (optional) True if clients should be able to write to the export (default false)

writethrough: boolean (optional) If true, caches are flushed after every write request to the export before completion is signalled. (since: 5.2; default: false)

iothread: string (optional) The name of the iothread object where the export will run. The default is to use the thread currently associated with the block node. (since: 5.2)

fixed-iothread: boolean (optional) True prevents the block node from being moved to another thread while the export is active. If true and `iothread` is given, export creation fails if the block node cannot be moved to the iothread. The default is false. (since: 5.2)

type: BlockExportType Not documented

The members of `BlockExportOptionsNbd` when `type` is "nbd"

The members of `BlockExportOptionsVhostUserBlk` when `type` is "vhost-user-blk"

The members of `BlockExportOptionsFuse` when `type` is "fuse" (If: `defined(CONFIG_FUSE)`)

Since

4.2

block-export-add (Command)

Creates a new block export.

Arguments

The members of `BlockExportOptions`

Since

5.2

block-export-del (Command)

Request to remove a block export. This drops the user's reference to the export, but the export may still stay around after this command returns until the shutdown of the export has completed.

Arguments

id: string Block export id.

mode: BlockExportRemoveMode (optional) Mode of command operation. See `BlockExportRemoveMode` description. Default is 'safe'.

Returns

Error if the export is not found or mode is 'safe' and the export is still in use (e.g. by existing client connections)

Since

5.2

BLOCK_EXPORT_DELETED (Event)

Emitted when a block export is removed and its id can be reused.

Arguments

id: string Block export id.

Since

5.2

BlockExportInfo (Object)

Information about a single block export.

Members

id: string The unique identifier for the block export

type: BlockExportType The block export type

node-name: string The node name of the block node that is exported

shutting-down: boolean True if the export is shutting down (e.g. after a block-export-del command, but before the shutdown has completed)

Since

5.2

`query-block-exports` (Command)

Returns

A list of `BlockExportInfo` describing all block exports

Since

5.2

4.8.9 Character devices

`ChardevInfo` (Object)

Information about a character device.

Members

label: `string` the label of the character device

filename: `string` the filename of the character device

frontend-open: `boolean` shows whether the frontend device attached to this backend (eg. with the `chardev=...` option) is in open or closed state (since 2.1)

Notes

`filename` is encoded using the QEMU command line character device encoding. See the QEMU man page for details.

Since

0.14

`query-chardev` (Command)

Returns information about current character devices.

Returns

a list of `ChardevInfo`

Since

0.14

Example

```

-> { "execute": "query-chardev" }
<- {
  "return": [
    {
      "label": "charchannel0",
      "filename": "unix:/var/lib/libvirt/qemu/seabios.rhel6.agent,server",
      "frontend-open": false
    },
    {
      "label": "charmonitor",
      "filename": "unix:/var/lib/libvirt/qemu/seabios.rhel6.monitor,server",
      "frontend-open": true
    },
    {
      "label": "charserial0",
      "filename": "pty:/dev/pts/2",
      "frontend-open": true
    }
  ]
}

```

ChardevBackendInfo (Object)

Information about a character device backend

Members

name: string The backend name

Since

2.0

query-chardev-backends (Command)

Returns information about character device backends.

Returns

a list of ChardevBackendInfo

Since

2.0

Example

```
-> { "execute": "query-chardev-backends" }
<- {
  "return": [
    {
      "name": "udp"
    },
    {
      "name": "tcp"
    },
    {
      "name": "unix"
    },
    {
      "name": "spiceport"
    }
  ]
}
```

DataFormat (Enum)

An enumeration of data format.

Values

utf8 Data is a UTF-8 string (RFC 3629)

base64 Data is Base64 encoded binary (RFC 3548)

Since

1.4

ringbuf-write (Command)

Write to a ring buffer character device.

Arguments

device: string the ring buffer character device name

data: string data to write

format: DataFormat (optional) data encoding (default 'utf8').

- **base64**: data must be base64 encoded text. Its binary decoding gets written.

- utf8: data's UTF-8 encoding is written
- data itself is always Unicode regardless of format, like any other string.

Returns

Nothing on success

Since

1.4

Example

```
-> { "execute": "ringbuf-write",
      "arguments": { "device": "foo",
                     "data": "abcdefgh",
                     "format": "utf8" } }
<- { "return": {} }
```

ringbuf-read (Command)

Read from a ring buffer character device.

Arguments

device: **string** the ring buffer character device name

size: **int** how many bytes to read at most

format: **DataFormat (optional)** data encoding (default 'utf8').

- base64: the data read is returned in base64 encoding.
- utf8: the data read is interpreted as UTF-8. Bug: can screw up when the buffer contains invalid UTF-8 sequences, NUL characters, after the ring buffer lost data, and when reading stops because the size limit is reached.
- The return value is always Unicode regardless of format, like any other string.

Returns

data read from the device

Since

1.4

Example

```
-> { "execute": "ringbuf-read",  
    "arguments": { "device": "foo",  
                  "size": 1000,  
                  "format": "utf8" } }  
<- { "return": "abcdefgh" }
```

ChardevCommon (Object)

Configuration shared across all chardev backends

Members

logfile: string (optional) The name of a logfile to save output

logappend: boolean (optional) true to append instead of truncate (default to false to truncate)

Since

2.6

ChardevFile (Object)

Configuration info for file chardevs.

Members

in: string (optional) The name of the input file

out: string The name of the output file

append: boolean (optional) Open the file in append mode (default false to truncate) (Since 2.6)

The members of ChardevCommon

Since

1.4

ChardevHostdev (Object)

Configuration info for device and pipe chardevs.

Members

device: string The name of the special file for the device, i.e. /dev/ttyS0 on Unix or COM1: on Windows

The members of ChardevCommon

Since

1.4

ChardevSocket (Object)

Configuration info for (stream) socket chardevs.

Members

addr: SocketAddressLegacy socket address to listen on (server=true) or connect to (server=false)

tls-creds: string (optional) the ID of the TLS credentials object (since 2.6)

tls-authz: string (optional) the ID of the QAuthZ authorization object against which the client's x509 distinguished name will be validated. This object is only resolved at time of use, so can be deleted and recreated on the fly while the chardev server is active. If missing, it will default to denying access (since 4.0)

server: boolean (optional) create server socket (default: true)

wait: boolean (optional) wait for incoming connection on server sockets (default: false). Silently ignored with server: false. This use is deprecated.

nodelay: boolean (optional) set TCP_NODELAY socket option (default: false)

telnet: boolean (optional) enable telnet protocol on server sockets (default: false)

tn3270: boolean (optional) enable tn3270 protocol on server sockets (default: false) (Since: 2.10)

websocket: boolean (optional) enable websocket protocol on server sockets (default: false) (Since: 3.1)

reconnect: int (optional) For a client socket, if a socket is disconnected, then attempt a reconnect after the given number of seconds. Setting this to zero disables this function. (default: 0) (Since: 2.2)

The members of ChardevCommon

Since

1.4

ChardevUdp (Object)

Configuration info for datagram socket chardevs.

Members

remote: SocketAddressLegacy remote address

local: SocketAddressLegacy (optional) local address

The members of ChardevCommon

Since

1.5

ChardevMux (Object)

Configuration info for mux chardevs.

Members

chardev: **string** name of the base chardev.

The members of ChardevCommon

Since

1.5

ChardevStdio (Object)

Configuration info for stdio chardevs.

Members

signal: **boolean (optional)** Allow signals (such as SIGINT triggered by ^C) be delivered to qemu. Default: true.

The members of ChardevCommon

Since

1.5

ChardevSpiceChannel (Object)

Configuration info for spice vm channel chardevs.

Members

type: **string** kind of channel (for example vdagent).

The members of ChardevCommon

Since

1.5

If

`defined(CONFIG_SPICE)`

ChardevSpicePort (Object)

Configuration info for spice port chardevs.

Members

fqdn: `string` name of the channel (see docs/spice-port-fqdn.txt)

The members of `ChardevCommon`

Since

1.5

If

`defined(CONFIG_SPICE)`

ChardevVC (Object)

Configuration info for virtual console chardevs.

Members

width: `int (optional)` console width, in pixels

height: `int (optional)` console height, in pixels

cols: `int (optional)` console width, in chars

rows: `int (optional)` console height, in chars

The members of `ChardevCommon`

Since

1.5

ChardevRingbuf (Object)

Configuration info for ring buffer chardevs.

Members

size: `int` (optional) ring buffer size, must be power of two, default is 65536

The members of `ChardevCommon`

Since

1.5

ChardevBackend (Object)

Configuration info for the new chardev backend.

Members

type One of `file`, `serial`, `parallel`, `pipe`, `socket`, `udp`, `pty`, `null`, `mux`, `msmouse`, `wctablet`, `braille`, `testdev`, `stdio`, `console`, `spicevmc`, `spiceport`, `vc`, `ringbuf`, `memory`

data: `ChardevFile` when **type** is `"file"`

data: `ChardevHostdev` when **type** is `"serial"`

data: `ChardevHostdev` when **type** is `"parallel"`

data: `ChardevHostdev` when **type** is `"pipe"`

data: `ChardevSocket` when **type** is `"socket"`

data: `ChardevUdp` when **type** is `"udp"`

data: `ChardevCommon` when **type** is `"pty"`

data: `ChardevCommon` when **type** is `"null"`

data: `ChardevMux` when **type** is `"mux"`

data: `ChardevCommon` when **type** is `"msmouse"`

data: `ChardevCommon` when **type** is `"wctablet"`

data: `ChardevCommon` when **type** is `"braille"`

data: `ChardevCommon` when **type** is `"testdev"`

data: `ChardevStdio` when **type** is `"stdio"`

data: `ChardevCommon` when **type** is `"console"`

data: `ChardevSpiceChannel` when **type** is `"spicevmc"` (If: `defined(CONFIG_SPICE)`)

data: `ChardevSpicePort` when **type** is `"spiceport"` (If: `defined(CONFIG_SPICE)`)

data: `ChardevVC` when **type** is `"vc"`

data: `ChardevRingbuf` when **type** is `"ringbuf"`

data: `ChardevRingbuf` when **type** is `"memory"`

Since

1.4 (testdev since 2.2, wctablet since 2.9)

ChardevReturn (Object)

Return info about the chardev backend just created.

Members

pty: string (optional) name of the slave pseudoterminal device, present if and only if a chardev of type 'pty' was created

Since

1.4

chardev-add (Command)

Add a character device backend

Arguments

id: string the chardev's ID, must be unique

backend: ChardevBackend backend type and parameters

Returns

ChardevReturn.

Since

1.4

Example

```

-> { "execute" : "chardev-add",
    "arguments" : { "id" : "foo",
                    "backend" : { "type" : "null", "data" : {} } } }
<- { "return": {} }

-> { "execute" : "chardev-add",
    "arguments" : { "id" : "bar",
                    "backend" : { "type" : "file",
                                  "data" : { "out" : "/tmp/bar.log" } } } }

```

(continues on next page)

(continued from previous page)

```
<- { "return": {} }

-> { "execute" : "chardev-add",
    "arguments" : { "id" : "baz",
                    "backend" : { "type" : "pty", "data" : {} } } }

<- { "return": { "pty" : "/dev/pty/42" } }
```

chardev-change (Command)

Change a character device backend

Arguments

id: string the chardev's ID, must exist

backend: ChardevBackend new backend type and parameters

Returns

ChardevReturn.

Since

2.10

Example

```
-> { "execute" : "chardev-change",
    "arguments" : { "id" : "baz",
                    "backend" : { "type" : "pty", "data" : {} } } }

<- { "return": { "pty" : "/dev/pty/42" } }

-> { "execute" : "chardev-change",
    "arguments" : {
        "id" : "charchannel2",
        "backend" : {
            "type" : "socket",
            "data" : {
                "addr" : {
                    "type" : "unix" ,
                    "data" : {
                        "path" : "/tmp/charchannel2.socket"
                    }
                }
            },
            "server" : true,
            "wait" : false }}}

<- { "return": {} }
```

chardev-remove (Command)

Remove a character device backend

Arguments

id: string the chardev's ID, must exist and not be in use

Returns

Nothing on success

Since

1.4

Example

```
-> { "execute": "chardev-remove", "arguments": { "id" : "foo" } }  
<- { "return": {} }
```

chardev-send-break (Command)

Send a break to a character device

Arguments

id: string the chardev's ID, must exist

Returns

Nothing on success

Since

2.10

Example

```
-> { "execute": "chardev-send-break", "arguments": { "id" : "foo" } }  
<- { "return": {} }
```

VSERPORT_CHANGE (Event)

Emitted when the guest opens or closes a virtio-serial port.

Arguments

id: **string** device identifier of the virtio-serial port

open: **boolean** true if the guest has opened the virtio-serial port

Note

This event is rate-limited.

Since

2.1

Example

```
<- { "event": "VSERPORT_CHANGE",  
      "data": { "id": "channel0", "open": true },  
      "timestamp": { "seconds": 1401385907, "microseconds": 422329 } }
```

4.8.10 Dump guest memory

DumpGuestMemoryFormat (Enum)

An enumeration of guest-memory-dump's format.

Values

elf elf format

kdump-zlib kdump-compressed format with zlib-compressed

kdump-lzo kdump-compressed format with lzo-compressed

kdump-snappy kdump-compressed format with snappy-compressed

win-dmp Windows full crashdump format, can be used instead of ELF converting (since 2.13)

Since

2.0

`dump-guest-memory` (Command)

Dump guest's memory to vmcore. It is a synchronous operation that can take very long depending on the amount of guest memory.

Arguments

paging: boolean if true, do paging to get guest's memory mapping. This allows using gdb to process the core file.

IMPORTANT: this option can make QEMU allocate several gigabytes of RAM. This can happen for a large guest, or a malicious guest pretending to be large.

Also, `paging=true` has the following limitations:

1. The guest may be in a catastrophic state or can have corrupted memory, which cannot be trusted
2. The guest can be in real-mode even if paging is enabled. For example, the guest uses ACPI to sleep, and ACPI sleep state goes in real-mode
3. Currently only supported on i386 and x86_64.

protocol: string the filename or file descriptor of the vmcore. The supported protocols are:

1. file: the protocol starts with "file:", and the following string is the file's path.
2. fd: the protocol starts with "fd:", and the following string is the fd's name.

detach: boolean (optional) if true, QMP will return immediately rather than waiting for the dump to finish. The user can track progress using "query-dump". (since 2.6).

begin: int (optional) if specified, the starting physical address.

length: int (optional) if specified, the memory size, in bytes. If you don't want to dump all guest's memory, please specify the start `begin` and `length`

format: DumpGuestMemoryFormat (optional) if specified, the format of guest memory dump. But non-elf format is conflict with paging and filter, ie. `paging`, `begin` and `length` is not allowed to be specified with non-elf `format` at the same time (since 2.0)

Note

All boolean arguments default to false

Returns

nothing on success

Since

1.2

Example

```
-> { "execute": "dump-guest-memory",  
    "arguments": { "protocol": "fd:dump" } }  
<- { "return": {} }
```

DumpStatus (Enum)

Describe the status of a long-running background guest memory dump.

Values

none no dump-guest-memory has started yet.

active there is one dump running in background.

completed the last dump has finished successfully.

failed the last dump has failed.

Since

2.6

DumpQueryResult (Object)

The result format for ‘query-dump’.

Members

status: DumpStatus enum of `DumpStatus`, which shows current dump status

completed: int bytes written in latest dump (uncompressed)

total: int total bytes to be written in latest dump (uncompressed)

Since

2.6

query-dump (Command)

Query latest dump status.

Returns

A `DumpStatus` object showing the dump status.

Since

2.6

Example

```
-> { "execute": "query-dump" }
<- { "return": { "status": "active", "completed": 1024000,
                  "total": 2048000 } }
```

DUMP_COMPLETED (Event)

Emitted when background dump has completed

Arguments

result: DumpQueryResult final dump status

error: string (optional) human-readable error string that provides hint on why dump failed. Only presents on failure. The user should not try to interpret the error string.

Since

2.6

Example

```
{ "event": "DUMP_COMPLETED",
  "data": { "result": { "total": 1090650112, "status": "completed",
                        "completed": 1090650112 } } }
```

DumpGuestMemoryCapability (Object)

A list of the available formats for dump-guest-memory

Members

formats: array of DumpGuestMemoryFormat Not documented

Since

2.0

`query-dump-guest-memory-capability` (Command)

Returns the available formats for dump-guest-memory

Returns

A `DumpGuestMemoryCapability` object listing available formats for dump-guest-memory

Since

2.0

Example

```
-> { "execute": "query-dump-guest-memory-capability" }
<- { "return": { "formats":
                  ["elf", "kdump-zlib", "kdump-lzo", "kdump-snappy"] } }
```

4.8.11 Net devices

`set_link` (Command)

Sets the link status of a virtual network adapter.

Arguments

name: `string` the device name of the virtual network adapter

up: `boolean` true to set the link status to be up

Returns

Nothing on success If name is not a valid network device, `DeviceNotFound`

Since

0.14

Notes

Not all network adapters support setting link status. This command will succeed even if the network adapter does not support link status notification.

Example

```
-> { "execute": "set_link",
      "arguments": { "name": "e1000.0", "up": false } }
<- { "return": {} }
```

netdev_add (Command)

Add a network backend.

Additional arguments depend on the type.

Arguments

The members of Netdev

Since

0.14

Returns

Nothing on success If type is not a valid network backend, DeviceNotFound

Example

```
-> { "execute": "netdev_add",
      "arguments": { "type": "user", "id": "netdev1",
                     "dnssearch": "example.org" } }
<- { "return": {} }
```

netdev_del (Command)

Remove a network backend.

Arguments

id: **string** the name of the network backend to remove

Returns

Nothing on success If id is not a valid network backend, DeviceNotFound

Since

0.14

Example

```
-> { "execute": "netdev_del", "arguments": { "id": "netdev1" } }  
<- { "return": {} }
```

NetLegacyNicOptions (Object)

Create a new Network Interface Card.

Members

netdev: **string** (optional) id of -netdev to connect to
macaddr: **string** (optional) MAC address
model: **string** (optional) device model (e1000, rtl8139, virtio etc.)
addr: **string** (optional) PCI device address
vectors: **int** (optional) number of MSI-x vectors, 0 to disable MSI-X

Since

1.2

NetdevUserOptions (Object)

Use the user mode network stack which requires no administrator privilege to run.

Members

hostname: **string** (optional) client hostname reported by the builtin DHCP server
restrict: **boolean** (optional) isolate the guest from the host
ipv4: **boolean** (optional) whether to support IPv4, default true for enabled (since 2.6)
ipv6: **boolean** (optional) whether to support IPv6, default true for enabled (since 2.6)
ip: **string** (optional) legacy parameter, use net= instead
net: **string** (optional) IP network address that the guest will see, in the form addr[/netmask] The netmask is optional, and can be either in the form a.b.c.d or as a number of valid top-most bits. Default is 10.0.2.0/24.
host: **string** (optional) guest-visible address of the host
tftp: **string** (optional) root directory of the built-in TFTP server
bootfile: **string** (optional) BOOTP filename, for use with tftp=

dhcpstart: string (optional) the first of the 16 IPs the built-in DHCP server can assign

dns: string (optional) guest-visible address of the virtual nameserver

dnssearch: array of String (optional) list of DNS suffixes to search, passed as DHCP option to the guest

domainname: string (optional) guest-visible domain name of the virtual nameserver (since 3.0)

ipv6-prefix: string (optional) IPv6 network prefix (default is fec0::) (since 2.6). The network prefix is given in the usual hexadecimal IPv6 address notation.

ipv6-prefixlen: int (optional) IPv6 network prefix length (default is 64) (since 2.6)

ipv6-host: string (optional) guest-visible IPv6 address of the host (since 2.6)

ipv6-dns: string (optional) guest-visible IPv6 address of the virtual nameserver (since 2.6)

smb: string (optional) root directory of the built-in SMB server

smbserver: string (optional) IP address of the built-in SMB server

hostfwd: array of String (optional) redirect incoming TCP or UDP host connections to guest endpoints

guestfwd: array of String (optional) forward guest TCP connections

tftp-server-name: string (optional) RFC2132 “TFTP server name” string (Since 3.1)

Since

1.2

NetdevTapOptions (Object)

Used to configure a host TAP network interface backend.

Members

ifname: string (optional) interface name

fd: string (optional) file descriptor of an already opened tap

fds: string (optional) multiple file descriptors of already opened multiqueue capable tap

script: string (optional) script to initialize the interface

downscript: string (optional) script to shut down the interface

br: string (optional) bridge name (since 2.8)

helper: string (optional) command to execute to configure bridge

sndbuf: int (optional) send buffer limit. Understands [TGMKkb] suffixes.

vnet_hdr: boolean (optional) enable the IFF_VNET_HDR flag on the tap interface

vhost: boolean (optional) enable vhost-net network accelerator

vhostfd: string (optional) file descriptor of an already opened vhost net device

vhostfds: string (optional) file descriptors of multiple already opened vhost net devices

vhostforce: boolean (optional) vhost on for non-MSIX virtio guests

queues: int (optional) number of queues to be created for multiqueue capable tap

poll-us: int (optional) maximum number of microseconds that could be spent on busy polling for tap (since 2.7)

Since

1.2

NetdevSocketOptions (Object)

Socket netdevs are used to establish a network connection to another QEMU virtual machine via a TCP socket.

Members

fd: string (optional) file descriptor of an already opened socket

listen: string (optional) port number, and optional hostname, to listen on

connect: string (optional) port number, and optional hostname, to connect to

mcast: string (optional) UDP multicast address and port number

localaddr: string (optional) source address and port for multicast and udp packets

udp: string (optional) UDP unicast address and port number

Since

1.2

NetdevL2TPv3Options (Object)

Configure an Ethernet over L2TPv3 tunnel.

Members

src: string source address

dst: string destination address

srcport: string (optional) source port - mandatory for udp, optional for ip

dstport: string (optional) destination port - mandatory for udp, optional for ip

ipv6: boolean (optional) force the use of ipv6

udp: boolean (optional) use the udp version of l2tpv3 encapsulation

cookie64: boolean (optional) use 64 bit cookies

counter: boolean (optional) have sequence counter

pincounter: boolean (optional) pin sequence counter to zero - workaround for buggy implementations or networks with packet reorder

txcookie: int (optional) 32 or 64 bit transmit cookie

rxcookie: int (optional) 32 or 64 bit receive cookie

txsession: int 32 bit transmit session

rxsession: int (optional) 32 bit receive session - if not specified set to the same value as transmit

offset: int (optional) additional offset - allows the insertion of additional application-specific data before the packet payload

Since

2.1

NetdevVdeOptions (Object)

Connect to a vde switch running on the host.

Members

sock: string (optional) socket path

port: int (optional) port number

group: string (optional) group owner of socket

mode: int (optional) permissions for socket

Since

1.2

NetdevBridgeOptions (Object)

Connect a host TAP network interface to a host bridge device.

Members

br: string (optional) bridge name

helper: string (optional) command to execute to configure bridge

Since

1.2

NetdevHubPortOptions (Object)

Connect two or more net clients through a software hub.

Members

hubid: int hub identifier number

netdev: string (optional) used to connect hub to a netdev instead of a device (since 2.12)

Since

1.2

NetdevNetmapOptions (Object)

Connect a client to a netmap-enabled NIC or to a VALE switch port

Members

ifname: string Either the name of an existing network interface supported by netmap, or the name of a VALE port (created on the fly). A VALE port name is in the form 'valeXXX:YYY', where XXX and YYY are non-negative integers. XXX identifies a switch and YYY identifies a port of the switch. VALE ports having the same XXX are therefore connected to the same switch.

devname: string (optional) path of the netmap device (default: '/dev/netmap').

Since

2.0

NetdevVhostUserOptions (Object)

Vhost-user network backend

Members

chardev: string name of a unix socket chardev

vhostforce: boolean (optional) vhost on for non-MSIX virtio guests (default: false).

queues: int (optional) number of queues to be created for multiqueue vhost-user (default: 1) (Since 2.5)

Since

2.1

NetdevVhostVDPASOptions (Object)

Vhost-vdpa network backend

vDPA device is a device that uses a datapath which complies with the virtio specifications with a vendor specific control path.

Members

vhostdev: **string** (optional) path of vhost-vdpa device (default: '/dev/vhost-vdpa-0')

queues: **int** (optional) number of queues to be created for multiqueue vhost-vdpa (default: 1)

Since

5.1

NetClientDriver (Enum)

Available netdev drivers.

Values

none Not documented

nic Not documented

user Not documented

tap Not documented

12tpv3 Not documented

socket Not documented

vde Not documented

bridge Not documented

hubport Not documented

netmap Not documented

vhost-user Not documented

vhost-vdpa Not documented

Since

2.7

vhost-vdpa since 5.1

Netdev (Object)

Captures the configuration of a network device.

Members

id: `string` identifier for monitor commands.

type: `NetClientDriver` Specify the driver used for interpreting remaining arguments.

The members of `NetLegacyNicOptions` when type is "nic"

The members of `NetdevUserOptions` when type is "user"

The members of `NetdevTapOptions` when type is "tap"

The members of `NetdevL2TPv3Options` when type is "l2tpv3"

The members of `NetdevSocketOptions` when type is "socket"

The members of `NetdevVdeOptions` when type is "vde"

The members of `NetdevBridgeOptions` when type is "bridge"

The members of `NetdevHubPortOptions` when type is "hubport"

The members of `NetdevNetmapOptions` when type is "netmap"

The members of `NetdevVhostUserOptions` when type is "vhost-user"

The members of `NetdevVhostVDPAOptions` when type is "vhost-vdpa"

Since

1.2

'l2tpv3' - since 2.1

NetFilterDirection (Enum)

Indicates whether a netfilter is attached to a netdev's transmit queue or receive queue or both.

Values

all the filter is attached both to the receive and the transmit queue of the netdev (default).

rx the filter is attached to the receive queue of the netdev, where it will receive packets sent to the netdev.

tx the filter is attached to the transmit queue of the netdev, where it will receive packets sent by the netdev.

Since

2.5

RxState (Enum)

Packets receiving state

Values

normal filter assigned packets according to the mac-table

none don't receive any assigned packet

all receive all assigned packets

Since

1.6

RxFilterInfo (Object)

Rx-filter information for a NIC.

Members

name: string net client name

promiscuous: boolean whether promiscuous mode is enabled

multicast: RxState multicast receive state

unicast: RxState unicast receive state

vlan: RxState vlan receive state (Since 2.0)

broadcast-allowed: boolean whether to receive broadcast

multicast-overflow: boolean multicast table is overflowed or not

unicast-overflow: boolean unicast table is overflowed or not

main-mac: string the main macaddr string

vlan-table: array of int a list of active vlan id

unicast-table: array of string a list of unicast macaddr string

multicast-table: array of string a list of multicast macaddr string

Since

1.6

query-rx-filter (Command)

Return rx-filter information for all NICs (or for the given NIC).

Arguments

name: string (optional) net client name

Returns

list of `RxFilterInfo` for all NICs (or for the given NIC). Returns an error if the given name doesn't exist, or given NIC doesn't support rx-filter querying, or given net client isn't a NIC.

Since

1.6

Example

```
-> { "execute": "query-rx-filter", "arguments": { "name": "vnet0" } }
<- { "return": [
  {
    "promiscuous": true,
    "name": "vnet0",
    "main-mac": "52:54:00:12:34:56",
    "unicast": "normal",
    "vlan": "normal",
    "vlan-table": [
      4,
      0
    ],
    "unicast-table": [
    ],
    "multicast": "normal",
    "multicast-overflow": false,
    "unicast-overflow": false,
    "multicast-table": [
      "01:00:5e:00:00:01",
      "33:33:00:00:00:01",
      "33:33:ff:12:34:56"
    ],
    "broadcast-allowed": false
  }
]
```

NIC_RX_FILTER_CHANGED (Event)

Emitted once until the 'query-rx-filter' command is executed, the first event will always be emitted

Arguments

name: string (optional) net client name

path: string device path

Since

1.6

Example

```
<- { "event": "NIC_RX_FILTER_CHANGED",
      "data": { "name": "vnet0",
                 "path": "/machine/peripheral/vnet0/virtio-backend" },
      "timestamp": { "seconds": 1368697518, "microseconds": 326866 } }
}
```

AnnounceParameters (Object)

Parameters for self-announce timers

Members

initial: int Initial delay (in ms) before sending the first GARP/RARP announcement

max: int Maximum delay (in ms) between GARP/RARP announcement packets

rounds: int Number of self-announcement attempts

step: int Delay increase (in ms) after each self-announcement attempt

interfaces: array of string (optional) An optional list of interface names, which restricts the announcement to the listed interfaces. (Since 4.1)

id: string (optional) A name to be used to identify an instance of announce-timers and to allow it to be modified later. Not for use as part of the migration parameters. (Since 4.1)

Since

4.0

announce-self (Command)

Trigger generation of broadcast RARP frames to update network switches. This can be useful when network bonds fail-over the active slave.

Arguments

The members of **AnnounceParameters**

Example

```
-> { "execute": "announce-self",
      "arguments": {
        "initial": 50, "max": 550, "rounds": 10, "step": 50,
        "interfaces": ["vn2", "vn3"], "id": "bob" } }
<- { "return": {} }
```

Since

4.0

FAILOVER_NEGOTIATED (Event)

Emitted when VIRTIO_NET_F_STANDBY was enabled during feature negotiation. Failover primary devices which were hidden (not hotplugged when requested) before will now be hotplugged by the virtio-net standby device.

device-id: QEMU device id of the unplugged device

Arguments

device-id: string Not documented

Since

4.2

Example

```
<- { "event": "FAILOVER_NEGOTIATED",  
      "data": "net1" }
```

4.8.12 RDMA device

RDMA_GID_STATUS_CHANGED (Event)

Emitted when guest driver adds/deletes GID to/from device

Arguments

netdev: string RoCE Network Device name

gid-status: boolean Add or delete indication

subnet-prefix: int Subnet Prefix

interface-id: int Not documented

interface-id: Interface ID

Since

4.0

Example

```
<- {"timestamp": {"seconds": 1541579657, "microseconds": 986760},
    "event": "RDMA_GID_STATUS_CHANGED",
    "data":
      {"netdev": "bridge0",
       "interface-id": 15880512517475447892,
       "gid-status": true,
       "subnet-prefix": 33022}}
```

4.8.13 Rocker switch device

RockerSwitch (Object)

Rocker switch information.

Members

name: `string` switch name

id: `int` switch ID

ports: `int` number of front-panel ports

Since

2.4

query-rocker (Command)

Return rocker switch information.

Arguments

name: `string` Not documented

Returns

Rocker information

Since

2.4

Example

```
-> { "execute": "query-rocker", "arguments": { "name": "sw1" } }  
<- { "return": { "name": "sw1", "ports": 2, "id": 1327446905938 }}
```

RockerPortDuplex (Enum)

An enumeration of port duplex states.

Values

half half duplex

full full duplex

Since

2.4

RockerPortAutoneg (Enum)

An enumeration of port autoneg states.

Values

off autoneg is off

on autoneg is on

Since

2.4

RockerPort (Object)

Rocker switch port information.

Members

name: string port name

enabled: boolean port is enabled for I/O

link-up: boolean physical link is UP on port

speed: int port link speed in Mbps

duplex: RockerPortDuplex port link duplex

autoneg: RockerPortAutoneg port link autoneg

Since

2.4

query-rocker-ports (Command)

Return rocker switch port information.

Arguments

name: string Not documented

Returns

a list of RockerPort information

Since

2.4

Example

```

-> { "execute": "query-rocker-ports", "arguments": { "name": "sw1" } }
<- { "return": [ { "duplex": "full", "enabled": true, "name": "sw1.1",
                  "autoneg": "off", "link-up": true, "speed": 10000},
                  { "duplex": "full", "enabled": true, "name": "sw1.2",
                  "autoneg": "off", "link-up": true, "speed": 10000}
    ] }

```

RockerOfDpaFlowKey (Object)

Rocker switch OF-DPA flow key

Members

priority: int key priority, 0 being lowest priority

tbl-id: int flow table ID

in-pport: int (optional) physical input port

tunnel-id: int (optional) tunnel ID

vlan-id: int (optional) VLAN ID

eth-type: int (optional) Ethernet header type

eth-src: string (optional) Ethernet header source MAC address

eth-dst: string (optional) Ethernet header destination MAC address

ip-proto: int (optional) IP Header protocol field
ip-tos: int (optional) IP header TOS field
ip-dst: string (optional) IP header destination address

Note

optional members may or may not appear in the flow key depending if they're relevant to the flow key.

Since

2.4

RockerOfDpaFlowMask (Object)

Rocker switch OF-DPA flow mask

Members

in-pport: int (optional) physical input port
tunnel-id: int (optional) tunnel ID
vlan-id: int (optional) VLAN ID
eth-src: string (optional) Ethernet header source MAC address
eth-dst: string (optional) Ethernet header destination MAC address
ip-proto: int (optional) IP Header protocol field
ip-tos: int (optional) IP header TOS field

Note

optional members may or may not appear in the flow mask depending if they're relevant to the flow mask.

Since

2.4

RockerOfDpaFlowAction (Object)

Rocker switch OF-DPA flow action

Members

goto-tbl: `int (optional)` next table ID
group-id: `int (optional)` group ID
tunnel-lport: `int (optional)` tunnel logical port ID
vlan-id: `int (optional)` VLAN ID
new-vlan-id: `int (optional)` new VLAN ID
out-pport: `int (optional)` physical output port

Note

optional members may or may not appear in the flow action depending if they're relevant to the flow action.

Since

2.4

RockerOfDpaFlow (Object)

Rocker switch OF-DPA flow

Members

cookie: `int` flow unique cookie ID
hits: `int` count of matches (hits) on flow
key: `RockerOfDpaFlowKey` flow key
mask: `RockerOfDpaFlowMask` flow mask
action: `RockerOfDpaFlowAction` flow action

Since

2.4

query-rocker-of-dpa-flows (Command)

Return rocker OF-DPA flow information.

Arguments

name: `string` switch name
tbl-id: `int (optional)` flow table ID. If tbl-id is not specified, returns flow information for all tables.

Returns

rocker OF-DPA flow information

Since

2.4

Example

```
-> { "execute": "query-rocker-of-dpa-flows",  
    "arguments": { "name": "sw1" } }  
<- { "return": [ {"key": {"in-pport": 0, "priority": 1, "tbl-id": 0},  
                  "hits": 138,  
                  "cookie": 0,  
                  "action": {"goto-tbl": 10},  
                  "mask": {"in-pport": 4294901760}  
                },  
        {...more...}  
    ] }
```

RockerOfDpaGroup (Object)

Rocker switch OF-DPA group

Members

id: **int** group unique ID

type: **int** group type

vlan-id: **int** (optional) VLAN ID

pport: **int** (optional) physical port number

index: **int** (optional) group index, unique with group type

out-pport: **int** (optional) output physical port number

group-id: **int** (optional) next group ID

set-vlan-id: **int** (optional) VLAN ID to set

pop-vlan: **int** (optional) pop VLAN headr from packet

group-ids: **array of int** (optional) list of next group IDs

set-eth-src: **string** (optional) set source MAC address in Ethernet header

set-eth-dst: **string** (optional) set destination MAC address in Ethernet header

ttl-check: **int** (optional) perform TTL check

Note

optional members may or may not appear in the group depending if they're relevant to the group type.

Since

2.4

query-rocker-of-dpa-groups (Command)

Return rocker OF-DPA group information.

Arguments

name: string switch name

type: int (optional) group type. If type is not specified, returns group information for all group types.

Returns

rocker OF-DPA group information

Since

2.4

Example

```

-> { "execute": "query-rocker-of-dpa-groups",
    "arguments": { "name": "sw1" } }
<- { "return": [ { "type": 0, "out-pport": 2,
                  "pport": 2, "vlan-id": 3841,
                  "pop-vlan": 1, "id": 251723778},
                { "type": 0, "out-pport": 0,
                  "pport": 0, "vlan-id": 3841,
                  "pop-vlan": 1, "id": 251723776},
                { "type": 0, "out-pport": 1,
                  "pport": 1, "vlan-id": 3840,
                  "pop-vlan": 1, "id": 251658241},
                { "type": 0, "out-pport": 0,
                  "pport": 0, "vlan-id": 3840,
                  "pop-vlan": 1, "id": 251658240}
      ] }

```

4.8.14 TPM (trusted platform module) devices

TpmModel (Enum)

An enumeration of TPM models

Values

tpm-tis TPM TIS model

tpm-crb TPM CRB model (since 2.12)

tpm-spapr TPM SPAPR model (since 5.0)

Since

1.5

query-tpm-models (Command)

Return a list of supported TPM models

Returns

a list of TpmModel

Since

1.5

Example

```
-> { "execute": "query-tpm-models" }
<- { "return": [ "tpm-tis", "tpm-crb", "tpm-spapr" ] }
```

TpmType (Enum)

An enumeration of TPM types

Values

passthrough TPM passthrough type

emulator Software Emulator TPM type Since: 2.11

Since

1.5

query-tpm-types (Command)

Return a list of supported TPM types

Returns

a list of TpmType

Since

1.5

Example

```
-> { "execute": "query-tpm-types" }  
<- { "return": [ "passthrough", "emulator" ] }
```

TPMPassthroughOptions (Object)

Information about the TPM passthrough type

Members

path: **string** (optional) string describing the path used for accessing the TPM device

cancel-path: **string** (optional) string showing the TPM's sysfs cancel file for cancellation of TPM commands while they are executing

Since

1.5

TPMEmulatorOptions (Object)

Information about the TPM emulator type

Members

chardev: **string** Name of a unix socket chardev

Since

2.11

TpmTypeOptions (Object)

A union referencing different TPM backend types' configuration options

Members

type

- 'passthrough' The configuration options for the TPM passthrough type
- 'emulator' The configuration options for TPM emulator backend type

data: TPMPassthroughOptions when type is "passthrough"

data: TPMEulatorOptions when type is "emulator"

Since

1.5

TPMInfo (Object)

Information about the TPM

Members

id: string The Id of the TPM

model: TpmModel The TPM frontend model

options: TpmTypeOptions The TPM (backend) type configuration options

Since

1.5

query-tpm (Command)

Return information about the TPM device

Returns

TPMInfo on success

Since

1.5

Example

```

-> { "execute": "query-tpm" }
<- { "return":
  [
    { "model": "tpm-tis",
      "options":
        { "type": "passthrough",
          "data":
            { "cancel-path": "/sys/class/misc/tpm0/device/cancel",
              "path": "/dev/tpm0"
            }
        },
      "id": "tpm0"
    }
  ]
}

```

4.8.15 Remote desktop

set_password (Command)

Sets the password of a remote display session.

Arguments

protocol: string

- ‘vnc’ to modify the VNC server password
- ‘spice’ to modify the Spice server password

password: string the new password

connected: string (optional) how to handle existing clients when changing the password. If nothing is specified, defaults to ‘keep’ ‘fail’ to fail the command if clients are connected ‘disconnect’ to disconnect existing clients ‘keep’ to maintain existing clients

Returns

- Nothing on success
- If Spice is not enabled, DeviceNotFound

Since

0.14

Example

```
-> { "execute": "set_password", "arguments": { "protocol": "vnc",  
                                              "password": "secret" } }  
<- { "return": {} }
```

expire_password (Command)

Expire the password of a remote display server.

Arguments

protocol: string the name of the remote display protocol ‘vnc’ or ‘spice’

time: string when to expire the password.

- ‘now’ to expire the password immediately
- ‘never’ to cancel password expiration
- ‘+INT’ where INT is the number of seconds from now (integer)
- ‘INT’ where INT is the absolute time in seconds

Returns

- Nothing on success
- If `protocol` is ‘spice’ and Spice is not active, `DeviceNotFound`

Since

0.14

Notes

Time is relative to the server and currently there is no way to coordinate server time with client time. It is not recommended to use the absolute time version of the `time` parameter unless you’re sure you are on the same machine as the QEMU instance.

Example

```
-> { "execute": "expire_password", "arguments": { "protocol": "vnc",  
                                              "time": "+60" } }  
<- { "return": {} }
```

screendump (Command)

Write a PPM of the VGA screen to a file.

Arguments

filename: string the path of a new PPM file to store the image

device: string (optional) ID of the display device that should be dumped. If this parameter is missing, the primary display will be used. (Since 2.12)

head: int (optional) head to use in case the device supports multiple heads. If this parameter is missing, head #0 will be used. Also note that the head can only be specified in conjunction with the device ID. (Since 2.12)

Returns

Nothing on success

Since

0.14

Example

```
-> { "execute": "screendump",
      "arguments": { "filename": "/tmp/image" } }
<- { "return": {} }
```

Spice

SpiceBasicInfo (Object)

The basic information for SPICE network connection

Members

host: string IP address

port: string port number

family: NetworkAddressFamily address family

Since

2.1

If

defined(CONFIG_SPICE)

SpiceServerInfo (Object)

Information about a SPICE server

Members

auth: string (optional) authentication method

The members of **SpiceBasicInfo**

Since

2.1

If

`defined(CONFIG_SPICE)`

SpiceChannel (Object)

Information about a SPICE client channel.

Members

connection-id: int SPICE connection id number. All channels with the same id belong to the same SPICE session.

channel-type: int SPICE channel type number. “1” is the main control channel, filter for this one if you want to track spice sessions only

channel-id: int SPICE channel ID number. Usually “0”, might be different when multiple channels of the same type exist, such as multiple display channels in a multihead setup

tls: boolean true if the channel is encrypted, false otherwise.

The members of **SpiceBasicInfo**

Since

0.14

If

`defined(CONFIG_SPICE)`

SpiceQueryMouseMode (Enum)

An enumeration of Spice mouse states.

Values

client Mouse cursor position is determined by the client.

server Mouse cursor position is determined by the server.

unknown No information is available about mouse mode used by the spice server.

Note

spice/enums.h has a `SpiceMouseMoveMode` already, hence the name.

Since

1.1

If

`defined(CONFIG_SPICE)`

SpiceInfo (Object)

Information about the SPICE session.

Members

enabled: boolean true if the SPICE server is enabled, false otherwise

migrated: boolean true if the last guest migration completed and spice migration had completed as well. false otherwise. (since 1.4)

host: string (optional) The hostname the SPICE server is bound to. This depends on the name resolution on the host and may be an IP address.

port: int (optional) The SPICE server's port number.

compiled-version: string (optional) SPICE server version.

tls-port: int (optional) The SPICE server's TLS port number.

auth: string (optional) the current authentication type used by the server

- 'none' if no authentication is being used
- 'spice' uses SASL or direct TLS authentication, depending on command line options

mouse-mode: SpiceQueryMouseMoveMode The mode in which the mouse cursor is displayed currently. Can be determined by the client or the server, or unknown if spice server doesn't provide this information. (since: 1.1)

channels: array of SpiceChannel (optional) a list of `SpiceChannel` for each active spice channel

Since

0.14

If

`defined(CONFIG_SPICE)`

query-spice (Command)

Returns information about the current SPICE server

Returns

`SpiceInfo`

Since

0.14

Example

```
-> { "execute": "query-spice" }
<- { "return": {
    "enabled": true,
    "auth": "spice",
    "port": 5920,
    "tls-port": 5921,
    "host": "0.0.0.0",
    "channels": [
        {
            "port": "54924",
            "family": "ipv4",
            "channel-type": 1,
            "connection-id": 1804289383,
            "host": "127.0.0.1",
            "channel-id": 0,
            "tls": true
        },
        {
            "port": "36710",
            "family": "ipv4",
            "channel-type": 4,
            "connection-id": 1804289383,
            "host": "127.0.0.1",
            "channel-id": 0,
            "tls": false
        },
        [ ... more channels follow ... ]
    ]
}
```

If

```
defined(CONFIG_SPICE)
```

SPICE_CONNECTED (Event)

Emitted when a SPICE client establishes a connection

Arguments

server: **SpiceBasicInfo** server information

client: **SpiceBasicInfo** client information

Since

0.14

Example

```
<- { "timestamp": {"seconds": 1290688046, "microseconds": 388707},
      "event": "SPICE_CONNECTED",
      "data": {
        "server": { "port": "5920", "family": "ipv4", "host": "127.0.0.1"},
        "client": {"port": "52873", "family": "ipv4", "host": "127.0.0.1"}
      }
}
```

If

```
defined(CONFIG_SPICE)
```

SPICE_INITIALIZED (Event)

Emitted after initial handshake and authentication takes place (if any) and the SPICE channel is up and running

Arguments

server: **SpiceServerInfo** server information

client: **SpiceChannel** client information

Since

0.14

Example

```
<- { "timestamp": {"seconds": 1290688046, "microseconds": 417172},
      "event": "SPICE_INITIALIZED",
      "data": {"server": {"auth": "spice", "port": "5921",
                          "family": "ipv4", "host": "127.0.0.1"},
               "client": {"port": "49004", "family": "ipv4", "channel-type": 3,
                           "connection-id": 1804289383, "host": "127.0.0.1",
                           "channel-id": 0, "tls": true}
            }}

```

If

defined(CONFIG_SPICE)

SPICE_DISCONNECTED (Event)

Emitted when the SPICE connection is closed

Arguments

server: `SpiceBasicInfo` server information

client: `SpiceBasicInfo` client information

Since

0.14

Example

```
<- { "timestamp": {"seconds": 1290688046, "microseconds": 388707},
      "event": "SPICE_DISCONNECTED",
      "data": {
        "server": { "port": "5920", "family": "ipv4", "host": "127.0.0.1"},
        "client": { "port": "52873", "family": "ipv4", "host": "127.0.0.1"}
      }
}

```

If

defined(CONFIG_SPICE)

SPICE_MIGRATE_COMPLETED (Event)

Emitted when SPICE migration has completed

Since

1.3

Example

```
<- { "timestamp": {"seconds": 1290688046, "microseconds": 417172},
      "event": "SPICE_MIGRATE_COMPLETED" }
```

If

defined(CONFIG_SPICE)

VNC

VncBasicInfo (Object)

The basic information for vnc network connection

Members

host: string IP address

service: string The service name of the vnc port. This may depend on the host system's service database so symbolic names should not be relied on.

family: NetworkAddressFamily address family

websocket: boolean true in case the socket is a websocket (since 2.3).

Since

2.1

If

defined(CONFIG_VNC)

VncServerInfo (Object)

The network connection information for server

Members

auth: string (optional) authentication method used for the plain (non-websocket) VNC server

The members of VncBasicInfo

Since

2.1

If

`defined(CONFIG_VNC)`

VncClientInfo (Object)

Information about a connected VNC client.

Members

x509_dname: string (optional) If x509 authentication is in use, the Distinguished Name of the client.

sasl_username: string (optional) If SASL authentication is in use, the SASL username used for authentication.

The members of `VncBasicInfo`

Since

0.14

If

`defined(CONFIG_VNC)`

VncInfo (Object)

Information about the VNC session.

Members

enabled: boolean true if the VNC server is enabled, false otherwise

host: string (optional) The hostname the VNC server is bound to. This depends on the name resolution on the host and may be an IP address.

family: NetworkAddressFamily (optional)

- 'ipv6' if the host is listening for IPv6 connections
- 'ipv4' if the host is listening for IPv4 connections
- 'unix' if the host is listening on a unix domain socket
- 'unknown' otherwise

service: string (optional) The service name of the server's port. This may depends on the host system's service database so symbolic names should not be relied on.

auth: string (optional) the current authentication type used by the server

- 'none' if no authentication is being used
- 'vnc' if VNC authentication is being used
- 'vencrypt+plain' if VEncrypt is used with plain text authentication
- 'vencrypt+tls+none' if VEncrypt is used with TLS and no authentication
- 'vencrypt+tls+vnc' if VEncrypt is used with TLS and VNC authentication
- 'vencrypt+tls+plain' if VEncrypt is used with TLS and plain text auth
- 'vencrypt+x509+none' if VEncrypt is used with x509 and no auth
- 'vencrypt+x509+vnc' if VEncrypt is used with x509 and VNC auth
- 'vencrypt+x509+plain' if VEncrypt is used with x509 and plain text auth
- 'vencrypt+tls+sasl' if VEncrypt is used with TLS and SASL auth
- 'vencrypt+x509+sasl' if VEncrypt is used with x509 and SASL auth

clients: array of VncClientInfo (optional) a list of `VncClientInfo` of all currently connected clients

Since

0.14

If

`defined(CONFIG_VNC)`

VncPrimaryAuth (Enum)

vnc primary authentication method.

Values

none Not documented

vnc Not documented

ra2 Not documented

ra2ne Not documented

tight Not documented

ultra Not documented

tls Not documented

vencrypt Not documented

sasl Not documented

Since

2.3

If

`defined(CONFIG_VNC)`

VncVencryptSubAuth (Enum)

vnc sub authentication method with vencrypt.

Values

plain Not documented

tls-none Not documented

x509-none Not documented

tls-vnc Not documented

x509-vnc Not documented

tls-plain Not documented

x509-plain Not documented

tls-sasl Not documented

x509-sasl Not documented

Since

2.3

If

`defined(CONFIG_VNC)`

VncServerInfo2 (Object)

The network connection information for server

Members

auth: VncPrimaryAuth The current authentication type used by the servers

vencrypt: VncVencryptSubAuth (optional) The vencrypt sub authentication type used by the servers, only specified in case `auth == vencrypt`.

The members of **VncBasicInfo**

Since

2.9

If`defined(CONFIG_VNC)`**VncInfo2 (Object)**

Information about a vnc server

Members**id:** `string` vnc server name.**server:** `array of VncServerInfo2` A list of `VncBasincInfo` describing all listening sockets. The list can be empty (in case the vnc server is disabled). It also may have multiple entries: normal + websocket, possibly also ipv4 + ipv6 in the future.**clients:** `array of VncClientInfo` A list of `VncClientInfo` of all currently connected clients. The list can be empty, for obvious reasons.**auth:** `VncPrimaryAuth` The current authentication type used by the non-websockets servers**vencrypt:** `VncVencryptSubAuth (optional)` The vencrypt authentication type used by the servers, only specified in case `auth == vencrypt`.**display:** `string (optional)` The display device the vnc server is linked to.**Since**

2.3

If`defined(CONFIG_VNC)`**query-vnc (Command)**

Returns information about the current VNC server

Returns`VncInfo`

Since

0.14

Example

```
-> { "execute": "query-vnc" }
<- { "return": {
    "enabled": true,
    "host": "0.0.0.0",
    "service": "50402",
    "auth": "vnc",
    "family": "ipv4",
    "clients": [
        {
            "host": "127.0.0.1",
            "service": "50401",
            "family": "ipv4"
        }
    ]
  }
}
```

If`defined(CONFIG_VNC)`**query-vnc-servers (Command)**

Returns a list of vnc servers. The list can be empty.

Returns

a list of `VncInfo2`

Since

2.3

If`defined(CONFIG_VNC)`**change-vnc-password (Command)**

Change the VNC server password.

Arguments

password: string the new password to use with VNC authentication

Since

1.1

Notes

An empty password in this command will set the password to the empty string. Existing clients are unaffected by executing this command.

If

defined(CONFIG_VNC)

VNC_CONNECTED (Event)

Emitted when a VNC client establishes a connection

Arguments

server: VncServerInfo server information

client: VncBasicInfo client information

Note

This event is emitted before any authentication takes place, thus the authentication ID is not provided

Since

0.13

Example

```
<- { "event": "VNC_CONNECTED",  
      "data": {  
        "server": { "auth": "sas1", "family": "ipv4",  
                    "service": "5901", "host": "0.0.0.0" },  
        "client": { "family": "ipv4", "service": "58425",  
                    "host": "127.0.0.1" } },  
      "timestamp": { "seconds": 1262976601, "microseconds": 975795 } }
```

If

```
defined(CONFIG_VNC)
```

VNC_INITIALIZED (Event)

Emitted after authentication takes place (if any) and the VNC session is made active

Arguments

server: `VncServerInfo` server information

client: `VncClientInfo` client information

Since

0.13

Example

```
<- { "event": "VNC_INITIALIZED",  
      "data": {  
        "server": { "auth": "sasl", "family": "ipv4",  
                    "service": "5901", "host": "0.0.0.0"},  
        "client": { "family": "ipv4", "service": "46089",  
                    "host": "127.0.0.1", "sasl_username": "luiz" } },  
      "timestamp": { "seconds": 1263475302, "microseconds": 150772 } }
```

If

```
defined(CONFIG_VNC)
```

VNC_DISCONNECTED (Event)

Emitted when the connection is closed

Arguments

server: `VncServerInfo` server information

client: `VncClientInfo` client information

Since

0.13

Example

```
<- { "event": "VNC_DISCONNECTED",
      "data": {
        "server": { "auth": "sasl", "family": "ipv4",
                     "service": "5901", "host": "0.0.0.0" },
        "client": { "family": "ipv4", "service": "58425",
                     "host": "127.0.0.1", "sasl_username": "luiz" } },
      "timestamp": { "seconds": 1262976601, "microseconds": 975795 } }
```

If

```
defined(CONFIG_VNC)
```

4.8.16 Input

MouseInfo (Object)

Information about a mouse device.

Members

name: string the name of the mouse device

index: int the index of the mouse device

current: boolean true if this device is currently receiving mouse events

absolute: boolean true if this device supports absolute coordinates as input

Since

0.14

query-mice (Command)

Returns information about each active mouse device

Returns

a list of `MouseInfo` for each device

Since

0.14

Example

```
-> { "execute": "query-mice" }
<- { "return": [
    {
      "name": "QEMU Microsoft Mouse",
      "index": 0,
      "current": false,
      "absolute": false
    },
    {
      "name": "QEMU PS/2 Mouse",
      "index": 1,
      "current": true,
      "absolute": true
    }
  ]
}
```

QKeyCode (Enum)

An enumeration of key name.

This is used by the `send-key` command.

Values

unmapped since 2.0

pause since 2.0

ro since 2.4

kp_comma since 2.4

kp_equals since 2.6

power since 2.6

hiragana since 2.9

henkan since 2.9

yen since 2.9

sleep since 2.10

wake since 2.10

audionext since 2.10

audioprev since 2.10

audiostop since 2.10

audioplay since 2.10

audiomute since 2.10

volumeup since 2.10

volumedown since 2.10
mediaselect since 2.10
mail since 2.10
calculator since 2.10
computer since 2.10
ac_home since 2.10
ac_back since 2.10
ac_forward since 2.10
ac_refresh since 2.10
ac_bookmarks since 2.10
muhenkan since 2.12
katakanahiragana since 2.12
shift Not documented
shift_r Not documented
alt Not documented
alt_r Not documented
ctrl Not documented
ctrl_r Not documented
menu Not documented
esc Not documented
1 Not documented
2 Not documented
3 Not documented
4 Not documented
5 Not documented
6 Not documented
7 Not documented
8 Not documented
9 Not documented
0 Not documented
minus Not documented
equal Not documented
backspace Not documented
tab Not documented
q Not documented
w Not documented

e Not documented
r Not documented
t Not documented
y Not documented
u Not documented
i Not documented
o Not documented
p Not documented
bracket_left Not documented
bracket_right Not documented
ret Not documented
a Not documented
s Not documented
d Not documented
f Not documented
g Not documented
h Not documented
j Not documented
k Not documented
l Not documented
semicolon Not documented
apostrophe Not documented
grave_accent Not documented
backslash Not documented
z Not documented
x Not documented
c Not documented
v Not documented
b Not documented
n Not documented
m Not documented
comma Not documented
dot Not documented
slash Not documented
asterisk Not documented
spc Not documented

caps_lock Not documented
f1 Not documented
f2 Not documented
f3 Not documented
f4 Not documented
f5 Not documented
f6 Not documented
f7 Not documented
f8 Not documented
f9 Not documented
f10 Not documented
num_lock Not documented
scroll_lock Not documented
kp_divide Not documented
kp_multiply Not documented
kp_subtract Not documented
kp_add Not documented
kp_enter Not documented
kp_decimal Not documented
sysrq Not documented
kp_0 Not documented
kp_1 Not documented
kp_2 Not documented
kp_3 Not documented
kp_4 Not documented
kp_5 Not documented
kp_6 Not documented
kp_7 Not documented
kp_8 Not documented
kp_9 Not documented
less Not documented
f11 Not documented
f12 Not documented
print Not documented
home Not documented
pgup Not documented

pgdn Not documented
end Not documented
left Not documented
up Not documented
down Not documented
right Not documented
insert Not documented
delete Not documented
stop Not documented
again Not documented
props Not documented
undo Not documented
front Not documented
copy Not documented
open Not documented
paste Not documented
find Not documented
cut Not documented
lf Not documented
help Not documented
meta_l Not documented
meta_r Not documented
compose Not documented

‘sysrq’ was mistakenly added to hack around the fact that the ps2 driver was not generating correct scancodes sequences when ‘alt+print’ was pressed. This flaw is now fixed and the ‘sysrq’ key serves no further purpose. Any further use of ‘sysrq’ will be transparently changed to ‘print’, so they are effectively synonyms.

Since

1.3

KeyValue (Object)

Represents a keyboard key.

Members

type One of `number`, `qcode`

data: int when **type** is `"number"`

data: QKeyCode when **type** is `"qcode"`

Since

1.3

send-key (Command)

Send keys to guest.

Arguments

keys: array of KeyValue An array of `KeyValue` elements. All `KeyValues` in this array are simultaneously sent to the guest. A `KeyValue.number` value is sent directly to the guest, while `KeyValue.qcode` must be a valid `QKeyCode` value

hold-time: int (optional) time to delay key up events, milliseconds. Defaults to 100

Returns

- Nothing on success
- If key is unknown or redundant, `InvalidParameter`

Since

1.3

Example

```
-> { "execute": "send-key",
    "arguments": { "keys": [ { "type": "qcode", "data": "ctrl" },
                             { "type": "qcode", "data": "alt" },
                             { "type": "qcode", "data": "delete" } ] } }
<- { "return": {} }
```

InputButton (Enum)

Button of a pointer input device (mouse, tablet).

Values

side front side button of a 5-button mouse (since 2.9)

extra rear side button of a 5-button mouse (since 2.9)

left Not documented

middle Not documented

right Not documented

wheel-up Not documented

wheel-down Not documented

Since

2.0

InputAxis (Enum)

Position axis of a pointer input device (mouse, tablet).

Values

x Not documented

y Not documented

Since

2.0

InputKeyEvent (Object)

Keyboard input event.

Members

key: KeyValue Which key this event is for.

down: boolean True for key-down and false for key-up events.

Since

2.0

InputBtnEvent (Object)

Pointer button input event.

Members

button: `InputButton` Which button this event is for.

down: `boolean` True for key-down and false for key-up events.

Since

2.0

`InputMoveEvent` (Object)

Pointer motion input event.

Members

axis: `InputAxis` Which axis is referenced by `value`.

value: `int` Pointer position. For absolute coordinates the valid range is 0 -> 0x7fff

Since

2.0

`InputEvent` (Object)

Input event union.

Members

type the input type, one of:

- 'key': Input event of Keyboard
- 'btn': Input event of pointer buttons
- 'rel': Input event of relative pointer motion
- 'abs': Input event of absolute pointer motion

data: `InputKeyEvent` when `type` is "key"

data: `InputBtnEvent` when `type` is "btn"

data: `InputMoveEvent` when `type` is "rel"

data: `InputMoveEvent` when `type` is "abs"

Since

2.0

input-send-event (Command)

Send input event(s) to guest.

The `device` and `head` parameters can be used to send the input event to specific input devices in case (a) multiple input devices of the same kind are added to the virtual machine and (b) you have configured input routing (see `docs/multiseat.txt`) for those input devices. The parameters work exactly like the device and head properties of input devices. If `device` is missing, only devices that have no input routing config are admissible. If `device` is specified, both input devices with and without input routing config are admissible, but devices with input routing config take precedence.

Arguments

device: string (optional) display device to send event(s) to.

head: int (optional) head to send event(s) to, in case the display device supports multiple scanouts.

events: array of InputEvent List of InputEvent union.

Returns

Nothing on success.

Since

2.6

Note

The consoles are visible in the qom tree, under `/backend/console[$index]`. They have a device link and head property, so it is possible to map which console belongs to which device and display.

Example

```
1. Press left mouse button.
-> { "execute": "input-send-event",
    "arguments": { "device": "video0",
                  "events": [ { "type": "btn",
                              "data" : { "down": true, "button": "left" } } ] } }
<- { "return": {} }

-> { "execute": "input-send-event",
    "arguments": { "device": "video0",
                  "events": [ { "type": "btn",
                              "data" : { "down": false, "button": "left" } } ] } }
<- { "return": {} }

2. Press ctrl-alt-del.
-> { "execute": "input-send-event",
```

(continues on next page)

(continued from previous page)

```

    "arguments": { "events": [
      { "type": "key", "data" : { "down": true,
        "key": { "type": "qcode", "data": "ctrl" } } },
      { "type": "key", "data" : { "down": true,
        "key": { "type": "qcode", "data": "alt" } } },
      { "type": "key", "data" : { "down": true,
        "key": { "type": "qcode", "data": "delete" } } } ] } }
<- { "return": {} }

3. Move mouse pointer to absolute coordinates (20000, 400).

-> { "execute": "input-send-event" ,
  "arguments": { "events": [
    { "type": "abs", "data" : { "axis": "x", "value" : 20000 } },
    { "type": "abs", "data" : { "axis": "y", "value" : 400 } } ] } }
<- { "return": {} }

```

GrabToggleKeys (Enum)

Keys to toggle input-linux between host and guest.

Values

ctrl-ctrl Not documented

alt-alt Not documented

shift-shift Not documented

meta-meta Not documented

scrolllock Not documented

ctrl-scrolllock Not documented

Since

4.0

DisplayGTK (Object)

GTK display options.

Members

grab-on-hover: boolean (optional) Grab keyboard input on mouse hover.

zoom-to-fit: boolean (optional) Zoom guest display to fit into the host window. When turned off the host window will be resized instead. In case the display device can notify the guest on window resizes (virtio-gpu) this will default to “on”, assuming the guest will resize the display to match the window size then. Otherwise it defaults to “off”. Since 3.1

Since

2.12

DisplayEGLHeadless (Object)

EGL headless display options.

Members

rendernode: string (optional) Which DRM render node should be used. Default is the first available node on the host.

Since

3.1

DisplayGLMode (Enum)

Display OpenGL mode.

Values

off Disable OpenGL (default).

on Use OpenGL, pick context type automatically. Would better be named ‘auto’ but is called ‘on’ for backward compatibility with bool type.

core Use OpenGL with Core (desktop) Context.

es Use OpenGL with ES (embedded systems) Context.

Since

3.0

DisplayCurses (Object)

Curses display options.

Members

charset: string (optional) Font charset used by guest (default: CP437).

Since

4.0

DisplayType (Enum)

Display (user interface) type.

Values

default The default user interface, selecting from the first available of gtk, sdl, cocoa, and vnc.

none No user interface or video output display. The guest will still see an emulated graphics card, but its output will not be displayed to the QEMU user.

gtk The GTK user interface.

sdl The SDL user interface.

egl-headless No user interface, offload GL operations to a local DRI device. Graphical display need to be paired with VNC or Spice. (Since 3.1)

curses Display video output via curses. For graphics device models which support a text mode, QEMU can display this output using a curses/ncurses interface. Nothing is displayed when the graphics device is in graphical mode or if the graphics device does not support a text mode. Generally only the VGA device models support text mode.

cocoa The Cocoa user interface.

spice-app Set up a Spice server and run the default associated application to connect to it. The server will redirect the serial console and QEMU monitors. (Since 4.0)

Since

2.12

DisplayOptions (Object)

Display (user interface) options.

Members

type: DisplayType Which DisplayType qemu should use.

full-screen: boolean (optional) Start user interface in fullscreen mode (default: off).

window-close: boolean (optional) Allow to quit qemu with window close button (default: on).

show-cursor: boolean (optional) Force showing the mouse cursor (default: off). (since: 5.0)

gl: DisplayGLMode (optional) Enable OpenGL support (default: off).

The members of DisplayGTK when type is "gtk"

The members of DisplayCurses when type is "curses"

The members of DisplayEGLHeadless when type is "egl-headless"

Since

2.12

`query-display-options` (Command)

Returns information about display configuration

Returns

DisplayOptions

Since

3.1

4.8.17 User authorization

`QAuthZListPolicy` (Enum)

The authorization policy result

Values

deny deny access

allow allow access

Since

4.0

`QAuthZListFormat` (Enum)

The authorization policy match format

Values

exact an exact string match

glob string with ? and * shell wildcard support

Since

4.0

QAuthZListRule (Object)

A single authorization rule.

Members

match: string a string or glob to match against a user identity

policy: QAuthZListPolicy the result to return if `match` evaluates to true

format: QAuthZListFormat (optional) the format of the `match` rule (default 'exact')

Since

4.0

QAuthZListRuleListHack (Object)

Not exposed via QMP; hack to generate QAuthZListRuleList for use internally by the code.

Members

unused: array of QAuthZListRule Not documented

Since

4.0

4.8.18 Migration

MigrationStats (Object)

Detailed migration status.

Members

transferred: int amount of bytes already transferred to the target VM

remaining: int amount of bytes remaining to be transferred to the target VM

total: int total amount of bytes involved in the migration process

duplicate: int number of duplicate (zero) pages (since 1.2)

skipped: int number of skipped zero pages (since 1.5)

normal: int number of normal pages (since 1.2)

normal-bytes: int number of normal bytes sent (since 1.2)

dirty-pages-rate: int number of pages dirtied by second by the guest (since 1.3)

mbps: number throughput in megabits/sec. (since 1.6)
dirty-sync-count: int number of times that dirty ram was synchronized (since 2.1)
postcopy-requests: int The number of page requests received from the destination (since 2.7)
page-size: int The number of bytes per page for the various page-based statistics (since 2.10)
multifd-bytes: int The number of bytes sent through multifd (since 3.0)
pages-per-second: int the number of memory pages transferred per second (Since 4.0)

Since

0.14

XBZRLECacheStats (Object)

Detailed XBZRLE migration cache statistics

Members

cache-size: int XBZRLE cache size
bytes: int amount of bytes already transferred to the target VM
pages: int amount of pages transferred to the target VM
cache-miss: int number of cache miss
cache-miss-rate: number rate of cache miss (since 2.1)
encoding-rate: number rate of encoded bytes (since 5.1)
overflow: int number of overflows

Since

1.2

CompressionStats (Object)

Detailed migration compression statistics

Members

pages: int amount of pages compressed and transferred to the target VM
busy: int count of times that no free thread was available to compress data
busy-rate: number rate of thread busy
compressed-size: int amount of bytes after compression
compression-rate: number rate of compressed size

Since

3.1

MigrationStatus (Enum)

An enumeration of migration status.

Values

none no migration has ever happened.

setup migration process has been initiated.

cancelling in the process of cancelling migration.

cancelled cancelling migration is finished.

active in the process of doing migration.

postcopy-active like active, but now in postcopy mode. (since 2.5)

postcopy-paused during postcopy but paused. (since 3.0)

postcopy-recover trying to recover from a paused postcopy. (since 3.0)

completed migration is finished.

failed some error occurred during migration process.

colo VM is in the process of fault tolerance, VM can not get into this state unless colo capability is enabled for migration. (since 2.8)

pre-switchover Paused before device serialisation. (since 2.11)

device During device serialisation when pause-before-switchover is enabled (since 2.11)

wait-unplug wait for device unplug request by guest OS to be completed. (since 4.2)

Since

2.3

VfioStats (Object)

Detailed VFIO devices migration statistics

Members

transferred: int amount of bytes transferred to the target VM by VFIO devices

Since

5.2

MigrationInfo (Object)

Information about current migration process.

Members

status: MigrationStatus (optional) MigrationStatus describing the current migration status. If this field is not returned, no migration process has been initiated

ram: MigrationStats (optional) MigrationStats containing detailed migration status, only returned if status is 'active' or 'completed'(since 1.2)

disk: MigrationStats (optional) MigrationStats containing detailed disk migration status, only returned if status is 'active' and it is a block migration

xbzrle-cache: XBZRLECacheStats (optional) XBZRLECacheStats containing detailed XBZRLE migration statistics, only returned if XBZRLE feature is on and status is 'active' or 'completed' (since 1.2)

total-time: int (optional) total amount of milliseconds since migration started. If migration has ended, it returns the total migration time. (since 1.2)

downtime: int (optional) only present when migration finishes correctly total downtime in milliseconds for the guest. (since 1.3)

expected-downtime: int (optional) only present while migration is active expected downtime in milliseconds for the guest in last walk of the dirty bitmap. (since 1.3)

setup-time: int (optional) amount of setup time in milliseconds *before* the iterations begin but *after* the QMP command is issued. This is designed to provide an accounting of any activities (such as RDMA pinning) which may be expensive, but do not actually occur during the iterative migration rounds themselves. (since 1.6)

cpu-throttle-percentage: int (optional) percentage of time guest cpus are being throttled during auto-converge. This is only present when auto-converge has started throttling guest cpus. (Since 2.7)

error-desc: string (optional) the human readable error description string, when status is 'failed'. Clients should not attempt to parse the error strings. (Since 2.7)

postcopy-blocktime: int (optional) total time when all vCPU were blocked during postcopy live migration. This is only present when the postcopy-blocktime migration capability is enabled. (Since 3.0)

postcopy-vcpu-blocktime: array of int (optional) list of the postcopy blocktime per vCPU. This is only present when the postcopy-blocktime migration capability is enabled. (Since 3.0)

compression: CompressionStats (optional) migration compression statistics, only returned if compression feature is on and status is 'active' or 'completed' (Since 3.1)

socket-address: array of SocketAddress (optional) Only used for tcp, to know what the real port is (Since 4.0)

vfio: VfioStats (optional) VfioStats containing detailed VFIO devices migration statistics, only returned if VFIO device is present, migration is supported by all VFIO devices and status is 'active' or 'completed' (since 5.2)

blocked: boolean True if outgoing migration is blocked (since 6.0)

blocked-reasons: array of string (optional) A list of reasons an outgoing migration is blocked (since 6.0)

Since

0.14

query-migrate (Command)

Returns information about current migration process. If migration is active there will be another json-object with RAM migration status and if block migration is active another one with block migration status.

Returns

MigrationInfo

Since

0.14

Example

```

1. Before the first migration

-> { "execute": "query-migrate" }
<- { "return": {} }

2. Migration is done and has succeeded

-> { "execute": "query-migrate" }
<- { "return": {
    "status": "completed",
    "total-time":12345,
    "setup-time":12345,
    "downtime":12345,
    "ram":{
        "transferred":123,
        "remaining":123,
        "total":246,
        "duplicate":123,
        "normal":123,
        "normal-bytes":123456,
        "dirty-sync-count":15
    }
  }
}

3. Migration is done and has failed

-> { "execute": "query-migrate" }
<- { "return": { "status": "failed" } }

4. Migration is being performed and is not a block migration:

-> { "execute": "query-migrate" }

```

(continues on next page)

(continued from previous page)

```

<- {
  "return":{
    "status":"active",
    "total-time":12345,
    "setup-time":12345,
    "expected-downtime":12345,
    "ram":{
      "transferred":123,
      "remaining":123,
      "total":246,
      "duplicate":123,
      "normal":123,
      "normal-bytes":123456,
      "dirty-sync-count":15
    }
  }
}

```

5. Migration **is** being performed **and is** a block migration:

```

-> { "execute": "query-migrate" }
<- {
  "return":{
    "status":"active",
    "total-time":12345,
    "setup-time":12345,
    "expected-downtime":12345,
    "ram":{
      "total":1057024,
      "remaining":1053304,
      "transferred":3720,
      "duplicate":123,
      "normal":123,
      "normal-bytes":123456,
      "dirty-sync-count":15
    },
    "disk":{
      "total":20971520,
      "remaining":20880384,
      "transferred":91136
    }
  }
}

```

6. Migration **is** being performed **and** XBZRLE **is** active:

```

-> { "execute": "query-migrate" }
<- {
  "return":{
    "status":"active",
    "total-time":12345,
    "setup-time":12345,
    "expected-downtime":12345,
    "ram":{
      "total":1057024,
      "remaining":1053304,
      "transferred":3720,

```

(continues on next page)

(continued from previous page)

```

        "duplicate":10,
        "normal":3333,
        "normal-bytes":3412992,
        "dirty-sync-count":15
    },
    "xbzrle-cache":{
        "cache-size":67108864,
        "bytes":20971520,
        "pages":2444343,
        "cache-miss":2244,
        "cache-miss-rate":0.123,
        "encoding-rate":80.1,
        "overflow":34434
    }
}
}

```

MigrationCapability (Enum)

Migration capabilities enumeration

Values

xbzrle Migration supports xbzrle (Xor Based Zero Run Length Encoding). This feature allows us to minimize migration traffic for certain work loads, by sending compressed difference of the pages

rdma-pin-all Controls whether or not the entire VM memory footprint is mlock()'d on demand or all at once. Refer to docs/rdma.txt for usage. Disabled by default. (since 2.0)

zero-blocks During storage migration encode blocks of zeroes efficiently. This essentially saves 1MB of zeroes per block on the wire. Enabling requires source and target VM to support this feature. To enable it is sufficient to enable the capability on the source VM. The feature is disabled by default. (since 1.6)

compress Use multiple compression threads to accelerate live migration. This feature can help to reduce the migration traffic, by sending compressed pages. Please note that if compress and xbzrle are both on, compress only takes effect in the ram bulk stage, after that, it will be disabled and only xbzrle takes effect, this can help to minimize migration traffic. The feature is disabled by default. (since 2.4)

events generate events for each migration state change (since 2.4)

auto-converge If enabled, QEMU will automatically throttle down the guest to speed up convergence of RAM migration. (since 1.6)

postcopy-ram Start executing on the migration target before all of RAM has been migrated, pulling the remaining pages along as needed. The capacity must have the same setting on both source and target or migration will not even start. NOTE: If the migration fails during postcopy the VM will fail. (since 2.6)

x-colo If enabled, migration will never end, and the state of the VM on the primary side will be migrated continuously to the VM on secondary side, this process is called COarse-Grain LOck Stepping (COLO) for Non-stop Service. (since 2.8)

release-ram if enabled, qemu will free the migrated ram pages on the source during postcopy-ram migration. (since 2.9)

block If enabled, QEMU will also migrate the contents of all block devices. Default is disabled. A possible alternative uses mirror jobs to a builtin NBD server on the destination, which offers more flexibility. (Since 2.10)

return-path If enabled, migration will use the return path even for precopy. (since 2.10)

pause-before-switchover Pause outgoing migration before serialising device state and before disabling block IO (since 2.11)

multifd Use more than one fd for migration (since 4.0)

dirty-bitmaps If enabled, QEMU will migrate named dirty bitmaps. (since 2.12)

postcopy-blocktime Calculate downtime for postcopy live migration (since 3.0)

late-block-activate If enabled, the destination will not activate block devices (and thus take locks) immediately at the end of migration. (since 3.0)

x-ignore-shared If enabled, QEMU will not migrate shared memory (since 4.0)

validate-uuid Send the UUID of the source to allow the destination to ensure it is the same. (since 4.2)

background-snapshot If enabled, the migration stream will be a snapshot of the VM exactly at the point when the migration procedure starts. The VM RAM is saved with running VM. (since 6.0)

Since

1.2

MigrationCapabilityStatus (Object)

Migration capability information

Members

capability: MigrationCapability capability enum

state: boolean capability state bool

Since

1.2

migrate-set-capabilities (Command)

Enable/Disable the following migration capabilities (like xbzrle)

Arguments

capabilities: array of MigrationCapabilityStatus json array of capability modifications to make

Since

1.2

Example

```
-> { "execute": "migrate-set-capabilities" , "arguments":
    { "capabilities": [ { "capability": "xbzrle", "state": true } ] } }
```

query-migrate-capabilities (Command)

Returns information about the current migration capabilities status

Returns

MigrationCapabilitiesStatus

Since

1.2

Example

```
-> { "execute": "query-migrate-capabilities" }
<- { "return": [
    { "state": false, "capability": "xbzrle"},
    { "state": false, "capability": "rdma-pin-all"},
    { "state": false, "capability": "auto-converge"},
    { "state": false, "capability": "zero-blocks"},
    { "state": false, "capability": "compress"},
    { "state": true, "capability": "events"},
    { "state": false, "capability": "postcopy-ram"},
    { "state": false, "capability": "x-colo"}
  ] }
```

MultiFDCompression (Enum)

An enumeration of multifd compression methods.

Values

none no compression.

zlib use zlib compression method.

zstd (If: defined(CONFIG_ZSTD)) use zstd compression method.

Since

5.0

BitmapMigrationBitmapAliasTransform (Object)

Members

persistent: boolean (optional) If present, the bitmap will be made persistent or transient depending on this parameter.

Since

6.0

BitmapMigrationBitmapAlias (Object)

Members

name: string The name of the bitmap.

alias: string An alias name for migration (for example the bitmap name on the opposite site).

transform: BitmapMigrationBitmapAliasTransform (optional) Allows the modification of the migrated bitmap. (since 6.0)

Since

5.2

BitmapMigrationNodeAlias (Object)

Maps a block node name and the bitmaps it has to aliases for dirty bitmap migration.

Members

node-name: string A block node name.

alias: string An alias block node name for migration (for example the node name on the opposite site).

bitmaps: array of BitmapMigrationBitmapAlias Mappings for the bitmaps on this node.

Since

5.2

MigrationParameter (Enum)

Migration parameters enumeration

Values

- announce-initial** Initial delay (in milliseconds) before sending the first announce (Since 4.0)
- announce-max** Maximum delay (in milliseconds) between packets in the announcement (Since 4.0)
- announce-rounds** Number of self-announce packets sent after migration (Since 4.0)
- announce-step** Increase in delay (in milliseconds) between subsequent packets in the announcement (Since 4.0)
- compress-level** Set the compression level to be used in live migration, the compression level is an integer between 0 and 9, where 0 means no compression, 1 means the best compression speed, and 9 means best compression ratio which will consume more CPU.
- compress-threads** Set compression thread count to be used in live migration, the compression thread count is an integer between 1 and 255.
- compress-wait-thread** Controls behavior when all compression threads are currently busy. If true (default), wait for a free compression thread to become available; otherwise, send the page uncompressed. (Since 3.1)
- decompress-threads** Set decompression thread count to be used in live migration, the decompression thread count is an integer between 1 and 255. Usually, decompression is at least 4 times as fast as compression, so set the decompress-threads to the number about 1/4 of compress-threads is adequate.
- throttle-trigger-threshold** The ratio of bytes_dirty_period and bytes_xfer_period to trigger throttling. It is expressed as percentage. The default value is 50. (Since 5.0)
- cpu-throttle-initial** Initial percentage of time guest cpus are throttled when migration auto-converge is activated. The default value is 20. (Since 2.7)
- cpu-throttle-increment** throttle percentage increase each time auto-converge detects that migration is not making progress. The default value is 10. (Since 2.7)
- cpu-throttle-tailslow** Make CPU throttling slower at tail stage At the tail stage of throttling, the Guest is very sensitive to CPU percentage while the `cpu-throttle-increment` is excessive usually at tail stage. If this parameter is true, we will compute the ideal CPU percentage used by the Guest, which may exactly make the dirty rate match the dirty rate threshold. Then we will choose a smaller throttle increment between the one specified by `cpu-throttle-increment` and the one generated by ideal CPU percentage. Therefore, it is compatible to traditional throttling, meanwhile the throttle increment won't be excessive at tail stage. The default value is false. (Since 5.1)
- tls-creds** ID of the 'tls-creds' object that provides credentials for establishing a TLS connection over the migration data channel. On the outgoing side of the migration, the credentials must be for a 'client' endpoint, while for the incoming side the credentials must be for a 'server' endpoint. Setting this will enable TLS for all migrations. The default is unset, resulting in unsecured migration at the QEMU level. (Since 2.7)
- tls-hostname** hostname of the target host for the migration. This is required when using x509 based TLS credentials and the migration URI does not already include a hostname. For example if using fd: or exec: based migration, the hostname must be provided so that the server's x509 certificate identity can be validated. (Since 2.7)
- tls-authz** ID of the 'authz' object subclass that provides access control checking of the TLS x509 certificate distinguished name. This object is only resolved at time of use, so can be deleted and recreated on the fly while the migration server is active. If missing, it will default to denying access (Since 4.0)
- max-bandwidth** to set maximum speed for migration. maximum speed in bytes per second. (Since 2.8)
- downtime-limit** set maximum tolerated downtime for migration. maximum downtime in milliseconds (Since 2.8)
- x-checkpoint-delay** The delay time (in ms) between two COLO checkpoints in periodic mode. (Since 2.8)

block-incremental Affects how much storage is migrated when the block migration capability is enabled. When false, the entire storage backing chain is migrated into a flattened image at the destination; when true, only the active qcow2 layer is migrated and the destination must already have access to the same backing chain as was used on the source. (since 2.10)

multifd-channels Number of channels used to migrate data in parallel. This is the same number that the number of sockets used for migration. The default value is 2 (since 4.0)

xbzrle-cache-size cache size to be used by XBZRLE migration. It needs to be a multiple of the target page size and a power of 2 (Since 2.11)

max-postcopy-bandwidth Background transfer bandwidth during postcopy. Defaults to 0 (unlimited). In bytes per second. (Since 3.0)

max-cpu-throttle maximum cpu throttle percentage. Defaults to 99. (Since 3.1)

multifd-compression Which compression method to use. Defaults to none. (Since 5.0)

multifd-zlib-level Set the compression level to be used in live migration, the compression level is an integer between 0 and 9, where 0 means no compression, 1 means the best compression speed, and 9 means best compression ratio which will consume more CPU. Defaults to 1. (Since 5.0)

multifd-zstd-level Set the compression level to be used in live migration, the compression level is an integer between 0 and 20, where 0 means no compression, 1 means the best compression speed, and 20 means best compression ratio which will consume more CPU. Defaults to 1. (Since 5.0)

block-bitmap-mapping Maps block nodes and bitmaps on them to aliases for the purpose of dirty bitmap migration. Such aliases may for example be the corresponding names on the opposite site. The mapping must be one-to-one, but not necessarily complete: On the source, unmapped bitmaps and all bitmaps on unmapped nodes will be ignored. On the destination, encountering an unmapped alias in the incoming migration stream will result in a report, and all further bitmap migration data will then be discarded. Note that the destination does not know about bitmaps it does not receive, so there is no limitation or requirement regarding the number of bitmaps received, or how they are named, or on which nodes they are placed. By default (when this parameter has never been set), bitmap names are mapped to themselves. Nodes are mapped to their block device name if there is one, and to their node name otherwise. (Since 5.2)

Since

2.4

MigrateSetParameters (Object)

Members

announce-initial: int (optional) Initial delay (in milliseconds) before sending the first announce (Since 4.0)

announce-max: int (optional) Maximum delay (in milliseconds) between packets in the announcement (Since 4.0)

announce-rounds: int (optional) Number of self-announce packets sent after migration (Since 4.0)

announce-step: int (optional) Increase in delay (in milliseconds) between subsequent packets in the announcement (Since 4.0)

compress-level: int (optional) compression level

compress-threads: int (optional) compression thread count

- compress-wait-thread: boolean (optional)** Controls behavior when all compression threads are currently busy. If true (default), wait for a free compression thread to become available; otherwise, send the page uncompressed. (Since 3.1)
- decompress-threads: int (optional)** decompression thread count
- throttle-trigger-threshold: int (optional)** The ratio of `bytes_dirty_period` and `bytes_xfer_period` to trigger throttling. It is expressed as percentage. The default value is 50. (Since 5.0)
- cpu-throttle-initial: int (optional)** Initial percentage of time guest cpus are throttled when migration auto-converge is activated. The default value is 20. (Since 2.7)
- cpu-throttle-increment: int (optional)** throttle percentage increase each time auto-converge detects that migration is not making progress. The default value is 10. (Since 2.7)
- cpu-throttle-tailslow: boolean (optional)** Make CPU throttling slower at tail stage At the tail stage of throttling, the Guest is very sensitive to CPU percentage while the `cpu-throttle-increment` is excessive usually at tail stage. If this parameter is true, we will compute the ideal CPU percentage used by the Guest, which may exactly make the dirty rate match the dirty rate threshold. Then we will choose a smaller throttle increment between the one specified by `cpu-throttle-increment` and the one generated by ideal CPU percentage. Therefore, it is compatible to traditional throttling, meanwhile the throttle increment won't be excessive at tail stage. The default value is false. (Since 5.1)
- tls-creds: StrOrNull (optional)** ID of the 'tls-creds' object that provides credentials for establishing a TLS connection over the migration data channel. On the outgoing side of the migration, the credentials must be for a 'client' endpoint, while for the incoming side the credentials must be for a 'server' endpoint. Setting this to a non-empty string enables TLS for all migrations. An empty string means that QEMU will use plain text mode for migration, rather than TLS (Since 2.9) Previously (since 2.7), this was reported by omitting `tls-creds` instead.
- tls-hostname: StrOrNull (optional)** hostname of the target host for the migration. This is required when using x509 based TLS credentials and the migration URI does not already include a hostname. For example if using `fd:` or `exec:` based migration, the hostname must be provided so that the server's x509 certificate identity can be validated. (Since 2.7) An empty string means that QEMU will use the hostname associated with the migration URI, if any. (Since 2.9) Previously (since 2.7), this was reported by omitting `tls-hostname` instead.
- max-bandwidth: int (optional)** to set maximum speed for migration. maximum speed in bytes per second. (Since 2.8)
- downtime-limit: int (optional)** set maximum tolerated downtime for migration. maximum downtime in milliseconds (Since 2.8)
- x-checkpoint-delay: int (optional)** the delay time between two COLO checkpoints. (Since 2.8)
- block-incremental: boolean (optional)** Affects how much storage is migrated when the block migration capability is enabled. When false, the entire storage backing chain is migrated into a flattened image at the destination; when true, only the active qcow2 layer is migrated and the destination must already have access to the same backing chain as was used on the source. (since 2.10)
- multifd-channels: int (optional)** Number of channels used to migrate data in parallel. This is the same number that the number of sockets used for migration. The default value is 2 (since 4.0)
- xbzrle-cache-size: int (optional)** cache size to be used by XBZRLE migration. It needs to be a multiple of the target page size and a power of 2 (Since 2.11)
- max-postcopy-bandwidth: int (optional)** Background transfer bandwidth during postcopy. Defaults to 0 (unlimited). In bytes per second. (Since 3.0)
- max-cpu-throttle: int (optional)** maximum cpu throttle percentage. The default value is 99. (Since 3.1)
- multifd-compression: MultiFDCompression (optional)** Which compression method to use. Defaults to none. (Since 5.0)

multifd-zlib-level: int (optional) Set the compression level to be used in live migration, the compression level is an integer between 0 and 9, where 0 means no compression, 1 means the best compression speed, and 9 means best compression ratio which will consume more CPU. Defaults to 1. (Since 5.0)

multifd-zstd-level: int (optional) Set the compression level to be used in live migration, the compression level is an integer between 0 and 20, where 0 means no compression, 1 means the best compression speed, and 20 means best compression ratio which will consume more CPU. Defaults to 1. (Since 5.0)

block-bitmap-mapping: array of BitmapMigrationNodeAlias (optional) Maps block nodes and bitmaps on them to aliases for the purpose of dirty bitmap migration. Such aliases may for example be the corresponding names on the opposite site. The mapping must be one-to-one, but not necessarily complete: On the source, unmapped bitmaps and all bitmaps on unmapped nodes will be ignored. On the destination, encountering an unmapped alias in the incoming migration stream will result in a report, and all further bitmap migration data will then be discarded. Note that the destination does not know about bitmaps it does not receive, so there is no limitation or requirement regarding the number of bitmaps received, or how they are named, or on which nodes they are placed. By default (when this parameter has never been set), bitmap names are mapped to themselves. Nodes are mapped to their block device name if there is one, and to their node name otherwise. (Since 5.2)

tls-authz: StrOrNull (optional) Not documented

Since

2.4

migrate-set-parameters (Command)

Set various migration parameters.

Arguments

The members of MigrateSetParameters

Since

2.4

Example

```
-> { "execute": "migrate-set-parameters" ,  
      "arguments": { "compress-level": 1 } }
```

MigrationParameters (Object)

The optional members aren't actually optional.

Members

- announce-initial: int (optional)** Initial delay (in milliseconds) before sending the first announce (Since 4.0)
- announce-max: int (optional)** Maximum delay (in milliseconds) between packets in the announcement (Since 4.0)
- announce-rounds: int (optional)** Number of self-announce packets sent after migration (Since 4.0)
- announce-step: int (optional)** Increase in delay (in milliseconds) between subsequent packets in the announcement (Since 4.0)
- compress-level: int (optional)** compression level
- compress-threads: int (optional)** compression thread count
- compress-wait-thread: boolean (optional)** Controls behavior when all compression threads are currently busy. If true (default), wait for a free compression thread to become available; otherwise, send the page uncompressed. (Since 3.1)
- decompress-threads: int (optional)** decompression thread count
- throttle-trigger-threshold: int (optional)** The ratio of bytes_dirty_period and bytes_xfer_period to trigger throttling. It is expressed as percentage. The default value is 50. (Since 5.0)
- cpu-throttle-initial: int (optional)** Initial percentage of time guest cpus are throttled when migration auto-converge is activated. (Since 2.7)
- cpu-throttle-increment: int (optional)** throttle percentage increase each time auto-converge detects that migration is not making progress. (Since 2.7)
- cpu-throttle-tailslow: boolean (optional)** Make CPU throttling slower at tail stage At the tail stage of throttling, the Guest is very sensitive to CPU percentage while the `cpu-throttle-increment` is excessive usually at tail stage. If this parameter is true, we will compute the ideal CPU percentage used by the Guest, which may exactly make the dirty rate match the dirty rate threshold. Then we will choose a smaller throttle increment between the one specified by `cpu-throttle-increment` and the one generated by ideal CPU percentage. Therefore, it is compatible to traditional throttling, meanwhile the throttle increment won't be excessive at tail stage. The default value is false. (Since 5.1)
- tls-creds: string (optional)** ID of the 'tls-creds' object that provides credentials for establishing a TLS connection over the migration data channel. On the outgoing side of the migration, the credentials must be for a 'client' endpoint, while for the incoming side the credentials must be for a 'server' endpoint. An empty string means that QEMU will use plain text mode for migration, rather than TLS (Since 2.7) Note: 2.8 reports this by omitting `tls-creds` instead.
- tls-hostname: string (optional)** hostname of the target host for the migration. This is required when using x509 based TLS credentials and the migration URI does not already include a hostname. For example if using `fd:` or `exec:` based migration, the hostname must be provided so that the server's x509 certificate identity can be validated. (Since 2.7) An empty string means that QEMU will use the hostname associated with the migration URI, if any. (Since 2.9) Note: 2.8 reports this by omitting `tls-hostname` instead.
- tls-authz: string (optional)** ID of the 'authz' object subclass that provides access control checking of the TLS x509 certificate distinguished name. (Since 4.0)
- max-bandwidth: int (optional)** to set maximum speed for migration. maximum speed in bytes per second. (Since 2.8)
- downtime-limit: int (optional)** set maximum tolerated downtime for migration. maximum downtime in milliseconds (Since 2.8)
- x-checkpoint-delay: int (optional)** the delay time between two COLO checkpoints. (Since 2.8)

block-incremental: boolean (optional) Affects how much storage is migrated when the block migration capability is enabled. When false, the entire storage backing chain is migrated into a flattened image at the destination; when true, only the active qcow2 layer is migrated and the destination must already have access to the same backing chain as was used on the source. (since 2.10)

multifd-channels: int (optional) Number of channels used to migrate data in parallel. This is the same number that the number of sockets used for migration. The default value is 2 (since 4.0)

xbzrle-cache-size: int (optional) cache size to be used by XBZRLE migration. It needs to be a multiple of the target page size and a power of 2 (Since 2.11)

max-postcopy-bandwidth: int (optional) Background transfer bandwidth during postcopy. Defaults to 0 (unlimited). In bytes per second. (Since 3.0)

max-cpu-throttle: int (optional) maximum cpu throttle percentage. Defaults to 99. (Since 3.1)

multifd-compression: MultiFDCompression (optional) Which compression method to use. Defaults to none. (Since 5.0)

multifd-zlib-level: int (optional) Set the compression level to be used in live migration, the compression level is an integer between 0 and 9, where 0 means no compression, 1 means the best compression speed, and 9 means best compression ratio which will consume more CPU. Defaults to 1. (Since 5.0)

multifd-zstd-level: int (optional) Set the compression level to be used in live migration, the compression level is an integer between 0 and 20, where 0 means no compression, 1 means the best compression speed, and 20 means best compression ratio which will consume more CPU. Defaults to 1. (Since 5.0)

block-bitmap-mapping: array of BitmapMigrationNodeAlias (optional) Maps block nodes and bitmaps on them to aliases for the purpose of dirty bitmap migration. Such aliases may for example be the corresponding names on the opposite site. The mapping must be one-to-one, but not necessarily complete: On the source, unmapped bitmaps and all bitmaps on unmapped nodes will be ignored. On the destination, encountering an unmapped alias in the incoming migration stream will result in a report, and all further bitmap migration data will then be discarded. Note that the destination does not know about bitmaps it does not receive, so there is no limitation or requirement regarding the number of bitmaps received, or how they are named, or on which nodes they are placed. By default (when this parameter has never been set), bitmap names are mapped to themselves. Nodes are mapped to their block device name if there is one, and to their node name otherwise. (Since 5.2)

Since

2.4

query-migrate-parameters (Command)

Returns information about the current migration parameters

Returns

MigrationParameters

Since

2.4

Example

```

-> { "execute": "query-migrate-parameters" }
<- { "return": {
    "decompress-threads": 2,
    "cpu-throttle-increment": 10,
    "compress-threads": 8,
    "compress-level": 1,
    "cpu-throttle-initial": 20,
    "max-bandwidth": 33554432,
    "downtime-limit": 300
  }
}

```

client_migrate_info (Command)

Set migration information for remote display. This makes the server ask the client to automatically reconnect using the new parameters once migration finished successfully. Only implemented for SPICE.

Arguments

protocol: string must be “spice”

hostname: string migration target hostname

port: int (optional) spice tcp port for plaintext channels

tls-port: int (optional) spice tcp port for tls-secured channels

cert-subject: string (optional) server certificate subject

Since

0.14

Example

```

-> { "execute": "client_migrate_info",
    "arguments": { "protocol": "spice",
                  "hostname": "virt42.lab.kraxel.org",
                  "port": 1234 } }
<- { "return": {} }

```

migrate-start-postcopy (Command)

Followup to a migration command to switch the migration to postcopy mode. The postcopy-ram capability must be set on both source and destination before the original migration command.

Since

2.5

Example

```
-> { "execute": "migrate-start-postcopy" }  
<- { "return": {} }
```

MIGRATION (Event)

Emitted when a migration event happens

Arguments

status: MigrationStatus MigrationStatus describing the current migration status.

Since

2.4

Example

```
<- { "timestamp": { "seconds": 1432121972, "microseconds": 744001 },  
      "event": "MIGRATION",  
      "data": { "status": "completed" } }
```

MIGRATION_PASS (Event)

Emitted from the source side of a migration at the start of each pass (when it syncs the dirty bitmap)

Arguments

pass: int An incrementing count (starting at 1 on the first pass)

Since

2.6

Example

```
{ "timestamp": { "seconds": 1449669631, "microseconds": 239225 },  
  "event": "MIGRATION_PASS", "data": { "pass": 2 } }
```


COLOMessage (Enum)

The message transmission between Primary side and Secondary side.

Values

checkpoint-ready Secondary VM (SVM) is ready for checkpointing

checkpoint-request Primary VM (PVM) tells SVM to prepare for checkpointing

checkpoint-reply SVM gets PVM's checkpoint request

vmstate-send VM's state will be sent by PVM.

vmstate-size The total size of VMstate.

vmstate-received VM's state has been received by SVM.

vmstate-loaded VM's state has been loaded by SVM.

Since

2.8

COLOMode (Enum)

The COLO current mode.

Values

none COLO is disabled.

primary COLO node in primary side.

secondary COLO node in slave side.

Since

2.8

FailoverStatus (Enum)

An enumeration of COLO failover status

Values

none no failover has ever happened

require got failover requirement but not handled

active in the process of doing failover

completed finish the process of failover

relaunch restart the failover process, from ‘none’ -> ‘completed’ (Since 2.9)

Since

2.8

COLO_EXIT (Event)

Emitted when VM finishes COLO mode due to some errors happening or at the request of users.

Arguments

mode: COLOMode report COLO mode when COLO exited.

reason: COLOExitReason describes the reason for the COLO exit.

Since

3.1

Example

```
<- { "timestamp": {"seconds": 2032141960, "microseconds": 417172},  
      "event": "COLO_EXIT", "data": {"mode": "primary", "reason": "request" } }
```

COLOExitReason (Enum)

The reason for a COLO exit.

Values

none failover has never happened. This state does not occur in the COLO_EXIT event, and is only visible in the result of query-colo-status.

request COLO exit is due to an external request.

error COLO exit is due to an internal error.

processing COLO is currently handling a failover (since 4.0).

Since

3.1

x-colo-lost-heartbeat (Command)

Tell qemu that heartbeat is lost, request it to do takeover procedures. If this command is sent to the PVM, the Primary side will exit COLO mode. If sent to the Secondary, the Secondary side will run failover work, then takes over server operation to become the service VM.

Since

2.8

Example

```
-> { "execute": "x-colo-lost-heartbeat" }
<- { "return": {} }
```

migrate_cancel (Command)

Cancel the current executing migration process.

Returns

nothing on success

Notes

This command succeeds even if there is no migration process running.

Since

0.14

Example

```
-> { "execute": "migrate_cancel" }
<- { "return": {} }
```

migrate-continue (Command)

Continue migration when it's in a paused state.

Arguments

state: **MigrationStatus** The state the migration is currently expected to be in

Returns

nothing on success

Since

2.11

Example

```
-> { "execute": "migrate-continue" , "arguments":  
    { "state": "pre-switchover" } }  
<- { "return": {} }
```

`migrate_set_downtime` (Command)

Set maximum tolerated downtime for migration.

Arguments

value: number maximum downtime in seconds

Features

deprecated This command is deprecated. Use ‘migrate-set-parameters’ instead.

Returns

nothing on success

Since

0.14

Example

```
-> { "execute": "migrate_set_downtime", "arguments": { "value": 0.1 } }  
<- { "return": {} }
```

`migrate_set_speed` (Command)

Set maximum speed for migration.

Arguments

value: int maximum speed in bytes per second.

Features

deprecated This command is deprecated. Use ‘migrate-set-parameters’ instead.

Returns

nothing on success

Since

0.14

Example

```
-> { "execute": "migrate_set_speed", "arguments": { "value": 1024 } }
<- { "return": {} }
```

migrate-set-cache-size (Command)

Set cache size to be used by XBZRLE migration

Arguments

value: int cache size in bytes

Features

deprecated This command is deprecated. Use ‘migrate-set-parameters’ instead.

The size will be rounded down to the nearest power of 2. The cache size can be modified before and during ongoing migration

Returns

nothing on success

Since

1.2

Example

```
-> { "execute": "migrate-set-cache-size",  
      "arguments": { "value": 536870912 } }  
<- { "return": {} }
```

query-migrate-cache-size (Command)

Query migration XBZRLE cache size

Features

deprecated This command is deprecated. Use ‘query-migrate-parameters’ instead.

Returns

XBZRLE cache size in bytes

Since

1.2

Example

```
-> { "execute": "query-migrate-cache-size" }  
<- { "return": 67108864 }
```

migrate (Command)

Migrates the current running guest to another Virtual Machine.

Arguments

uri: string the Uniform Resource Identifier of the destination VM

blk: boolean (optional) do block migration (full disk copy)

inc: boolean (optional) incremental disk copy migration

detach: boolean (optional) this argument exists only for compatibility reasons and is ignored by QEMU

resume: boolean (optional) resume one paused migration, default “off”. (since 3.0)

Returns

nothing on success

Since

0.14

Notes

1. The ‘query-migrate’ command should be used to check migration’s progress and final result (this information is provided by the ‘status’ member)
2. All boolean arguments default to false
3. The user Monitor’s “detach” argument is invalid in QMP and should not be used

Example

```
-> { "execute": "migrate", "arguments": { "uri": "tcp:0:4446" } }  
<- { "return": {} }
```

migrate-incoming (Command)

Start an incoming migration, the qemu must have been started with -incoming defer

Arguments

uri: string The Uniform Resource Identifier identifying the source or address to listen on

Returns

nothing on success

Since

2.3

Notes

1. It’s a bad idea to use a string for the uri, but it needs to stay compatible with -incoming and the format of the uri is already exposed above libvirt.
2. QEMU must be started with -incoming defer to allow migrate-incoming to be used.
3. The uri format is the same as for -incoming

Example

```
-> { "execute": "migrate-incoming",  
    "arguments": { "uri": "tcp::4446" } }  
<- { "return": {} }
```

xen-save-devices-state (Command)

Save the state of all devices to file. The RAM and the block devices of the VM are not saved by this command.

Arguments

filename: string the file to save the state of the devices to as binary data. See xen-save-devices-state.txt for a description of the binary format.

live: boolean (optional) Optional argument to ask QEMU to treat this command as part of a live migration. Default to true. (since 2.11)

Returns

Nothing on success

Since

1.1

Example

```
-> { "execute": "xen-save-devices-state",  
    "arguments": { "filename": "/tmp/save" } }  
<- { "return": {} }
```

xen-set-global-dirty-log (Command)

Enable or disable the global dirty log mode.

Arguments

enable: boolean true to enable, false to disable.

Returns

nothing

Since

1.3

Example

```
-> { "execute": "xen-set-global-dirty-log",
      "arguments": { "enable": true } }
<- { "return": {} }
```

xen-load-devices-state (Command)

Load the state of all devices from file. The RAM and the block devices of the VM are not loaded by this command.

Arguments

filename: string the file to load the state of the devices from as binary data. See xen-save-devices-state.txt for a description of the binary format.

Since

2.7

Example

```
-> { "execute": "xen-load-devices-state",
      "arguments": { "filename": "/tmp/resume" } }
<- { "return": {} }
```

xen-set-replication (Command)

Enable or disable replication.

Arguments

enable: boolean true to enable, false to disable.

primary: boolean true for primary or false for secondary.

failover: boolean (optional) true to do failover, false to stop. but cannot be specified if ‘enable’ is true. default value is false.

Returns

nothing.

Example

```
-> { "execute": "xen-set-replication",  
      "arguments": {"enable": true, "primary": false} }  
<- { "return": {} }
```

Since

2.9

If

defined(CONFIG_REPLICATION)

ReplicationStatus (Object)

The result format for 'query-xen-replication-status'.

Members

error: boolean true if an error happened, false if replication is normal.

desc: string (optional) the human readable error description string, when `error` is 'true'.

Since

2.9

If

defined(CONFIG_REPLICATION)

query-xen-replication-status (Command)

Query replication status while the vm is running.

Returns

A `ReplicationResult` object showing the status.

Example

```
-> { "execute": "query-xen-replication-status" }  
<- { "return": { "error": false } }
```

Since

2.9

If`defined(CONFIG_REPLICATION)`**xen-colo-do-checkpoint (Command)**

Xen uses this command to notify replication to trigger a checkpoint.

Returns

nothing.

Example

```
-> { "execute": "xen-colo-do-checkpoint" }
<- { "return": {} }
```

Since

2.9

If`defined(CONFIG_REPLICATION)`**COLOStatus (Object)**

The result format for ‘query-colo-status’.

Members

mode: COLOMode COLO running mode. If COLO is running, this field will return ‘primary’ or ‘secondary’.

last-mode: COLOMode COLO last running mode. If COLO is running, this field will return same like mode field, after failover we can use this field to get last colo mode. (since 4.0)

reason: COLOExitReason describes the reason for the COLO exit.

Since

3.1

query-colo-status (Command)

Query COLO status while the vm is running.

Returns

A COLOStatus object showing the status.

Example

```
-> { "execute": "query-colo-status" }
<- { "return": { "mode": "primary", "reason": "request" } }
```

Since

3.1

migrate-recover (Command)

Provide a recovery migration stream URI.

Arguments

uri: string the URI to be used for the recovery of migration stream.

Returns

nothing.

Example

```
-> { "execute": "migrate-recover",
      "arguments": { "uri": "tcp:192.168.1.200:12345" } }
<- { "return": {} }
```

Since

3.0

migrate-pause (Command)

Pause a migration. Currently it only supports postcopy.

Returns

nothing.

Example

```
-> { "execute": "migrate-pause" }
<- { "return": {} }
```

Since

3.0

UNPLUG_PRIMARY (Event)

Emitted from source side of a migration when migration state is WAIT_UNPLUG. Device was unplugged by guest operating system. Device resources in QEMU are kept on standby to be able to re-plug it in case of migration failure.

Arguments

device-id: `string` QEMU device id of the unplugged device

Since

4.2

Example

```
{ "event": "UNPLUG_PRIMARY", "data": { "device-id": "hostdev0" } }
```

DirtyRateStatus (Enum)

An enumeration of dirtyrate status.

Values

unstarted the dirtyrate thread has not been started.

measuring the dirtyrate thread is measuring.

measured the dirtyrate thread has measured and results are available.

Since

5.2

DirtyRateInfo (Object)

Information about current dirty page rate of vm.

Members

dirty-rate: int (optional) an estimate of the dirty page rate of the VM in units of MB/s, present only when estimating the rate has completed.

status: DirtyRateStatus status containing dirtyrate query status includes ‘unstarted’ or ‘measuring’ or ‘measured’

start-time: int start time in units of second for calculation

calc-time: int time in units of second for sample dirty pages

Since

5.2

calc-dirty-rate (Command)

start calculating dirty page rate for vm

Arguments

calc-time: int time in units of second for sample dirty pages

Since

5.2

Example

```
{ "command": "calc-dirty-rate", "data": { "calc-time": 1 } }
```

query-dirty-rate (Command)

query dirty page rate in units of MB/s for vm

Since

5.2

snapshot-save (Command)

Save a VM snapshot

Arguments

job-id: string identifier for the newly created job

tag: string name of the snapshot to create

vmstate: string block device node name to save vmstate to

devices: array of string list of block device node names to save a snapshot to

Applications should not assume that the snapshot save is complete when this command returns. The job commands / events must be used to determine completion and to fetch details of any errors that arise.

Note that execution of the guest CPUs may be stopped during the time it takes to save the snapshot. A future version of QEMU may ensure CPUs are executing continuously.

It is strongly recommended that `devices` contain all writable block device nodes if a consistent snapshot is required.

If `tag` already exists, an error will be reported

Returns

nothing

Example

```
-> { "execute": "snapshot-save",
    "data": {
        "job-id": "snapsave0",
        "tag": "my-snap",
        "vmstate": "disk0",
        "devices": ["disk0", "disk1"]
    }
}
<- { "return": { } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "created", "id": "snapsave0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "running", "id": "snapsave0" } }
<- { "event": "STOP" }
<- { "event": "RESUME" }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "waiting", "id": "snapsave0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "pending", "id": "snapsave0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "concluded", "id": "snapsave0" } }
-> { "execute": "query-jobs" }
<- { "return": [ { "current-progress": 1,
    "status": "concluded",
    "total-progress": 1,
    "type": "snapshot-save",
    "id": "snapsave0" } ] }
```

Since

6.0

snapshot-load (Command)

Load a VM snapshot

Arguments

job-id: string identifier for the newly created job

tag: string name of the snapshot to load.

vmstate: string block device node name to load vmstate from

devices: array of string list of block device node names to load a snapshot from

Applications should not assume that the snapshot load is complete when this command returns. The job commands / events must be used to determine completion and to fetch details of any errors that arise.

Note that execution of the guest CPUs will be stopped during the time it takes to load the snapshot.

It is strongly recommended that `devices` contain all writable block device nodes that can have changed since the original `snapshot-save` command execution.

Returns

nothing

Example

```
-> { "execute": "snapshot-load",
    "data": {
        "job-id": "snapload0",
        "tag": "my-snap",
        "vmstate": "disk0",
        "devices": ["disk0", "disk1"]
    }
}
<- { "return": { } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "created", "id": "snapload0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "running", "id": "snapload0" } }
<- { "event": "STOP" }
<- { "event": "RESUME" }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "waiting", "id": "snapload0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "pending", "id": "snapload0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "concluded", "id": "snapload0" } }
-> { "execute": "query-jobs" }
```

(continues on next page)

(continued from previous page)

```
<- { "return": [{ "current-progress": 1,
                  "status": "concluded",
                  "total-progress": 1,
                  "type": "snapshot-load",
                  "id": "snapload0"}] }
```

Since

6.0

snapshot-delete (Command)

Delete a VM snapshot

Arguments

job-id: string identifier for the newly created job

tag: string name of the snapshot to delete.

devices: array of string list of block device node names to delete a snapshot from

Applications should not assume that the snapshot delete is complete when this command returns. The job commands / events must be used to determine completion and to fetch details of any errors that arise.

Returns

nothing

Example

```
-> { "execute": "snapshot-delete",
    "data": {
        "job-id": "snapdelete0",
        "tag": "my-snap",
        "devices": ["disk0", "disk1"]
    }
}
<- { "return": { } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "created", "id": "snapdelete0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "running", "id": "snapdelete0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "waiting", "id": "snapdelete0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "pending", "id": "snapdelete0" } }
<- { "event": "JOB_STATUS_CHANGE",
    "data": { "status": "concluded", "id": "snapdelete0" } }
-> { "execute": "query-jobs" }
```

(continues on next page)

(continued from previous page)

```
<- {"return": [{"current-progress": 1,
               "status": "concluded",
               "total-progress": 1,
               "type": "snapshot-delete",
               "id": "snapdelete0"}]}
```

Since

6.0

4.8.19 Transactions

Abort (Object)

This action can be used to test transaction failure.

Since

1.6

ActionCompletionMode (Enum)

An enumeration of Transactional completion modes.

Values

individual Do not attempt to cancel any other Actions if any Actions fail after the Transaction request succeeds. All Actions that can complete successfully will do so without waiting on others. This is the default.

grouped If any Action fails after the Transaction succeeds, cancel all Actions. Actions do not complete until all Actions are ready to complete. May be rejected by Actions that do not support this completion mode.

Since

2.5

TransactionAction (Object)

A discriminated record of operations that can be performed with `transaction`. Action type can be:

- `abort`: since 1.6
- `block-dirty-bitmap-add`: since 2.5
- `block-dirty-bitmap-remove`: since 4.2
- `block-dirty-bitmap-clear`: since 2.5
- `block-dirty-bitmap-enable`: since 4.0

- `block-dirty-bitmap-disable`: since 4.0
- `block-dirty-bitmap-merge`: since 4.0
- `blockdev-backup`: since 2.3
- `blockdev-snapshot`: since 2.5
- `blockdev-snapshot-internal-sync`: since 1.7
- `blockdev-snapshot-sync`: since 1.1
- `drive-backup`: since 1.6

Members

type One of `abort`, `block-dirty-bitmap-add`, `block-dirty-bitmap-remove`, `block-dirty-bitmap-clear`, `block-dirty-bitmap-enable`, `block-dirty-bitmap-disable`, `block-dirty-bitmap-merge`, `blockdev-backup`, `blockdev-snapshot`, `blockdev-snapshot-internal-sync`, `blockdev-snapshot-sync`, `drive-backup`

data: `Abort` when type is `"abort"`

data: `BlockDirtyBitmapAdd` when type is `"block-dirty-bitmap-add"`

data: `BlockDirtyBitmap` when type is `"block-dirty-bitmap-remove"`

data: `BlockDirtyBitmap` when type is `"block-dirty-bitmap-clear"`

data: `BlockDirtyBitmap` when type is `"block-dirty-bitmap-enable"`

data: `BlockDirtyBitmap` when type is `"block-dirty-bitmap-disable"`

data: `BlockDirtyBitmapMerge` when type is `"block-dirty-bitmap-merge"`

data: `BlockdevBackup` when type is `"blockdev-backup"`

data: `BlockdevSnapshot` when type is `"blockdev-snapshot"`

data: `BlockdevSnapshotInternal` when type is `"blockdev-snapshot-internal-sync"`

data: `BlockdevSnapshotSync` when type is `"blockdev-snapshot-sync"`

data: `DriveBackup` when type is `"drive-backup"`

Since

1.1

TransactionProperties (Object)

Optional arguments to modify the behavior of a Transaction.

Members

completion-mode: `ActionCompletionMode` (optional) Controls how jobs launched asynchronously by Actions will complete or fail as a group. See `ActionCompletionMode` for details.

Since

2.5

`transaction` (Command)

Executes a number of transactionable QMP commands atomically. If any operation fails, then the entire set of actions will be abandoned and the appropriate error returned.

For external snapshots, the dictionary contains the device, the file to use for the new snapshot, and the format. The default format, if not specified, is `qcow2`.

Each new snapshot defaults to being created by QEMU (wiping any contents if the file already exists), but it is also possible to reuse an externally-created file. In the latter case, you should ensure that the new image file has the same contents as the current one; QEMU cannot perform any meaningful check. Typically this is achieved by using the current image file as the backing file for the new image.

On failure, the original disks pre-snapshot attempt will be used.

For internal snapshots, the dictionary contains the device and the snapshot's name. If an internal snapshot matching name already exists, the request will be rejected. Only some image formats support it, for example, `qcow2`, `rbd`, and `sheepdog`.

On failure, qemu will try delete the newly created internal snapshot in the transaction. When an I/O error occurs during deletion, the user needs to fix it later with `qemu-img` or other command.

Arguments

actions: array of TransactionAction List of `TransactionAction`; information needed for the respective operations.

properties: TransactionProperties (optional) structure of additional options to control the execution of the transaction. See `TransactionProperties` for additional detail.

Returns

nothing on success

Errors depend on the operations of the transaction

Note

The transaction aborts on the first failure. Therefore, there will be information on only one failed operation returned in an error condition, and subsequent actions will not have been attempted.

Since

1.1

Example

```
-> { "execute": "transaction",
    "arguments": { "actions": [
        { "type": "blockdev-snapshot-sync", "data" : { "device": "ide-hd0",
            "snapshot-file": "/some/place/my-image",
            "format": "qcow2" } },
        { "type": "blockdev-snapshot-sync", "data" : { "node-name": "myfile",
            "snapshot-file": "/some/place/my-image2",
            "snapshot-node-name": "node3432",
            "mode": "existing",
            "format": "qcow2" } },
        { "type": "blockdev-snapshot-sync", "data" : { "device": "ide-hd1",
            "snapshot-file": "/some/place/my-image2",
            "mode": "existing",
            "format": "qcow2" } },
        { "type": "blockdev-snapshot-internal-sync", "data" : {
            "device": "ide-hd2",
            "name": "snapshot0" } } ] } }
<- { "return": {} }
```

4.8.20 Tracing

TraceEventState (Enum)

State of a tracing event.

Values

unavailable The event is statically disabled.

disabled The event is dynamically disabled.

enabled The event is dynamically enabled.

Since

2.2

TraceEventInfo (Object)

Information of a tracing event.

Members

name: string Event name.

state: TraceEventState Tracing state.

vcpu: boolean Whether this is a per-vCPU event (since 2.7).

An event is per-vCPU if it has the “vcpu” property in the “trace-events” files.

Since

2.2

`trace-event-get-state` (Command)

Query the state of events.

Arguments

name: string Event name pattern (case-sensitive glob).

vcpu: int (optional) The vCPU to query (any by default; since 2.7).

Returns

a list of `TraceEventInfo` for the matching events

An event is returned if:

- its name matches the `name` pattern, and
- if `vcpu` is given, the event has the “`vcpu`” property.

Therefore, if `vcpu` is given, the operation will only match per-vCPU events, returning their state on the specified vCPU. Special case: if `name` is an exact match, `vcpu` is given and the event does not have the “`vcpu`” property, an error is returned.

Since

2.2

Example

```
-> { "execute": "trace-event-get-state",  
    "arguments": { "name": "qemu_memalign" } }  
<- { "return": [ { "name": "qemu_memalign", "state": "disabled" } ] }
```

`trace-event-set-state` (Command)

Set the dynamic tracing state of events.

Arguments

name: string Event name pattern (case-sensitive glob).

enable: boolean Whether to enable tracing.

ignore-unavailable: boolean (optional) Do not match unavailable events with `name`.

vcpu: int (optional) The vCPU to act upon (all by default; since 2.7).

An event's state is modified if: - its name matches the `name` pattern, and - if `vcpu` is given, the event has the "vcpu" property.

Therefore, if `vcpu` is given, the operation will only match per-vCPU events, setting their state on the specified vCPU. Special case: if `name` is an exact match, `vcpu` is given and the event does not have the "vcpu" property, an error is returned.

Since

2.2

Example

```
-> { "execute": "trace-event-set-state",
      "arguments": { "name": "qemu_memalign", "enable": "true" } }
<- { "return": {} }
```

4.8.21 QMP monitor control

`qmp_capabilities` (Command)

Enable QMP capabilities.

Arguments:

Arguments

enable: array of QMPCapability (optional) An optional list of QMPCapability values to enable. The client must not enable any capability that is not mentioned in the QMP greeting message. If the field is not provided, it means no QMP capabilities will be enabled. (since 2.12)

Example

```
-> { "execute": "qmp_capabilities",
      "arguments": { "enable": [ "oob" ] } }
<- { "return": {} }
```

Notes

This command is valid exactly when first connecting: it must be issued before any other command will be accepted, and will fail once the monitor is accepting other commands. (see `qemu docs/interop/qmp-spec.txt`)

The QMP client needs to explicitly enable QMP capabilities, otherwise all the QMP capabilities will be turned off by default.

Since

0.13

QMPCapability (Enum)

Enumeration of capabilities to be advertised during initial client connection, used for agreeing on particular QMP extension behaviors.

Values

oob QMP ability to support out-of-band requests. (Please refer to qmp-spec.txt for more information on OOB)

Since

2.12

VersionTriple (Object)

A three-part version number.

Members

major: int The major version number.

minor: int The minor version number.

micro: int The micro version number.

Since

2.4

VersionInfo (Object)

A description of QEMU's version.

Members

qemu: VersionTriple The version of QEMU. By current convention, a micro version of 50 signifies a development branch. A micro version greater than or equal to 90 signifies a release candidate for the next minor version. A micro version of less than 50 signifies a stable release.

package: string QEMU will always set this field to an empty string. Downstream versions of QEMU should set this to a non-empty string. The exact format depends on the downstream however it is highly recommended that a unique name is used.

Since

0.14

query-version (Command)

Returns the current version of QEMU.

Returns

A `VersionInfo` object describing the current version of QEMU.

Since

0.14

Example

```
-> { "execute": "query-version" }
<- {
  "return": {
    "qemu": {
      "major": 0,
      "minor": 11,
      "micro": 5
    },
    "package": ""
  }
}
```

CommandInfo (Object)

Information about a QMP command

Members

name: string The command name

Since

0.14

query-commands (Command)

Return a list of supported QMP commands by this server

Returns

A list of `CommandInfo` for all supported commands

Since

0.14

Example

```
-> { "execute": "query-commands" }
<- {
  "return": [
    {
      "name": "query-balloon"
    },
    {
      "name": "system_powerdown"
    }
  ]
}
```

Note

This example has been shortened as the real response is too long.

EventInfo (Object)

Information about a QMP event

Members

name: string The event name

Since

1.2

query-events (Command)

Return information on QMP events.

Features

deprecated This command is deprecated, because its output doesn't reflect compile-time configuration. Use 'query-qmp-schema' instead.

Returns

A list of EventInfo.

Since

1.2

Example

```
-> { "execute": "query-events" }
<- {
  "return": [
    {
      "name": "SHUTDOWN"
    },
    {
      "name": "RESET"
    }
  ]
}
```

Note

This example has been shortened as the real response is too long.

quit (Command)

This command will cause the QEMU process to exit gracefully. While every attempt is made to send the QMP response before terminating, this is not guaranteed. When using this interface, a premature EOF would not be unexpected.

Since

0.14

Example

```
-> { "execute": "quit" }
<- { "return": {} }
```

MonitorMode (Enum)

An enumeration of monitor modes.

Values

readline HMP monitor (human-oriented command line interface)

control QMP monitor (JSON-based machine interface)

Since

5.0

MonitorOptions (Object)

Options to be used for adding a new monitor.

Members

id: **string** (optional) Name of the monitor

mode: **MonitorMode** (optional) Selects the monitor mode (default: readline in the system emulator, control in qemu-storage-daemon)

pretty: **boolean** (optional) Enables pretty printing (QMP only)

chardev: **string** Name of a character device to expose the monitor on

Since

5.0

4.8.22 QMP introspection

query-qmp-schema (Command)

Command query-qmp-schema exposes the QMP wire ABI as an array of SchemaInfo. This lets QMP clients figure out what commands and events are available in this QEMU, and their parameters and results.

However, the SchemaInfo can't reflect all the rules and restrictions that apply to QMP. It's interface introspection (figuring out what's there), not interface specification. The specification is in the QAPI schema.

Furthermore, while we strive to keep the QMP wire format backwards-compatible across qemu versions, the introspection output is not guaranteed to have the same stability. For example, one version of qemu may list an object member as an optional non-variant, while another lists the same member only through the object's variants; or the type of a member may change from a generic string into a specific enum or from one specific type into an alternate that includes the original type alongside something else.

Returns

array of SchemaInfo, where each element describes an entity in the ABI: command, event, type, ...

The order of the various SchemaInfo is unspecified; however, all names are guaranteed to be unique (no name will be duplicated with different meta-types).

Note

the QAPI schema is also used to help define *internal* interfaces, by defining QAPI types. These are not part of the QMP wire ABI, and therefore not returned by this command.

Since

2.5

SchemaMetaType (Enum)

This is a `SchemaInfo`'s meta type, i.e. the kind of entity it describes.

Values

builtin a predefined type such as 'int' or 'bool'.

enum an enumeration type

array an array type

object an object type (struct or union)

alternate an alternate type

command a QMP command

event a QMP event

Since

2.5

SchemaInfo (Object)

Members

name: string the entity's name, inherited from `base`. The `SchemaInfo` is always referenced by this name. Commands and events have the name defined in the QAPI schema. Unlike command and event names, type names are not part of the wire ABI. Consequently, type names are meaningless strings here, although they are still guaranteed unique regardless of `meta-type`.

meta-type: SchemaMetaType the entity's meta type, inherited from `base`.

features: array of string (optional) names of features associated with the entity, in no particular order. (since 4.1 for object types, 4.2 for commands, 5.0 for the rest)

The members of `SchemaInfoBuiltin` when `meta-type` is "builtin"

The members of `SchemaInfoEnum` when `meta-type` is "enum"

The members of `SchemaInfoArray` when `meta-type` is "array"

The members of `SchemaInfoObject` when `meta-type` is "object"

The members of `SchemaInfoAlternate` when `meta-type` is "alternate"

The members of `SchemaInfoCommand` when `meta-type` is "command"

The members of `SchemaInfoEvent` when `meta-type` is "event"

Additional members depend on the value of `meta-type`.

Since

2.5

SchemaInfoBuiltin (Object)

Additional SchemaInfo members for meta-type ‘builtin’.

Members

json-type: JSONType the JSON type used for this type on the wire.

Since

2.5

JSONType (Enum)

The four primitive and two structured types according to RFC 8259 section 1, plus ‘int’ (split off ‘number’), plus the obvious top type ‘value’.

Values

string Not documented

number Not documented

int Not documented

boolean Not documented

null Not documented

object Not documented

array Not documented

value Not documented

Since

2.5

SchemaInfoEnum (Object)

Additional SchemaInfo members for meta-type ‘enum’.

Members

values: array of string the enumeration type's values, in no particular order.

Values of this type are JSON string on the wire.

Since

2.5

SchemaInfoArray (Object)

Additional SchemaInfo members for meta-type 'array'.

Members

element-type: string the array type's element type.

Values of this type are JSON array on the wire.

Since

2.5

SchemaInfoObject (Object)

Additional SchemaInfo members for meta-type 'object'.

Members

members: array of SchemaInfoObjectMember the object type's (non-variant) members, in no particular order.

tag: string (optional) the name of the member serving as type tag. An element of `members` with this name must exist.

variants: array of SchemaInfoObjectVariant (optional) variant members, i.e. additional members that depend on the type tag's value. Present exactly when `tag` is present. The variants are in no particular order, and may even differ from the order of the values of the enum type of the `tag`.

Values of this type are JSON object on the wire.

Since

2.5

SchemaInfoObjectMember (Object)

An object member.

Members

name: string the member's name, as defined in the QAPI schema.

type: string the name of the member's type.

default: value (optional) default when used as command parameter. If absent, the parameter is mandatory. If present, the value must be null. The parameter is optional, and behavior when it's missing is not specified here. Future extension: if present and non-null, the parameter is optional, and defaults to this value.

features: array of string (optional) names of features associated with the member, in no particular order. (since 5.0)

Since

2.5

SchemaInfoObjectVariant (Object)

The variant members for a value of the type tag.

Members

case: string a value of the type tag.

type: string the name of the object type that provides the variant members when the type tag has value `case`.

Since

2.5

SchemaInfoAlternate (Object)

Additional SchemaInfo members for meta-type 'alternate'.

Members

members: array of SchemaInfoAlternateMember the alternate type's members, in no particular order. The members' wire encoding is distinct, see docs/devel/qapi-code-gen.txt section Alternate types.

On the wire, this can be any of the members.

Since

2.5

SchemaInfoAlternateMember (Object)

An alternate member.

Members

type: string the name of the member's type.

Since

2.5

SchemaInfoCommand (Object)

Additional SchemaInfo members for meta-type 'command'.

Members

arg-type: string the name of the object type that provides the command's parameters.

ret-type: string the name of the command's result type.

allow-oob: boolean (optional) whether the command allows out-of-band execution, defaults to false (Since: 2.12)

TODO

`success-response` (currently irrelevant, because it's QGA, not QMP)

Since

2.5

SchemaInfoEvent (Object)

Additional SchemaInfo members for meta-type 'event'.

Members

arg-type: string the name of the object type that provides the event's parameters.

Since

2.5

4.8.23 QEMU Object Model (QOM)

ObjectPropertyInfo (Object)

Members

name: string the name of the property

type: string the type of the property. This will typically come in one of four forms:

- 1) A primitive type such as 'u8', 'u16', 'bool', 'str', or 'double'. These types are mapped to the appropriate JSON type.
- 2) A child type in the form 'child<subtype>' where subtype is a qdev device type name. Child properties create the composition tree.
- 3) A link type in the form 'link<subtype>' where subtype is a qdev device type name. Link properties form the device model graph.

description: string (optional) if specified, the description of the property.

default-value: value (optional) the default value, if any (since 5.0)

Since

1.2

qom-list (Command)

This command will list any properties of an object given a path in the object model.

Arguments

path: string the path within the object model. See `qom-get` for a description of this parameter.

Returns

a list of `ObjectPropertyInfo` that describe the properties of the object.

Since

1.2

Example

```
-> { "execute": "qom-list",
      "arguments": { "path": "/chardevs" } }
<- { "return": [ { "name": "type", "type": "string" },
                  { "name": "parallel0", "type": "child<chardev-vc>" },
                  { "name": "serial0", "type": "child<chardev-vc>" },
                  { "name": "mon0", "type": "child<chardev-stdio>" } ] }
```

qom-get (Command)

This command will get a property from a object model path and return the value.

Arguments

path: string The path within the object model. There are two forms of supported paths—absolute and partial paths.

Absolute paths are derived from the root object and can follow `child<>` or `link<>` properties. Since they can follow `link<>` properties, they can be arbitrarily long. Absolute paths look like absolute filenames and are prefixed with a leading slash.

Partial paths look like relative filenames. They do not begin with a prefix. The matching rules for partial paths are subtle but designed to make specifying objects easy. At each level of the composition tree, the partial path is matched as an absolute path. The first match is not returned. At least two matches are searched for. A successful result is only returned if only one match is found. If more than one match is found, a flag is return to indicate that the match was ambiguous.

property: string The property name to read

Returns

The property value. The type depends on the property type. `child<>` and `link<>` properties are returned as `#str` pathnames. All integer property types (`u8`, `u16`, etc) are returned as `#int`.

Since

1.2

Example

```

1. Use absolute path
-> { "execute": "qom-get",
    "arguments": { "path": "/machine/unattached/device[0]",
                  "property": "hotplugged" } }
<- { "return": false }

2. Use partial path
-> { "execute": "qom-get",
    "arguments": { "path": "unattached/sysbus",
                  "property": "type" } }
<- { "return": "System" }

```

qom-set (Command)

This command will set a property from a object model path.

Arguments

path: string see `qom-get` for a description of this parameter

property: string the property name to set

value: value a value whose type is appropriate for the property type. See `qom-get` for a description of type mapping.

Since

1.2

Example

```
-> { "execute": "qom-set",  
      "arguments": { "path": "/machine",  
                     "property": "graphics",  
                     "value": false } }  
<- { "return": {} }
```

ObjectTypeInfo (Object)

This structure describes a search result from `qom-list-types`

Members

name: string the type name found in the search

abstract: boolean (optional) the type is abstract and can't be directly instantiated. Omitted if false. (since 2.10)

parent: string (optional) Name of parent type, if any (since 2.10)

Since

1.1

qom-list-types (Command)

This command will return a list of types given search parameters

Arguments

implements: string (optional) if specified, only return types that implement this type name

abstract: boolean (optional) if true, include abstract types in the results

Returns

a list of `ObjectTypeInfo` or an empty list if no results are found

Since

1.1

`qom-list-properties` (Command)

List properties associated with a QOM object.

Arguments

typename: `string` the type name of an object

Note

objects can create properties at runtime, for example to describe links between different devices and/or objects. These properties are not included in the output of this command.

Returns

a list of `ObjectPropertyInfo` describing object properties

Since

2.12

`object-add` (Command)

Create a QOM object.

Arguments

qom-type: `string` the class name for the object to be created

id: `string` the name of the new object

props: `value` (optional) a dictionary of properties to be passed to the backend. Deprecated since 5.0, specify the properties on the top level instead. It is an error to specify the same option both on the top level and in `props`.

Additional arguments depend on `qom-type` and are passed to the backend unchanged.

Returns

Nothing on success Error if `qom-type` is not a valid class name

Since

2.0

Example

```
-> { "execute": "object-add",  
      "arguments": { "qom-type": "rng-random", "id": "rng1",  
                     "filename": "/dev/hwrng" } }  
<- { "return": {} }
```

`object-del` (Command)

Remove a QOM object.

Arguments

id: string the name of the QOM object to remove

Returns

Nothing on success Error if `id` is not a valid id for a QOM object

Since

2.0

Example

```
-> { "execute": "object-del", "arguments": { "id": "rng1" } }  
<- { "return": {} }
```

4.8.24 Device infrastructure (qdev)

`device-list-properties` (Command)

List properties associated with a device.

Arguments

typename: string the type name of a device

Returns

a list of `ObjectPropertyInfo` describing a devices properties

Note

objects can create properties at runtime, for example to describe links between different devices and/or objects. These properties are not included in the output of this command.

Since

1.2

`device_add` (Command)

Arguments

driver: string the name of the new device's driver

bus: string (optional) the device's parent bus (device tree path)

id: string (optional) the device's ID, must be unique

Additional arguments depend on the type.

Add a device.

Notes

1. For detailed information about this command, please refer to the 'docs/qdev-device-use.txt' file.
2. It's possible to list device properties by running QEMU with the "--device DEVICE,help" command-line argument, where DEVICE is the device's name

Example

```
-> { "execute": "device_add",
    "arguments": { "driver": "e1000", "id": "net1",
                  "bus": "pci.0",
                  "mac": "52:54:00:12:34:56" } }
<- { "return": {} }
```

TODO

This command effectively bypasses QAPI completely due to its “additional arguments” business. It shouldn't have been added to the schema in this form. It should be qapified properly, or replaced by a properly qapified command.

Since

0.13

`device_del` (Command)

Remove a device from a guest

Arguments

id: `string` the device's ID or QOM path

Returns

Nothing on success If `id` is not a valid device, `DeviceNotFound`

Notes

When this command completes, the device may not be removed from the guest. Hot removal is an operation that requires guest cooperation. This command merely requests that the guest begin the hot removal process. Completion of the device removal process is signaled with a `DEVICE_DELETED` event. Guest reset will automatically complete removal for all devices.

Since

0.14

Example

```
-> { "execute": "device_del",  
      "arguments": { "id": "net1" } }  
<- { "return": {} }  
  
-> { "execute": "device_del",  
      "arguments": { "id": "/machine/peripheral-anon/device[0]" } }  
<- { "return": {} }
```

`DEVICE_DELETED` (Event)

Emitted whenever the device removal completion is acknowledged by the guest. At this point, it's safe to reuse the specified device ID. Device removal can be initiated by the guest or by HMP/QMP commands.

Arguments

device: `string (optional)` device name

path: `string` device path

Since

1.5

Example

```
<- { "event": "DEVICE_DELETED",
      "data": { "device": "virtio-net-pci-0",
                 "path": "/machine/peripheral/virtio-net-pci-0" },
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

4.8.25 Machines

SysEmuTarget (Enum)

The comprehensive enumeration of QEMU system emulation (“softmmu”) targets. Run “./configure –help” in the project root directory, and look for the *-softmmu targets near the “–target-list” option. The individual target constants are not documented here, for the time being.

Values

rx since 5.0

avr since 5.1

aarch64 Not documented

alpha Not documented

arm Not documented

cris Not documented

hppa Not documented

i386 Not documented

lm32 Not documented

m68k Not documented

microblaze Not documented

microblazeel Not documented

mips Not documented

mips64 Not documented

mips64el Not documented

mipsel Not documented

moxie Not documented

nios2 Not documented

or1k Not documented

ppc Not documented

ppc64 Not documented
riscv32 Not documented
riscv64 Not documented
s390x Not documented
sh4 Not documented
sh4eb Not documented
sparc Not documented
sparc64 Not documented
tricore Not documented
unicore32 Not documented
x86_64 Not documented
xtensa Not documented
xtensaeb Not documented

Notes

The resulting QMP strings can be appended to the “qemu-system-” prefix to produce the corresponding QEMU executable name. This is true even for “qemu-system-x86_64”.

Since

3.0

CpuInfoArch (Enum)

An enumeration of cpu types that enable additional information during `query-cpus` and `query-cpus-fast`.

Values

s390 since 2.12
riscv since 2.12
x86 Not documented
sparc Not documented
ppc Not documented
mips Not documented
tricore Not documented
other Not documented

Since

2.6

CpuInfo (Object)

Information about a virtual CPU

Members

cpu: **int** the index of the virtual CPU

current: **boolean** this only exists for backwards compatibility and should be ignored

halted: **boolean** true if the virtual CPU is in the halt state. Halt usually refers to a processor specific low power mode.

qom_path: **string** path to the CPU object in the QOM tree (since 2.4)

thread_id: **int** ID of the underlying host thread

props: **CpuInstanceProperties (optional)** properties describing to which node/socket/core/thread virtual CPU belongs to, provided if supported by board (since 2.10)

arch: **CpuInfoArch** architecture of the cpu, which determines which additional fields will be listed (since 2.6)

The members of **CpuInfoX86** when **arch** is "**x86**"

The members of **CpuInfoSPARC** when **arch** is "**sparc**"

The members of **CpuInfoPPC** when **arch** is "**ppc**"

The members of **CpuInfoMIPS** when **arch** is "**mips**"

The members of **CpuInfoTricore** when **arch** is "**tricore**"

The members of **CpuInfoS390** when **arch** is "**s390**"

The members of **CpuInfoRISCV** when **arch** is "**riscv**"

Since

0.14

Notes

halted is a transient state that changes frequently. By the time the data is sent to the client, the guest may no longer be halted.

CpuInfoX86 (Object)

Additional information about a virtual i386 or x86_64 CPU

Members

pc: int the 64-bit instruction pointer

Since

2.6

CpuInfoSPARC (Object)

Additional information about a virtual SPARC CPU

Members

pc: int the PC component of the instruction pointer

npc: int the NPC component of the instruction pointer

Since

2.6

CpuInfoPPC (Object)

Additional information about a virtual PPC CPU

Members

nip: int the instruction pointer

Since

2.6

CpuInfoMIPS (Object)

Additional information about a virtual MIPS CPU

Members

PC: int the instruction pointer

Since

2.6

CpuInfoTricore (Object)

Additional information about a virtual Tricore CPU

Members

PC: int the instruction pointer

Since

2.6

CpuInfoRISCV (Object)

Additional information about a virtual RISCV CPU

Members

pc: int the instruction pointer

Since 2.12

CpuS390State (Enum)

An enumeration of cpu states that can be assumed by a virtual S390 CPU

Values

uninitialized Not documented

stopped Not documented

check-stop Not documented

operating Not documented

load Not documented

Since

2.12

CpuInfoS390 (Object)

Additional information about a virtual S390 CPU

Members

cpu-state: CpuS390State the virtual CPU's state

Since

2.12

query-cpus (Command)

Returns a list of information about each virtual CPU.

This command causes vCPU threads to exit to userspace, which causes a small interruption to guest CPU execution. This will have a negative impact on realtime guests and other latency sensitive guest workloads.

Features

deprecated This command is deprecated, because it interferes with the guest. Use ‘query-cpus-fast’ instead to avoid the vCPU interruption.

Returns

a list of CpuInfo for each virtual CPU

Since

0.14

Example

```
-> { "execute": "query-cpus" }
<- { "return": [
  {
    "CPU":0,
    "current":true,
    "halted":false,
    "qom_path":"/machine/unattached/device[0]",
    "arch":"x86",
    "pc":3227107138,
    "thread_id":3134
  },
  {
    "CPU":1,
    "current":false,
    "halted":true,
    "qom_path":"/machine/unattached/device[2]",
    "arch":"x86",
    "pc":7108165,
    "thread_id":3135
  }
]
```

CpuInfoFast (Object)

Information about a virtual CPU

Members

cpu-index: `int` index of the virtual CPU

qom-path: `string` path to the CPU object in the QOM tree

thread-id: `int` ID of the underlying host thread

props: `CpuInstanceProperties` (optional) properties describing to which node/socket/core/thread virtual CPU belongs to, provided if supported by board

arch: `CpuInfoArch` base architecture of the cpu

target: `SysEmuTarget` the QEMU system emulation target, which determines which additional fields will be listed (since 3.0)

The members of `CpuInfoS390` when `target` is "s390x"

Features

deprecated Member `arch` is deprecated. Use `target` instead.

Since

2.12

query-cpus-fast (Command)

Returns information about all virtual CPUs. This command does not incur a performance penalty and should be used in production instead of `query-cpus`.

Returns

list of `CpuInfoFast`

Since

2.12

Example

```
-> { "execute": "query-cpus-fast" }
<- { "return": [
    {
        "thread-id": 25627,
```

(continues on next page)

(continued from previous page)

```

        "props": {
            "core-id": 0,
            "thread-id": 0,
            "socket-id": 0
        },
        "qom-path": "/machine/unattached/device[0]",
        "arch": "x86",
        "target": "x86_64",
        "cpu-index": 0
    },
    {
        "thread-id": 25628,
        "props": {
            "core-id": 0,
            "thread-id": 0,
            "socket-id": 1
        },
        "qom-path": "/machine/unattached/device[2]",
        "arch": "x86",
        "target": "x86_64",
        "cpu-index": 1
    }
]
}

```

MachineInfo (Object)

Information describing a machine.

Members

name: string the name of the machine

alias: string (optional) an alias for the machine name

is-default: boolean (optional) whether the machine is default

cpu-max: int maximum number of CPUs supported by the machine type (since 1.5)

hotpluggable-cpus: boolean cpu hotplug via -device is supported (since 2.7)

numa-mem-supported: boolean true if '-numa node,mem' option is supported by the machine type and false otherwise (since 4.1)

deprecated: boolean if true, the machine type is deprecated and may be removed in future versions of QEMU according to the QEMU deprecation policy (since 4.1)

default-cpu-type: string (optional) default CPU model typename if none is requested via the -cpu argument. (since 4.2)

default-ram-id: string (optional) the default ID of initial RAM memory backend (since 5.2)

Since

1.2

query-machines (Command)

Return a list of supported machines

Returns

a list of MachineInfo

Since

1.2

CurrentMachineParams (Object)

Information describing the running machine parameters.

Members

wakeup-suspend-support: boolean true if the machine supports wake up from suspend

Since

4.0

query-current-machine (Command)

Return information on the current virtual machine.

Returns

CurrentMachineParams

Since

4.0

TargetInfo (Object)

Information describing the QEMU target.

Members

arch: SysEmuTarget the target architecture

Since

1.2

query-target (Command)

Return information about the target for this QEMU

Returns

TargetInfo

Since

1.2

UuidInfo (Object)

Guest UUID information (Universally Unique Identifier).

Members

UUID: string the UUID of the guest

Since

0.14

Notes

If no UUID was specified for the guest, a null UUID is returned.

query-uuid (Command)

Query the guest UUID information.

Returns

The `UuidInfo` for the guest

Since

0.14

Example

```
-> { "execute": "query-uuid" }  
<- { "return": { "UUID": "550e8400-e29b-41d4-a716-446655440000" } }
```

GuidInfo (Object)

GUID information.

Members

guid: `string` the globally unique identifier

Since

2.9

query-vm-generation-id (Command)

Show Virtual Machine Generation ID

Since

2.9

system_reset (Command)

Performs a hard reset of a guest.

Since

0.14

Example

```
-> { "execute": "system_reset" }  
<- { "return": {} }
```

system_powerdown (Command)

Requests that a guest perform a powerdown operation.

Since

0.14

Notes

A guest may or may not respond to this command. This command returning does not indicate that a guest has accepted the request or that it has shut down. Many guests will respond to this command by prompting the user in some way.

Example

```
-> { "execute": "system_powerdown" }
<- { "return": {} }
```

system_wakeup (Command)

Wake up guest from suspend. If the guest has wake-up from suspend support enabled (wakeup-suspend-support flag from query-current-machine), wake-up guest from suspend if the guest is in SUSPENDED state. Return an error otherwise.

Since

1.1

Returns

nothing.

Note

prior to 4.0, this command does nothing in case the guest isn't suspended.

Example

```
-> { "execute": "system_wakeup" }
<- { "return": {} }
```

LostTickPolicy (Enum)

Policy for handling lost ticks in timer devices. Ticks end up getting lost when, for example, the guest is paused.

Values

discard throw away the missed ticks and continue with future injection normally. The guest OS will see the timer jump ahead by a potentially quite significant amount all at once, as if the intervening chunk of time had simply not existed; needless to say, such a sudden jump can easily confuse a guest OS which is not specifically prepared to deal with it. Assuming the guest OS can deal correctly with the time jump, the time in the guest and in the host should now match.

delay continue to deliver ticks at the normal rate. The guest OS will not notice anything is amiss, as from its point of view time will have continued to flow normally. The time in the guest should now be behind the time in the host by exactly the amount of time during which ticks have been missed.

slew deliver ticks at a higher rate to catch up with the missed ticks. The guest OS will not notice anything is amiss, as from its point of view time will have continued to flow normally. Once the timer has managed to catch up with all the missing ticks, the time in the guest and in the host should match.

Since

2.0

`inject-nmi` (Command)

Injects a Non-Maskable Interrupt into the default CPU (x86/s390) or all CPUs (ppc64). The command fails when the guest doesn't support injecting.

Returns

If successful, nothing

Since

0.14

Note

prior to 2.1, this command was only supported for x86 and s390 VMs

Example

```
-> { "execute": "inject-nmi" }
<- { "return": {} }
```

`KvmInfo` (Object)

Information about support for KVM acceleration

Members

enabled: boolean true if KVM acceleration is active

present: boolean true if KVM acceleration is built into this executable

Since

0.14

query-kvm (Command)

Returns information about KVM acceleration

Returns

KvmInfo

Since

0.14

Example

```
-> { "execute": "query-kvm" }  
<- { "return": { "enabled": true, "present": true } }
```

NumaOptionsType (Enum)

Values

node NUMA nodes configuration

dist NUMA distance configuration (since 2.10)

cpu property based CPU(s) to node mapping (Since: 2.10)

hmat-lb memory latency and bandwidth information (Since: 5.0)

hmat-cache memory side cache information (Since: 5.0)

Since

2.1

NumaOptions (Object)

A discriminated record of NUMA options. (for OptsVisitor)

Members

type: `NumaOptionsType` Not documented

The members of `NumaNodeOptions` when `type` is "node"

The members of `NumaDistOptions` when `type` is "dist"

The members of `NumaCpuOptions` when `type` is "cpu"

The members of `NumaHmatLBOptions` when `type` is "hmat-lb"

The members of `NumaHmatCacheOptions` when `type` is "hmat-cache"

Since

2.1

NumaNodeOptions (Object)

Create a guest NUMA node. (for `OptsVisitor`)

Members

nodeid: `int` (optional) NUMA node ID (increase by 1 from 0 if omitted)

cpus: `array of int` (optional)

VCPUs belonging to this node (assign VCPUS round-robin if omitted)

mem: `int` (optional) memory size of this node; mutually exclusive with `memdev`. Equally divide total memory among nodes if both `mem` and `memdev` are omitted.

memdev: `string` (optional) memory backend object. If specified for one node, it must be specified for all nodes.

initiator: `int` (optional) defined in ACPI 6.3 Chapter 5.2.27.3 Table 5-145, points to the `nodeid` which has the memory controller responsible for this NUMA node. This field provides additional information as to the initiator node that is closest (as in directly attached) to this node, and therefore has the best performance (since 5.0)

Since

2.1

NumaDistOptions (Object)

Set the distance between 2 NUMA nodes.

Members

src: `int` source NUMA node.

dst: `int` destination NUMA node.

val: `int` NUMA distance from source node to destination node. When a node is unreachable from another node, set the distance between them to 255.

Since

2.10

X86CPURegister32 (Enum)

A X86 32-bit register

Values

EAX Not documented

EBX Not documented

ECX Not documented

EDX Not documented

ESP Not documented

EBP Not documented

ESI Not documented

EDI Not documented

Since

1.5

X86CPUFeatureWordInfo (Object)

Information about a X86 CPU feature word

Members

cpuid-input-eax: int Input EAX value for CPUID instruction for that feature word

cpuid-input-ecx: int (optional) Input ECX value for CPUID instruction for that feature word

cpuid-register: X86CPURegister32 Output register containing the feature bits

features: int value of output register, containing the feature bits

Since

1.5

DummyForceArrays (Object)

Not used by QMP; hack to let us use X86CPUFeatureWordInfoList internally

Members

unused: array of **X86CPUFeatureWordInfo** Not documented

Since

2.5

NumaCpuOptions (Object)

Option “-numa cpu” overrides default cpu to node mapping. It accepts the same set of cpu properties as returned by query-hotpluggable-cpus[].props, where node-id could be used to override default node mapping.

Members

The members of **CpuInstanceProperties**

Since

2.10

HmatLBMemoryHierarchy (Enum)

The memory hierarchy in the System Locality Latency and Bandwidth Information Structure of HMAT (Heterogeneous Memory Attribute Table)

For more information about **HmatLBMemoryHierarchy**, see chapter 5.2.27.4: Table 5-146: Field “Flags” of ACPI 6.3 spec.

Values

memory the structure represents the memory performance

first-level first level of memory side cache

second-level second level of memory side cache

third-level third level of memory side cache

Since

5.0

HmatLBDataType (Enum)

Data type in the System Locality Latency and Bandwidth Information Structure of HMAT (Heterogeneous Memory Attribute Table)

For more information about **HmatLBDataType**, see chapter 5.2.27.4: Table 5-146: Field “Data Type” of ACPI 6.3 spec.

Values

access-latency access latency (nanoseconds)
read-latency read latency (nanoseconds)
write-latency write latency (nanoseconds)
access-bandwidth access bandwidth (Bytes per second)
read-bandwidth read bandwidth (Bytes per second)
write-bandwidth write bandwidth (Bytes per second)

Since

5.0

NumaHmatLBOptions (Object)

Set the system locality latency and bandwidth information between Initiator and Target proximity Domains.
For more information about NumaHmatLBOptions, see chapter 5.2.27.4: Table 5-146 of ACPI 6.3 spec.

Members

initiator: int the Initiator Proximity Domain.
target: int the Target Proximity Domain.
hierarchy: HmatLBMemoryHierarchy the Memory Hierarchy. Indicates the performance of memory or side cache.
data-type: HmatLBDataType presents the type of data, access/read/write latency or hit latency.
latency: int (optional) the value of latency from `initiator` to `target` proximity domain, the latency unit is “ns(nanosecond)”.
bandwidth: int (optional) the value of bandwidth between `initiator` and `target` proximity domain, the bandwidth unit is “Bytes per second”.

Since

5.0

HmatCacheAssociativity (Enum)

Cache associativity in the Memory Side Cache Information Structure of HMAT
For more information of HmatCacheAssociativity, see chapter 5.2.27.5: Table 5-147 of ACPI 6.3 spec.

Values

none

None (no memory side cache in this proximity domain, or cache associativity unknown)

direct Direct Mapped

complex Complex Cache Indexing (implementation specific)

Since

5.0

HmatCacheWritePolicy (Enum)

Cache write policy in the Memory Side Cache Information Structure of HMAT

For more information of HmatCacheWritePolicy, see chapter 5.2.27.5: Table 5-147: Field “Cache Attributes” of ACPI 6.3 spec.

Values

none None (no memory side cache in this proximity domain, or cache write policy unknown)

write-back Write Back (WB)

write-through Write Through (WT)

Since

5.0

NumaHmatCacheOptions (Object)

Set the memory side cache information for a given memory domain.

For more information of NumaHmatCacheOptions, see chapter 5.2.27.5: Table 5-147: Field “Cache Attributes” of ACPI 6.3 spec.

Members

node-id: int the memory proximity domain to which the memory belongs.

size: int the size of memory side cache in bytes.

level: int the cache level described in this structure.

associativity: HmatCacheAssociativity the cache associativity, none/direct-mapped/complex(complex cache indexing).

policy: HmatCacheWritePolicy the write policy, none/write-back/write-through.

line: int the cache Line size in bytes.

Since

5.0

HostMemPolicy (Enum)

Host memory policy types

Values

default restore default policy, remove any nondefault policy

preferred set the preferred host nodes for allocation

bind a strict policy that restricts memory allocation to the host nodes specified

interleave memory allocations are interleaved across the set of host nodes specified

Since

2.1

memsave (Command)

Save a portion of guest memory to a file.

Arguments

val: int the virtual address of the guest to start from

size: int the size of memory region to save

filename: string the file to save the memory to as binary data

cpu-index: int (optional) the index of the virtual CPU to use for translating the virtual address (defaults to CPU 0)

Returns

Nothing on success

Since

0.14

Notes

Errors were not reliably returned until 1.1

Example

```
-> { "execute": "memsave",
      "arguments": { "val": 10,
                     "size": 100,
                     "filename": "/tmp/virtual-mem-dump" } }
<- { "return": {} }
```

pmemsave (Command)

Save a portion of guest physical memory to a file.

Arguments

val: int the physical address of the guest to start from

size: int the size of memory region to save

filename: string the file to save the memory to as binary data

Returns

Nothing on success

Since

0.14

Notes

Errors were not reliably returned until 1.1

Example

```
-> { "execute": "pmemsave",
      "arguments": { "val": 10,
                     "size": 100,
                     "filename": "/tmp/physical-mem-dump" } }
<- { "return": {} }
```

Memdev (Object)

Information about memory backend

Members

id: string (optional) backend's ID if backend has 'id' property (since 2.9)

size: int memory backend size

merge: boolean enables or disables memory merge support

dump: boolean includes memory backend's memory in a core dump or not

prealloc: boolean enables or disables memory preallocation

host-nodes: array of int host nodes for its memory policy

policy: HostMemPolicy memory policy of memory backend

Since

2.1

query-memdev (Command)

Returns information for all memory backends.

Returns

a list of Memdev.

Since

2.1

Example

```
-> { "execute": "query-memdev" }
<- { "return": [
  {
    "id": "mem1",
    "size": 536870912,
    "merge": false,
    "dump": true,
    "prealloc": false,
    "host-nodes": [0, 1],
    "policy": "bind"
  },
  {
    "size": 536870912,
    "merge": false,
    "dump": true,
    "prealloc": true,
    "host-nodes": [2, 3],
    "policy": "preferred"
  }
]
```

(continues on next page)

(continued from previous page)

```

    }
  ]
}
```

CpuInstanceProperties (Object)

List of properties to be used for hotplugging a CPU instance, it should be passed by management with `device_add` command when a CPU is being hotplugged.

Members

node-id: int (optional) NUMA node ID the CPU belongs to

socket-id: int (optional) socket number within node/board the CPU belongs to

die-id: int (optional) die number within node/board the CPU belongs to (Since 4.1)

core-id: int (optional) core number within die the CPU belongs to

thread-id: int (optional) thread number within core the CPU belongs to

Note

currently there are 5 properties that could be present but management should be prepared to pass through other properties with `device_add` command to allow for future interface extension. This also requires the filed names to be kept in sync with the properties passed to `-device/device_add`.

Since

2.7

HotpluggableCPU (Object)

Members

type: string CPU object type for usage with `device_add` command

props: CpuInstanceProperties list of properties to be used for hotplugging CPU

vcpus-count: int number of logical VCPU threads `HotpluggableCPU` provides

qom-path: string (optional) link to existing CPU object if CPU is present or omitted if CPU is not present.

Since

2.7

query-hotpluggable-cpus (Command)

TODO

Better documentation; currently there is none.

Returns

a list of HotpluggableCPU objects.

Since

2.7

Example

For pseries machine type started with `-smp 2,cores=2,maxcpus=4 -cpu POWER8`:

```
-> { "execute": "query-hotpluggable-cpus" }
<- { "return": [
  { "props": { "core": 8 }, "type": "POWER8-spapr-cpu-core",
    "vcpus-count": 1 },
  { "props": { "core": 0 }, "type": "POWER8-spapr-cpu-core",
    "vcpus-count": 1, "qom-path": "/machine/unattached/device[0]" }
]}
```

For pc machine type started with `-smp 1,maxcpus=2`:

```
-> { "execute": "query-hotpluggable-cpus" }
<- { "return": [
  {
    "type": "qemu64-x86_64-cpu", "vcpus-count": 1,
    "props": { "core-id": 0, "socket-id": 1, "thread-id": 0 }
  },
  {
    "qom-path": "/machine/unattached/device[0]",
    "type": "qemu64-x86_64-cpu", "vcpus-count": 1,
    "props": { "core-id": 0, "socket-id": 0, "thread-id": 0 }
  }
]}
```

For s390x-virtio-ccw machine type started with `-smp 1,maxcpus=2 -cpu qemu`
(Since: 2.11):

```
-> { "execute": "query-hotpluggable-cpus" }
<- { "return": [
  {
    "type": "qemu-s390x-cpu", "vcpus-count": 1,
    "props": { "core-id": 1 }
  },
  {
    "qom-path": "/machine/unattached/device[0]",
    "type": "qemu-s390x-cpu", "vcpus-count": 1,
```

(continues on next page)

(continued from previous page)

```

    "props": { "core-id": 0 }
  }
}]

```

set-numa-node (Command)

Runtime equivalent of ‘-numa’ CLI option, available at preconfigure stage to configure numa mapping before initializing machine.

Since 3.0

Arguments

The members of `NumaOptions`

balloon (Command)

Request the balloon driver to change its balloon size.

Arguments

value: int the target logical size of the VM in bytes. We can deduce the size of the balloon using this formula:

$$\text{logical_vm_size} = \text{vm_ram_size} - \text{balloon_size}$$

$$\text{From it we have: } \text{balloon_size} = \text{vm_ram_size} - \text{value}$$

Returns

- Nothing on success
- If the balloon driver is enabled but not functional because the KVM kernel module cannot support it, `KvmMissingCap`
- If no balloon device is present, `DeviceNotActive`

Notes

This command just issues a request to the guest. When it returns, the balloon size may not have changed. A guest can change the balloon size independent of this command.

Since

0.14

Example

```
-> { "execute": "balloon", "arguments": { "value": 536870912 } }  
<- { "return": {} }
```

With a 2.5GiB guest this command inflated the ballon to 3GiB.

BalloonInfo (Object)

Information about the guest balloon device.

Members

actual: int the logical size of the VM in bytes Formula used: $\text{logical_vm_size} = \text{vm_ram_size} - \text{balloon_size}$

Since

0.14

query-balloon (Command)

Return information about the balloon device.

Returns

- `BalloonInfo` on success
- If the balloon driver is enabled but not functional because the KVM kernel module cannot support it, `KvmMissingCap`
- If no balloon device is present, `DeviceNotActive`

Since

0.14

Example

```
-> { "execute": "query-balloon" }  
<- { "return": {  
    "actual": 1073741824,  
  }  
}
```

BALLOON_CHANGE (Event)

Emitted when the guest changes the actual BALLOON level. This value is equivalent to the `actual` field return by the ‘query-balloon’ command

Arguments

actual: int the logical size of the VM in bytes Formula used: `logical_vm_size = vm_ram_size - balloon_size`

Note

this event is rate-limited.

Since

1.2

Example

```
<- { "event": "BALLOON_CHANGE",
      "data": { "actual": 944766976 },
      "timestamp": { "seconds": 1267020223, "microseconds": 435656 } }
```

MemoryInfo (Object)

Actual memory information in bytes.

Members

base-memory: int size of “base” memory specified with command line option -m.

plugged-memory: int (optional) size of memory that can be hot-unplugged. This field is omitted if target doesn’t support memory hotplug (i.e. CONFIG_MEM_DEVICE not defined at build time).

Since

2.11

query-memory-size-summary (Command)

Return the amount of initially allocated and present hotpluggable (if enabled) memory in bytes.

Example

```
-> { "execute": "query-memory-size-summary" }  
<- { "return": { "base-memory": 4294967296, "plugged-memory": 0 } }
```

Since

2.11

PCDIMMDeviceInfo (Object)

PCDIMMDevice state information

Members

id: string (optional) device's ID

addr: int physical address, where device is mapped

size: int size of memory that the device provides

slot: int slot number at which device is plugged in

node: int NUMA node number where device is plugged in

memdev: string memory backend linked with device

hotplugged: boolean true if device was hotplugged

hotpluggable: boolean true if device if could be added/removed while machine is running

Since

2.1

VirtioPMEMDeviceInfo (Object)

VirtioPMEM state information

Members

id: string (optional) device's ID

memaddr: int physical address in memory, where device is mapped

size: int size of memory that the device provides

memdev: string memory backend linked with device

Since

4.1

VirtioMEMDeviceInfo (Object)

VirtioMEMDevice state information

Members

id: **string** (optional) device's ID

memaddr: **int** physical address in memory, where device is mapped

requested-size: **int** the user requested size of the device

size: **int** the (current) size of memory that the device provides

max-size: **int** the maximum size of memory that the device can provide

block-size: **int** the block size of memory that the device provides

node: **int** NUMA node number where device is assigned to

memdev: **string** memory backend linked with the region

Since

5.1

MemoryDeviceInfo (Object)

Union containing information about a memory device

nvdimm is included since 2.12. virtio-pmem is included since 4.1. virtio-mem is included since 5.1.

Members

type One of `dim`, `nvdimm`, `virtio-pmem`, `virtio-mem`

data: `PCDIMMDeviceInfo` when **type** is `"dim`

`"nvdimm"`

data: `VirtioPMEMDeviceInfo` when **type** is `"virtio-pmem"`

data: `VirtioMEMDeviceInfo` when **type** is `"virtio-mem"`

Since

2.1

query-memory-devices (Command)

Lists available memory devices and their state

Since

2.1

Example

```
-> { "execute": "query-memory-devices" }
<- { "return": [ { "data":
    { "addr": 5368709120,
      "hotpluggable": true,
      "hotplugged": true,
      "id": "d1",
      "memdev": "/objects/memX",
      "node": 0,
      "size": 1073741824,
      "slot": 0},
    "type": "dimm"
  } ] }
```

MEMORY_DEVICE_SIZE_CHANGE (Event)

Emitted when the size of a memory device changes. Only emitted for memory devices that can actually change the size (e.g., virtio-mem due to guest action).

Arguments

id: **string** (optional) device's ID

size: **int** the new size of memory that the device provides

Note

this event is rate-limited.

Since

5.1

Example

```
<- { "event": "MEMORY_DEVICE_SIZE_CHANGE",
      "data": { "id": "vm0", "size": 1073741824},
      "timestamp": { "seconds": 1588168529, "microseconds": 201316 } }
```

MEM_UNPLUG_ERROR (Event)

Emitted when memory hot unplug error occurs.

Arguments

device: string device name

msg: string Informative message

Since

2.4

Example

```

<- { "event": "MEM_UNPLUG_ERROR"
      "data": { "device": "dimml",
                 "msg": "acpi: device unplug for unsupported device"
      },
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }

```

CpuModelInfo (Object)

Virtual CPU model.

A CPU model consists of the name of a CPU definition, to which delta changes are applied (e.g. features added/removed). Most magic values that an architecture might require should be hidden behind the name. However, if required, architectures can expose relevant properties.

Members

name: string the name of the CPU definition the model is based on

props: value (optional) a dictionary of QOM properties to be applied

Since

2.8

CpuModelExpansionType (Enum)

An enumeration of CPU model expansion types.

Values

static Expand to a static CPU model, a combination of a static base model name and property delta changes. As the static base model will never change, the expanded CPU model will be the same, independent of QEMU version, machine type, machine options, and accelerator options. Therefore, the resulting model can be used by tooling without having to specify a compatibility machine - e.g. when displaying the “host” model. The `static` CPU models are migration-safe.

full Expand all properties. The produced model is not guaranteed to be migration-safe, but allows tooling to get an insight and work with model details.

Note

When a non-migration-safe CPU model is expanded in static mode, some features enabled by the CPU model may be omitted, because they can't be implemented by a static CPU model definition (e.g. cache info passthrough and PMU passthrough in x86). If you need an accurate representation of the features enabled by a non-migration-safe CPU model, use `full`. If you need a static representation that will keep ABI compatibility even when changing QEMU version or machine-type, use `static` (but keep in mind that some features may be omitted).

Since

2.8

CpuModelCompareResult (Enum)

An enumeration of CPU model comparison results. The result is usually calculated using e.g. CPU features or CPU generations.

Values

incompatible If model A is incompatible to model B, model A is not guaranteed to run where model B runs and the other way around.

identical If model A is identical to model B, model A is guaranteed to run where model B runs and the other way around.

superset If model A is a superset of model B, model B is guaranteed to run where model A runs. There are no guarantees about the other way.

subset If model A is a subset of model B, model A is guaranteed to run where model B runs. There are no guarantees about the other way.

Since

2.8

CpuModelBaselineInfo (Object)

The result of a CPU model baseline.

Members

model: `CpuModelInfo` the baselined `CpuModelInfo`.

Since

2.8

If`defined(TARGET_S390X)`**CpuModelCompareInfo (Object)**

The result of a CPU model comparison.

Members

result: CpuModelCompareResult The result of the compare operation.

responsible-properties: array of string List of properties that led to the comparison result not being identical.

`responsible-properties` is a list of QOM property names that led to both CPUs not being detected as identical. For identical models, this list is empty. If a QOM property is read-only, that means there's no known way to make the CPU models identical. If the special property name "type" is included, the models are by definition not identical and cannot be made identical.

Since

2.8

If`defined(TARGET_S390X)`**query-cpu-model-comparison (Command)**

Compares two CPU models, returning how they compare in a specific configuration. The results indicates how both models compare regarding runnability. This result can be used by tooling to make decisions if a certain CPU model will run in a certain configuration or if a compatible CPU model has to be created by baselining.

Usually, a CPU model is compared against the maximum possible CPU model of a certain configuration (e.g. the "host" model for KVM). If that CPU model is identical or a subset, it will run in that configuration.

The result returned by this command may be affected by:

- QEMU version: CPU models may look different depending on the QEMU version. (Except for CPU models reported as "static" in `query-cpu-definitions`.)
- machine-type: CPU model may look different depending on the machine-type. (Except for CPU models reported as "static" in `query-cpu-definitions`.)
- machine options (including accelerator): in some architectures, CPU models may look different depending on machine and accelerator options. (Except for CPU models reported as "static" in `query-cpu-definitions`.)

- “-cpu” arguments and global properties: arguments to the -cpu option and global properties may affect expansion of CPU models. Using query-cpu-model-expansion while using these is not advised.

Some architectures may not support comparing CPU models. s390x supports comparing CPU models.

Arguments

modela: `CpuModelInfo` Not documented

modelb: `CpuModelInfo` Not documented

Returns

a `CpuModelBaselineInfo`. Returns an error if comparing CPU models is not supported, if a model cannot be used, if a model contains an unknown cpu definition name, unknown properties or properties with wrong types.

Note

this command isn't specific to s390x, but is only implemented on this architecture currently.

Since

2.8

If

`defined(TARGET_S390X)`

query-cpu-model-baseline (Command)

Baseline two CPU models, creating a compatible third model. The created model will always be a static, migration-safe CPU model (see “static” CPU model expansion for details).

This interface can be used by tooling to create a compatible CPU model out two CPU models. The created CPU model will be identical to or a subset of both CPU models when comparing them. Therefore, the created CPU model is guaranteed to run where the given CPU models run.

The result returned by this command may be affected by:

- QEMU version: CPU models may look different depending on the QEMU version. (Except for CPU models reported as “static” in query-cpu-definitions.)
- machine-type: CPU model may look different depending on the machine-type. (Except for CPU models reported as “static” in query-cpu-definitions.)
- machine options (including accelerator): in some architectures, CPU models may look different depending on machine and accelerator options. (Except for CPU models reported as “static” in query-cpu-definitions.)
- “-cpu” arguments and global properties: arguments to the -cpu option and global properties may affect expansion of CPU models. Using query-cpu-model-expansion while using these is not advised.

Some architectures may not support baselining CPU models. s390x supports baselining CPU models.

Arguments

modela: `CpuModelInfo` Not documented

modelb: `CpuModelInfo` Not documented

Returns

a `CpuModelBaselineInfo`. Returns an error if baselining CPU models is not supported, if a model cannot be used, if a model contains an unknown cpu definition name, unknown properties or properties with wrong types.

Note

this command isn't specific to s390x, but is only implemented on this architecture currently.

Since

2.8

If

`defined(TARGET_S390X)`

`CpuModelExpansionInfo` (Object)

The result of a cpu model expansion.

Members

model: `CpuModelInfo` the expanded `CpuModelInfo`.

Since

2.8

If

`defined(TARGET_S390X) || defined(TARGET_I386) || defined(TARGET_ARM)`

`query-cpu-model-expansion` (Command)

Expands a given CPU model (or a combination of CPU model + additional options) to different granularities, allowing tooling to get an understanding what a specific CPU model looks like in QEMU under a certain configuration.

This interface can be used to query the “host” CPU model.

The data returned by this command may be affected by:

- QEMU version: CPU models may look different depending on the QEMU version. (Except for CPU models reported as “static” in query-cpu-definitions.)
- machine-type: CPU model may look different depending on the machine-type. (Except for CPU models reported as “static” in query-cpu-definitions.)
- machine options (including accelerator): in some architectures, CPU models may look different depending on machine and accelerator options. (Except for CPU models reported as “static” in query-cpu-definitions.)
- “-cpu” arguments and global properties: arguments to the -cpu option and global properties may affect expansion of CPU models. Using query-cpu-model-expansion while using these is not advised.

Some architectures may not support all expansion types. s390x supports “full” and “static”. Arm only supports “full”.

Arguments

type: `CpuModelExpansionType` Not documented

model: `CpuModelInfo` Not documented

Returns

a `CpuModelExpansionInfo`. Returns an error if expanding CPU models is not supported, if the model cannot be expanded, if the model contains an unknown CPU definition name, unknown properties or properties with a wrong type. Also returns an error if an expansion type is not supported.

Since

2.8

If

```
defined(TARGET_S390X) || defined(TARGET_I386) || defined(TARGET_ARM)
```

CpuDefinitionInfo (Object)

Virtual CPU definition.

Members

name: `string` the name of the CPU definition

migration-safe: `boolean (optional)` whether a CPU definition can be safely used for migration in combination with a QEMU compatibility machine when migrating between different QEMU versions and between hosts with different sets of (hardware or software) capabilities. If not provided, information is not available and callers should not assume the CPU definition to be migration-safe. (since 2.8)

static: `boolean` whether a CPU definition is static and will not change depending on QEMU version, machine type, machine options and accelerator options. A static model is always migration-safe. (since 2.8)

unavailable-features: `array of string (optional)` List of properties that prevent the CPU model from running in the current host. (since 2.8)

typename: string Type name that can be used as argument to `device-list-properties`, to introspect properties configurable using `-cpu` or `-global`. (since 2.9)

alias-of: string (optional) Name of CPU model this model is an alias for. The target of the CPU model alias may change depending on the machine type. Management software is supposed to translate CPU model aliases in the VM configuration, because aliases may stop being migration-safe in the future (since 4.1)

deprecated: boolean If true, this CPU model is deprecated and may be removed in in some future version of QEMU according to the QEMU deprecation policy. (since 5.2)

`unavailable-features` is a list of QOM property names that represent CPU model attributes that prevent the CPU from running. If the QOM property is read-only, that means there's no known way to make the CPU model run in the current host. Implementations that choose not to provide specific information return the property name "type". If the property is read-write, it means that it MAY be possible to run the CPU model in the current host if that property is changed. Management software can use it as hints to suggest or choose an alternative for the user, or just to generate meaningful error messages explaining why the CPU model can't be used. If `unavailable-features` is an empty list, the CPU model is runnable using the current host and machine-type. If `unavailable-features` is not present, runnability information for the CPU is not available.

Since

1.2

If

```
defined(TARGET_PPC) || defined(TARGET_ARM) || defined(TARGET_I386) ||
defined(TARGET_S390X) || defined(TARGET_MIPS)
```

query-cpu-definitions (Command)

Return a list of supported virtual CPU definitions

Returns

a list of `CpuDefInfo`

Since

1.2

If

```
defined(TARGET_PPC) || defined(TARGET_ARM) || defined(TARGET_I386) ||
defined(TARGET_S390X) || defined(TARGET_MIPS)
```

4.8.26 Record/replay

ReplayMode (Enum)

Mode of the replay subsystem.

Values

none normal execution mode. Replay or record are not enabled.

record record mode. All non-deterministic data is written into the replay log.

play replay mode. Non-deterministic data required for system execution is read from the log.

Since

2.5

ReplayInfo (Object)

Record/replay information.

Members

mode: ReplayMode current mode.

filename: string (optional) name of the record/replay log file. It is present only in record or replay modes, when the log is recorded or replayed.

icount: int current number of executed instructions.

Since

5.2

query-replay (Command)

Retrieve the record/replay information. It includes current instruction count which may be used for `replay-break` and `replay-seek` commands.

Returns

record/replay information.

Since

5.2

Example

```
-> { "execute": "query-replay" }
<- { "return": { "mode": "play", "filename": "log.rr", "icount": 220414 } }
```

replay-break (Command)

Set replay breakpoint at instruction count `icount`. Execution stops when the specified instruction is reached. There can be at most one breakpoint. When breakpoint is set, any prior one is removed. The breakpoint may be set only in replay mode and only “in the future”, i.e. at instruction counts greater than the current one. The current instruction count can be observed with `query-replay`.

Arguments

icount: int instruction count to stop at

Since

5.2

Example

```
-> { "execute": "replay-break", "data": { "icount": 220414 } }
```

replay-delete-break (Command)

Remove replay breakpoint which was set with `replay-break`. The command is ignored when there are no replay breakpoints.

Since

5.2

Example

```
-> { "execute": "replay-delete-break" }
```

replay-seek (Command)

Automatically proceed to the instruction count `icount`, when replaying the execution. The command automatically loads nearest snapshot and replays the execution to find the desired instruction. When there is no preceding snapshot or the execution is not replayed, then the command fails. `icount` for the reference may be obtained with `query-replay` command.

Arguments

icount: int target instruction count

Since

5.2

Example

```
-> { "execute": "replay-seek", "data": { "icount": 220414 } }
```

4.8.27 Yank feature

YankInstanceType (Enum)

An enumeration of yank instance types. See `YankInstance` for more information.

Values

block-node Not documented

chardev Not documented

migration Not documented

Since

6.0

YankInstanceBlockNode (Object)

Specifies which block graph node to yank. See `YankInstance` for more information.

Members

node-name: string the name of the block graph node

Since

6.0

YankInstanceChardev (Object)

Specifies which character device to yank. See `YankInstance` for more information.

Members

id: `string` the chardev's ID

Since

6.0

YankInstance (Object)

A yank instance can be yanked with the `yank qmp` command to recover from a hanging QEMU.

Currently implemented yank instances:

- nbd block device: Yanking it will shut down the connection to the nbd server without attempting to reconnect.
- socket chardev: Yanking it will shut down the connected socket.
- migration: Yanking it will shut down all migration connections. Unlike `migrate_cancel`, it will not notify the migration process, so migration will go into `failed` state, instead of `cancelled` state. `yank` should be used to recover from hangs.

Members

type: `YankInstanceType` Not documented

The members of `YankInstanceBlockNode` when `type` is `"block-node"`

The members of `YankInstanceChardev` when `type` is `"chardev"`

Since

6.0

yank (Command)

Try to recover from hanging QEMU by yanking the specified instances. See `YankInstance` for more information.

Takes a list of `YankInstance` as argument.

Arguments

instances: `array of YankInstance` Not documented

Returns

- Nothing on success
- `DeviceNotFound` error, if any of the `YankInstances` doesn't exist

Example

```
-> { "execute": "yank",
      "arguments": {
        "instances": [
          { "type": "block-node",
            "node-name": "nbd0" }
        ] } }
<- { "return": {} }
```

Since

6.0

query-yank (Command)

Query yank instances. See `YankInstance` for more information.

Returns

list of `YankInstance`

Example

```
-> { "execute": "query-yank" }
<- { "return": [
      { "type": "block-node",
        "node-name": "nbd0" }
    ] }
```

Since

6.0

4.8.28 Miscellanea

add_client (Command)

Allow client connections for VNC, Spice and socket based character devices to be passed in to QEMU via `SCM_RIGHTS`.

Arguments

protocol: `string` protocol name. Valid names are “vnc”, “spice” or the name of a character device (eg. from `-chardev id=XXXX`)

fdname: `string` file descriptor name previously passed via ‘getfd’ command

skipauth: boolean (optional) whether to skip authentication. Only applies to “vnc” and “spice” protocols

tls: boolean (optional) whether to perform TLS. Only applies to the “spice” protocol

Returns

nothing on success.

Since

0.14

Example

```
-> { "execute": "add_client", "arguments": { "protocol": "vnc",  
                                           "fdname": "myclient" } }  
<- { "return": {} }
```

NameInfo (Object)

Guest name information.

Members

name: string (optional) The name of the guest

Since

0.14

query-name (Command)

Return the name information of a guest.

Returns

NameInfo of the guest

Since

0.14

Example

```
-> { "execute": "query-name" }  
<- { "return": { "name": "qemu-name" } }
```

IOThreadInfo (Object)

Information about an iothread

Members

id: `string` the identifier of the iothread

thread-id: `int` ID of the underlying host thread

poll-max-ns: `int` maximum polling time in ns, 0 means polling is disabled (since 2.9)

poll-grow: `int` how many ns will be added to polling time, 0 means that it's not configured (since 2.9)

poll-shrink: `int` how many ns will be removed from polling time, 0 means that it's not configured (since 2.9)

Since

2.0

query-iothreads (Command)

Returns a list of information about each iothread.

Note

this list excludes the QEMU main loop thread, which is not declared using the `-object iothread` command-line option. It is always the main thread of the process.

Returns

a list of `IOThreadInfo` for each iothread

Since

2.0

Example

```
-> { "execute": "query-iothreads" }
<- { "return": [
    {
        "id": "iothread0",
        "thread-id": 3134
    },
    {
        "id": "iothread1",
        "thread-id": 3135
    }
  ]
}
```

stop (Command)

Stop all guest VCPU execution.

Since

0.14

Notes

This function will succeed even if the guest is already in the stopped state. In “inmigrate” state, it will ensure that the guest remains paused once migration finishes, as if the -S option was passed on the command line.

Example

```
-> { "execute": "stop" }
<- { "return": {} }
```

cont (Command)

Resume guest VCPU execution.

Since

0.14

Returns

If successful, nothing

Notes

This command will succeed if the guest is currently running. It will also succeed if the guest is in the “inmigrate” state; in this case, the effect of the command is to make sure the guest starts once migration finishes, removing the effect of the -S command line option if it was passed.

Example

```
-> { "execute": "cont" }
<- { "return": {} }
```

x-exit-preconfig (Command)

Exit from “preconfig” state

This command makes QEMU exit the preconfig state and proceed with VM initialization using configuration data provided on the command line and via the QMP monitor during the preconfig state. The command is only available during the preconfig state (i.e. when the `-preconfig` command line option was in use).

Since 3.0

Returns

nothing

Example

```
-> { "execute": "x-exit-preconfig" }
<- { "return": {} }
```

human-monitor-command (Command)

Execute a command on the human monitor and return the output.

Arguments

command-line: string the command to execute in the human monitor

cpu-index: int (optional) The CPU to use for commands that require an implicit CPU

Features

savevm-monitor-nodes If present, HMP command `savevm` only snapshots monitor-owned nodes if they have no parents. This allows the use of ‘`savevm`’ with `-blockdev`. (since 4.2)

Returns

the output of the command as a string

Since

0.14

Notes

This command only exists as a stop-gap. Its use is highly discouraged. The semantics of this command are not guaranteed: this means that command names, arguments and responses can change or be removed at ANY time. Applications that rely on long term stability guarantees should NOT use this command.

Known limitations:

- This command is stateless, this means that commands that depend on state information (such as `getfd`) might not work
- Commands that prompt the user for data don't currently work

Example

```
-> { "execute": "human-monitor-command",
      "arguments": { "command-line": "info kvm" } }
<- { "return": "kvm support: enabled\r\n" }
```

`getfd` (Command)

Receive a file descriptor via SCM rights and assign it a name

Arguments

fdname: `string` file descriptor name

Returns

Nothing on success

Since

0.14

Notes

If `fdname` already exists, the file descriptor assigned to it will be closed and replaced by the received file descriptor. The `'closefd'` command can be used to explicitly close the file descriptor when it is no longer needed.

Example

```
-> { "execute": "getfd", "arguments": { "fdname": "fd1" } }  
<- { "return": {} }
```

closefd (Command)

Close a file descriptor previously passed via SCM rights

Arguments

fdname: string file descriptor name

Returns

Nothing on success

Since

0.14

Example

```
-> { "execute": "closefd", "arguments": { "fdname": "fd1" } }  
<- { "return": {} }
```

AddfdInfo (Object)

Information about a file descriptor that was added to an fd set.

Members

fdset-id: int The ID of the fd set that `fd` was added to.

fd: int The file descriptor that was received via SCM rights and added to the fd set.

Since

1.2

add-fd (Command)

Add a file descriptor, that was passed via SCM rights, to an fd set.

Arguments

fdset-id: int (optional) The ID of the fd set to add the file descriptor to.

opaque: string (optional) A free-form string that can be used to describe the fd.

Returns

- AddfdInfo on success
- If file descriptor was not received, FdNotSupplied
- If fdset-id is a negative value, InvalidParameterValue

Notes

The list of fd sets is shared by all monitor connections.

If fdset-id is not specified, a new fd set will be created.

Since

1.2

Example

```
-> { "execute": "add-fd", "arguments": { "fdset-id": 1 } }
<- { "return": { "fdset-id": 1, "fd": 3 } }
```

remove-fd (Command)

Remove a file descriptor from an fd set.

Arguments

fdset-id: int The ID of the fd set that the file descriptor belongs to.

fd: int (optional) The file descriptor that is to be removed.

Returns

- Nothing on success
- If fdset-id or fd is not found, FdNotFound

Since

1.2

Notes

The list of fd sets is shared by all monitor connections.

If `fd` is not specified, all file descriptors in `fdset-id` will be removed.

Example

```
-> { "execute": "remove-fd", "arguments": { "fdset-id": 1, "fd": 3 } }
<- { "return": {} }
```

FdsetFdInfo (Object)

Information about a file descriptor that belongs to an fd set.

Members

fd: int The file descriptor value.

opaque: string (optional) A free-form string that can be used to describe the fd.

Since

1.2

FdsetInfo (Object)

Information about an fd set.

Members

fdset-id: int The ID of the fd set.

fds: array of FdsetFdInfo A list of file descriptors that belong to this fd set.

Since

1.2

query-fdsets (Command)

Return information describing all fd sets.

Returns

A list of `FdsetInfo`

Since

1.2

Note

The list of fd sets is shared by all monitor connections.

Example

```

-> { "execute": "query-fdsets" }
<- { "return": [
  {
    "fds": [
      {
        "fd": 30,
        "opaque": "rdonly:/path/to/file"
      },
      {
        "fd": 24,
        "opaque": "rdwr:/path/to/file"
      }
    ],
    "fdset-id": 1
  },
  {
    "fds": [
      {
        "fd": 28
      },
      {
        "fd": 29
      }
    ],
    "fdset-id": 0
  }
]
}

```

CommandLineParameterType (Enum)

Possible types for an option parameter.

Values

string accepts a character string

boolean accepts “on” or “off”

number accepts a number

size accepts a number followed by an optional suffix (K)ilo, (M)ega, (G)iga, (T)era

Since

1.5

CommandLineParameterInfo (Object)

Details about a single parameter of a command line option.

Members

name: `string` parameter name

type: `CommandLineParameterType` parameter `CommandLineParameterType`

help: `string (optional)` human readable text string, not suitable for parsing.

default: `string (optional)` default value string (since 2.1)

Since

1.5

CommandLineOptionInfo (Object)

Details about a command line option, including its list of parameter details

Members

option: `string` option name

parameters: `array of CommandLineParameterInfo` an array of `CommandLineParameterInfo`

Since

1.5

query-command-line-options (Command)

Query command line option schema.

Arguments

option: `string (optional)` option name

Returns

list of `CommandLineOptionInfo` for all options (or for the given `option`). Returns an error if the given `option` doesn't exist.

Since

1.5

Example

```

-> { "execute": "query-command-line-options",
    "arguments": { "option": "option-rom" } }
<- { "return": [
    {
        "parameters": [
            {
                "name": "romfile",
                "type": "string"
            },
            {
                "name": "bootindex",
                "type": "number"
            }
        ],
        "option": "option-rom"
    }
]
}

```

RTC_CHANGE (Event)

Emitted when the guest changes the RTC time.

Arguments

offset: **int** offset between base RTC clock (as specified by -rtc base), and new RTC clock value

Note

This event is rate-limited.

Since

0.13

Example

```

<- { "event": "RTC_CHANGE",
    "data": { "offset": 78 },
    "timestamp": { "seconds": 1267020223, "microseconds": 435656 } }

```

If

```
defined(TARGET_ALPHA) || defined(TARGET_ARM) || defined(TARGET_HPPA) ||  
defined(TARGET_I386) || defined(TARGET_MIPS) || defined(TARGET_MIPS64) ||  
defined(TARGET_MOXIE) || defined(TARGET_PPC) || defined(TARGET_PPC64) ||  
defined(TARGET_S390X) || defined(TARGET_SH4) || defined(TARGET_SPARC)
```

rtc-reset-reinjection (Command)

This command will reset the RTC interrupt reinjection backlog. Can be used if another mechanism to synchronize guest time is in effect, for example QEMU guest agent's guest-set-time command.

Since

2.1

Example

```
-> { "execute": "rtc-reset-reinjection" }  
<- { "return": {} }
```

If

```
defined(TARGET_I386)
```

SevState (Enum)

An enumeration of SEV state information used during `query-sev`.

Values

uninit The guest is uninitialized.

launch-update The guest is currently being launched; plaintext data and register state is being imported.

launch-secret The guest is currently being launched; ciphertext data is being imported.

running The guest is fully launched or migrated in.

send-update The guest is currently being migrated out to another machine.

receive-update The guest is currently being migrated from another machine.

Since

2.12

If

```
defined(TARGET_I386)
```

SevInfo (Object)

Information about Secure Encrypted Virtualization (SEV) support

Members

enabled: boolean true if SEV is active
api-major: int SEV API major version
api-minor: int SEV API minor version
build-id: int SEV FW build id
policy: int SEV policy value
state: SevState SEV guest state
handle: int SEV firmware handle

Since

2.12

If

```
defined(TARGET_I386)
```

query-sev (Command)

Returns information about SEV

Returns

SevInfo

Since

2.12

Example

```
-> { "execute": "query-sev" }
<- { "return": { "enabled": true, "api-major" : 0, "api-minor" : 0,
                  "build-id" : 0, "policy" : 0, "state" : "running",
                  "handle" : 1 } }
```

If

defined(TARGET_I386)

SevLaunchMeasureInfo (Object)

SEV Guest Launch measurement information

Members

data: string the measurement value encoded in base64

Since

2.12

If

defined(TARGET_I386)

query-sev-launch-measure (Command)

Query the SEV guest launch information.

Returns

The SevLaunchMeasureInfo for the guest

Since

2.12

Example

```
-> { "execute": "query-sev-launch-measure" }
<- { "return": { "data": "4l8LXeNlSPUDlXPJG5966/8%YZ" } }
```


If

```
defined(TARGET_I386)
```

SevCapability (Object)

The struct describes capability for a Secure Encrypted Virtualization feature.

Members

pdh: string Platform Diffie-Hellman key (base64 encoded)

cert-chain: string PDH certificate chain (base64 encoded)

cbitpos: int C-bit location in page table entry

reduced-phys-bits: int Number of physical Address bit reduction when SEV is enabled

Since

2.12

If

```
defined(TARGET_I386)
```

query-sev-capabilities (Command)

This command is used to get the SEV capabilities, and is supported on AMD X86 platforms only.

Returns

SevCapability objects.

Since

2.12

Example

```
-> { "execute": "query-sev-capabilities" }
<- { "return": { "pdh": "8CCDD8DDD", "cert-chain": "888CCDDDEE",
                  "cbitpos": 47, "reduced-phys-bits": 5}}
```

If

```
defined(TARGET_I386)
```

sev-inject-launch-secret (Command)

This command injects a secret blob into memory of SEV guest.

Arguments

packet-header: string the launch secret packet header encoded in base64

secret: string the launch secret data to be injected encoded in base64

gpa: int (optional) the guest physical address where secret will be injected.

Since

6.0

If

`defined(TARGET_I386)`

dump-keys (Command)

Dump guest's storage keys

Arguments

filename: string the path to the file to dump to

This command is only supported on s390 architecture.

Since

2.5

Example

```
-> { "execute": "dump-keys",  
    "arguments": { "filename": "/tmp/skeys" } }  
<- { "return": {} }
```

If

`defined(TARGET_S390X)`

GICCapability (Object)

The struct describes capability for a specific GIC (Generic Interrupt Controller) version. These bits are not only decided by QEMU/KVM software version, but also decided by the hardware that the program is running upon.

Members

version: int version of GIC to be described. Currently, only 2 and 3 are supported.

emulated: boolean whether current QEMU/hardware supports emulated GIC device in user space.

kernel: boolean whether current QEMU/hardware supports hardware accelerated GIC device in kernel.

Since

2.6

If

defined(TARGET_ARM)

query-gic-capabilities (Command)

This command is ARM-only. It will return a list of GICCapability objects that describe its capability bits.

Returns

a list of GICCapability objects.

Since

2.6

Example

```
-> { "execute": "query-gic-capabilities" }
<- { "return": [{ "version": 2, "emulated": true, "kernel": false },
                  { "version": 3, "emulated": false, "kernel": true } ] }
```

If

defined(TARGET_ARM)

4.8.29 Audio

AudiodevPerDirectionOptions (Object)

General audio backend options that are used for both playback and recording.

Members

mixing-engine: boolean (optional) use QEMU's mixing engine to mix all streams inside QEMU and convert audio formats when not supported by the backend. When set to off, fixed-settings must be also off (default on, since 4.2)

fixed-settings: boolean (optional) use fixed settings for host input/output. When off, frequency, channels and format must not be specified (default true)

frequency: int (optional) frequency to use when using fixed settings (default 44100)

channels: int (optional) number of channels when using fixed settings (default 2)

voices: int (optional) number of voices to use (default 1)

format: AudioFormat (optional) sample format to use when using fixed settings (default s16)

buffer-length: int (optional) the buffer length in microseconds

Since

4.0

AudiodevGenericOptions (Object)

Generic driver-specific options.

Members

in: AudiodevPerDirectionOptions (optional) options of the capture stream

out: AudiodevPerDirectionOptions (optional) options of the playback stream

Since

4.0

AudiodevAlsaPerDirectionOptions (Object)

Options of the ALSA backend that are used for both playback and recording.

Members

dev: string (optional) the name of the ALSA device to use (default ‘default’)

period-length: int (optional) the period length in microseconds

try-poll: boolean (optional) attempt to use poll mode, falling back to non-polling access on failure (default true)

The members of **AudiodevPerDirectionOptions**

Since

4.0

AudiodevAlsaOptions (Object)

Options of the ALSA audio backend.

Members

in: AudiodevAlsaPerDirectionOptions (optional) options of the capture stream

out: AudiodevAlsaPerDirectionOptions (optional) options of the playback stream

threshold: int (optional) set the threshold (in microseconds) when playback starts

Since

4.0

AudiodevCoreaudioPerDirectionOptions (Object)

Options of the Core Audio backend that are used for both playback and recording.

Members

buffer-count: int (optional) number of buffers

The members of **AudiodevPerDirectionOptions**

Since

4.0

AudiodevCoreaudioOptions (Object)

Options of the coreaudio audio backend.

Members

in: AudiodevCoreaudioPerDirectionOptions (optional) options of the capture stream

out: AudiodevCoreaudioPerDirectionOptions (optional) options of the playback stream

Since

4.0

AudiodevDsoundOptions (Object)

Options of the DirectSound audio backend.

Members

in: AudiodevPerDirectionOptions (optional) options of the capture stream

out: AudiodevPerDirectionOptions (optional) options of the playback stream

latency: int (optional) add extra latency to playback in microseconds (default 10000)

Since

4.0

AudiodevJackPerDirectionOptions (Object)

Options of the JACK backend that are used for both playback and recording.

Members

server-name: string (optional) select from among several possible concurrent server instances (default: environment variable \$JACK_DEFAULT_SERVER if set, else “default”)

client-name: string (optional) the client name to use. The server will modify this name to create a unique variant, if needed unless `exact-name` is true (default: the guest’s name)

connect-ports: string (optional) if set, a regular expression of JACK client port name(s) to monitor for and automatically connect to

start-server: boolean (optional) start a jack server process if one is not already present (default: false)

exact-name: boolean (optional) use the exact name requested otherwise JACK automatically generates a unique one, if needed (default: false)

The members of AudiodevPerDirectionOptions

Since

5.1

AudiodevJackOptions (Object)

Options of the JACK audio backend.

Members

in: AudiodevJackPerDirectionOptions (optional) options of the capture stream

out: AudiodevJackPerDirectionOptions (optional) options of the playback stream

Since

5.1

AudiodevOssPerDirectionOptions (Object)

Options of the OSS backend that are used for both playback and recording.

Members

dev: string (optional) file name of the OSS device (default '/dev/dsp')

buffer-count: int (optional) number of buffers

try-poll: boolean (optional) attempt to use poll mode, falling back to non-polling access on failure (default true)

The members of AudiodevPerDirectionOptions

Since

4.0

AudiodevOssOptions (Object)

Options of the OSS audio backend.

Members

in: AudiodevOssPerDirectionOptions (optional) options of the capture stream

out: AudiodevOssPerDirectionOptions (optional) options of the playback stream

try-mmap: boolean (optional) try using memory-mapped access, falling back to non-memory-mapped access on failure (default true)

exclusive: boolean (optional) open device in exclusive mode (vmix won't work) (default false)

dsp-policy: int (optional) set the timing policy of the device (between 0 and 10, where smaller number means smaller latency but higher CPU usage) or -1 to use fragment mode (option ignored on some platforms) (default 5)

Since

4.0

AudiodevPaPerDirectionOptions (Object)

Options of the Pulseaudio backend that are used for both playback and recording.

Members

name: string (optional) name of the sink/source to use

stream-name: string (optional) name of the PulseAudio stream created by qemu. Can be used to identify the stream in PulseAudio when you create multiple PulseAudio devices or run multiple qemu instances (default: audiodev's id, since 4.2)

latency: int (optional) latency you want PulseAudio to achieve in microseconds (default 15000)

The members of AudiodevPerDirectionOptions

Since

4.0

AudiodevPaOptions (Object)

Options of the PulseAudio audio backend.

Members

in: AudiodevPaPerDirectionOptions (optional) options of the capture stream

out: AudiodevPaPerDirectionOptions (optional) options of the playback stream

server: string (optional) PulseAudio server address (default: let PulseAudio choose)

Since

4.0

AudiodevSdlPerDirectionOptions (Object)

Options of the SDL audio backend that are used for both playback and recording.

Members

buffer-count: int (optional) number of buffers (default 4)

The members of AudiodevPerDirectionOptions

Since

6.0

AudiodevSdlOptions (Object)

Options of the SDL audio backend.

Members

in: AudiodevSdlPerDirectionOptions (optional) options of the recording stream

out: AudiodevSdlPerDirectionOptions (optional) options of the playback stream

Since

6.0

AudiodevWavOptions (Object)

Options of the wav audio backend.

Members

in: AudiodevPerDirectionOptions (optional) options of the capture stream

out: AudiodevPerDirectionOptions (optional) options of the playback stream

path: string (optional) name of the wav file to record (default 'qemu.wav')

Since

4.0

AudioFormat (Enum)

An enumeration of possible audio formats.

Values

u8 unsigned 8 bit integer

s8 signed 8 bit integer

u16 unsigned 16 bit integer

s16 signed 16 bit integer

u32 unsigned 32 bit integer

s32 signed 32 bit integer

f32 single precision floating-point (since 5.0)

Since

4.0

AudiodevDriver (Enum)

An enumeration of possible audio backend drivers.

Values

jack JACK audio backend (since 5.1)

none Not documented

alsa Not documented

coreaudio Not documented

dsound Not documented

oss Not documented

pa Not documented

sdl Not documented

spice Not documented

wav Not documented

Since

4.0

Audiodev (Object)

Options of an audio backend.

Members

id: string identifier of the backend

driver: AudiodevDriver the backend driver to use

timer-period: int (optional) timer period (in microseconds, 0: use lowest possible)

The members of **AudiodevGenericOptions** when **driver** is "none"

The members of **AudiodevAlsaOptions** when **driver** is "alsa"

The members of **AudiodevCoreaudioOptions** when **driver** is "coreaudio"

The members of **AudiodevDsoundOptions** when **driver** is "dsound"

The members of **AudiodevJackOptions** when **driver** is "jack"

The members of `AudiodevOssOptions` when driver is "oss"

The members of `AudiodevPaOptions` when driver is "pa"

The members of `AudiodevSdlOptions` when driver is "sdl"

The members of `AudiodevGenericOptions` when driver is "spice"

The members of `AudiodevWavOptions` when driver is "wav"

Since

4.0

4.8.30 ACPI

`AcpiTableOptions` (Object)

Specify an ACPI table on the command line to load.

At most one of `file` and `data` can be specified. The list of files specified by any one of them is loaded and concatenated in order. If both are omitted, `data` is implied.

Other fields / optargs can be used to override fields of the generic ACPI table header; refer to the ACPI specification 5.0, section 5.2.6 System Description Table Header. If a header field is not overridden, then the corresponding value from the concatenated blob is used (in case of `file`), or it is filled in with a hard-coded value (in case of `data`).

String fields are copied into the matching ACPI member from lowest address upwards, and silently truncated / NUL-padded to length.

Members

sig: string (optional) table signature / identifier (4 bytes)

rev: int (optional) table revision number (dependent on signature, 1 byte)

oem_id: string (optional) OEM identifier (6 bytes)

oem_table_id: string (optional) OEM table identifier (8 bytes)

oem_rev: int (optional) OEM-supplied revision number (4 bytes)

asl_compiler_id: string (optional) identifier of the utility that created the table (4 bytes)

asl_compiler_rev: int (optional) revision number of the utility that created the table (4 bytes)

file: string (optional) colon (:) separated list of pathnames to load and concatenate as table data. The resultant binary blob is expected to have an ACPI table header. At least one file is required. This field excludes `data`.

data: string (optional) colon (:) separated list of pathnames to load and concatenate as table data. The resultant binary blob must not have an ACPI table header. At least one file is required. This field excludes `file`.

Since

1.5

ACPISlotType (Enum)

Values

DIMM memory slot

CPU logical CPU slot (since 2.7)

ACPIOSTInfo (Object)

OSPM Status Indication for a device For description of possible values of `source` and `status` fields see “_OST (OSPM Status Indication)” chapter of ACPI5.0 spec.

Members

device: string (optional) device ID associated with slot

slot: string slot ID, unique per slot of a given `slot-type`

slot-type: ACPISlotType type of the slot

source: int an integer containing the source event

status: int an integer containing the status code

Since

2.1

query-acpi-ospm-status (Command)

Return a list of ACPIOSTInfo for devices that support status reporting via ACPI _OST method.

Since

2.1

Example

```
-> { "execute": "query-acpi-ospm-status" }
<- { "return": [ { "device": "d1", "slot": "0", "slot-type": "DIMM", "source": 1,
↪ "status": 0 },
                { "slot": "1", "slot-type": "DIMM", "source": 0, "status": 0 },
                { "slot": "2", "slot-type": "DIMM", "source": 0, "status": 0 },
                { "slot": "3", "slot-type": "DIMM", "source": 0, "status": 0 }
      ] }
```

ACPI_DEVICE_OST (Event)

Emitted when guest executes ACPI _OST method.

Arguments

info: `ACPIOSTInfo` OSPM Status Indication

Since

2.1

Example

```
<- { "event": "ACPI_DEVICE_OST",
      "data": { "device": "d1", "slot": "0",
                 "slot-type": "DIMM", "source": 1, "status": 0 } }
```

4.8.31 PCI

`PciMemoryRange` (Object)

A PCI device memory region

Members

base: `int` the starting address (guest physical)

limit: `int` the ending address (guest physical)

Since

0.14

`PciMemoryRegion` (Object)

Information about a PCI device I/O region.

Members

bar: `int` the index of the Base Address Register for this region

type: `string`

- 'io' if the region is a PIO region
- 'memory' if the region is a MMIO region

size: `int` memory size

prefetch: `boolean` (optional) if `type` is 'memory', true if the memory is prefetchable

mem_type_64: `boolean` (optional) if `type` is 'memory', true if the BAR is 64-bit

address: `int` Not documented

Since

0.14

PciBusInfo (Object)

Information about a bus of a PCI Bridge device

Members

number: int primary bus interface number. This should be the number of the bus the device resides on.

secondary: int secondary bus interface number. This is the number of the main bus for the bridge

subordinate: int This is the highest number bus that resides below the bridge.

io_range: PciMemoryRange The PIO range for all devices on this bridge

memory_range: PciMemoryRange The MMIO range for all devices on this bridge

prefetchable_range: PciMemoryRange The range of prefetchable MMIO for all devices on this bridge

Since

2.4

PciBridgeInfo (Object)

Information about a PCI Bridge device

Members

bus: PciBusInfo information about the bus the device resides on

devices: array of PciDeviceInfo (optional) a list of `PciDeviceInfo` for each device on this bridge

Since

0.14

PciDeviceClass (Object)

Information about the Class of a PCI device

Members

desc: string (optional) a string description of the device's class

class: int the class code of the device

Since

2.4

PciDeviceId (Object)

Information about the Id of a PCI device

Members

device: int the PCI device id

vendor: int the PCI vendor id

subsystem: int (optional) the PCI subsystem id (since 3.1)

subsystem-vendor: int (optional) the PCI subsystem vendor id (since 3.1)

Since

2.4

PciDeviceInfo (Object)

Information about a PCI device

Members

bus: int the bus number of the device

slot: int the slot the device is located in

function: int the function of the slot used by the device

class_info: PciDeviceClass the class of the device

id: PciDeviceId the PCI device id

irq: int (optional) if an IRQ is assigned to the device, the IRQ number

irq_pin: int the IRQ pin, zero means no IRQ (since 5.1)

qdev_id: string the device name of the PCI device

pci_bridge: PciBridgeInfo (optional) if the device is a PCI bridge, the bridge information

regions: array of PciMemoryRegion a list of the PCI I/O regions associated with the device

Notes

the contents of `class_info.desc` are not stable and should only be treated as informational.

Since

0.14

PciInfo (Object)

Information about a PCI bus

Members

bus: int the bus index

devices: array of PciDeviceInfo a list of devices on this bus

Since

0.14

query-pci (Command)

Return information about the PCI bus topology of the guest.

Returns

a list of `PciInfo` for each PCI bus. Each bus is represented by a json-object, which has a key with a json-array of all PCI devices attached to it. Each device is represented by a json-object.

Since

0.14

Example

```
-> { "execute": "query-pci" }
<- { "return": [
  {
    "bus": 0,
    "devices": [
      {
        "bus": 0,
        "qdev_id": "",
        "slot": 0,
        "class_info": {
          "class": 1536,
          "desc": "Host bridge"
        },
        "id": {
          "device": 32902,
```

(continues on next page)

(continued from previous page)

```

        "vendor": 4663
    },
    "function": 0,
    "regions": [
    ]
},
{
    "bus": 0,
    "qdev_id": "",
    "slot": 1,
    "class_info": {
        "class": 1537,
        "desc": "ISA bridge"
    },
    "id": {
        "device": 32902,
        "vendor": 28672
    },
    "function": 0,
    "regions": [
    ]
},
{
    "bus": 0,
    "qdev_id": "",
    "slot": 1,
    "class_info": {
        "class": 257,
        "desc": "IDE controller"
    },
    "id": {
        "device": 32902,
        "vendor": 28688
    },
    "function": 1,
    "regions": [
        {
            "bar": 4,
            "size": 16,
            "address": 49152,
            "type": "io"
        }
    ]
},
{
    "bus": 0,
    "qdev_id": "",
    "slot": 2,
    "class_info": {
        "class": 768,
        "desc": "VGA controller"
    },
    "id": {
        "device": 4115,
        "vendor": 184
    },
    "function": 0,

```

(continues on next page)

(continued from previous page)

```

        "regions": [
            {
                "prefetch": true,
                "mem_type_64": false,
                "bar": 0,
                "size": 33554432,
                "address": 4026531840,
                "type": "memory"
            },
            {
                "prefetch": false,
                "mem_type_64": false,
                "bar": 1,
                "size": 4096,
                "address": 4060086272,
                "type": "memory"
            },
            {
                "prefetch": false,
                "mem_type_64": false,
                "bar": 6,
                "size": 65536,
                "address": -1,
                "type": "memory"
            }
        ],
    },
    {
        "bus": 0,
        "qdev_id": "",
        "irq": 11,
        "slot": 4,
        "class_info": {
            "class": 1280,
            "desc": "RAM controller"
        },
        "id": {
            "device": 6900,
            "vendor": 4098
        },
        "function": 0,
        "regions": [
            {
                "bar": 0,
                "size": 32,
                "address": 49280,
                "type": "io"
            }
        ]
    }
]
}

```

Note

This example has been shortened as the real response is too long.

4.9 QEMU Storage Daemon QMP Reference Manual

4.9.1 Block devices

Block core (VM unrelated)

4.9.2 Common data types

`IoOperationType` (Enum)

An enumeration of the I/O operation types

Values

read read operation

write write operation

Since

2.1

`OnOffAuto` (Enum)

An enumeration of three options: on, off, and auto

Values

auto QEMU selects the value between on and off

on Enabled

off Disabled

Since

2.2

`OnOffSplit` (Enum)

An enumeration of three values: on, off, and split

Values

on Enabled

off Disabled

split Mixed

Since

2.6

String (Object)

A fat type wrapping 'str', to be embedded in lists.

Members

str: string Not documented

Since

1.2

StrOrNull (Alternate)

This is a string value or the explicit lack of a string (null pointer in C). Intended for cases when 'optional absent' already has a different meaning.

Members

s: string the string value

n: null no string value

Since

2.10

OffAutoPCIBAR (Enum)

An enumeration of options for specifying a PCI BAR

Values

- off** The specified feature is disabled
- auto** The PCI BAR for the feature is automatically selected
- bar0** PCI BAR0 is used for the feature
- bar1** PCI BAR1 is used for the feature
- bar2** PCI BAR2 is used for the feature
- bar3** PCI BAR3 is used for the feature
- bar4** PCI BAR4 is used for the feature
- bar5** PCI BAR5 is used for the feature

Since

2.12

PCIELinkSpeed (Enum)

An enumeration of PCIe link speeds in units of GT/s

Values

- 2_5** 2.5GT/s
- 5** 5.0GT/s
- 8** 8.0GT/s
- 16** 16.0GT/s

Since

4.0

PCIELinkWidth (Enum)

An enumeration of PCIe link width

Values

- 1** x1
- 2** x2
- 4** x4
- 8** x8
- 12** x12

16 x16

32 x32

Since

4.0

4.9.3 Cryptography

`QCryptoTLSCredsEndpoint` (Enum)

The type of network endpoint that will be using the credentials. Most types of credential require different setup / structures depending on whether they will be used in a server versus a client.

Values

client the network endpoint is acting as the client

server the network endpoint is acting as the server

Since

2.5

`QCryptoSecretFormat` (Enum)

The data format that the secret is provided in

Values

raw raw bytes. When encoded in JSON only valid UTF-8 sequences can be used

base64 arbitrary base64 encoded binary data

Since

2.6

`QCryptoHashAlgorithm` (Enum)

The supported algorithms for computing content digests

Values

md5 MD5. Should not be used in any new code, legacy compat only

sha1 SHA-1. Should not be used in any new code, legacy compat only

sha224 SHA-224. (since 2.7)

sha256 SHA-256. Current recommended strong hash.

sha384 SHA-384. (since 2.7)

sha512 SHA-512. (since 2.7)

ripemd160 RIPEMD-160. (since 2.7)

Since

2.6

QCryptoCipherAlgorithm (Enum)

The supported algorithms for content encryption ciphers

Values

aes-128 AES with 128 bit / 16 byte keys

aes-192 AES with 192 bit / 24 byte keys

aes-256 AES with 256 bit / 32 byte keys

des-rgb RFB specific variant of single DES. Do not use except in VNC.

3des 3DES(EDE) with 192 bit / 24 byte keys (since 2.9)

cast5-128 Cast5 with 128 bit / 16 byte keys

serpent-128 Serpent with 128 bit / 16 byte keys

serpent-192 Serpent with 192 bit / 24 byte keys

serpent-256 Serpent with 256 bit / 32 byte keys

twofish-128 Twofish with 128 bit / 16 byte keys

twofish-192 Twofish with 192 bit / 24 byte keys

twofish-256 Twofish with 256 bit / 32 byte keys

Since

2.6

QCryptoCipherMode (Enum)

The supported modes for content encryption ciphers

Values

ecb Electronic Code Book
cbc Cipher Block Chaining
xts XEX with tweaked code book and ciphertext stealing
ctr Counter (Since 2.8)

Since

2.6

QCryptoIVGenAlgorithm (Enum)

The supported algorithms for generating initialization vectors for full disk encryption. The 'plain' generator should not be used for disks with sector numbers larger than 2^{32} , except where compatibility with pre-existing Linux dm-crypt volumes is required.

Values

plain 64-bit sector number truncated to 32-bits
plain64 64-bit sector number
essiv 64-bit sector number encrypted with a hash of the encryption key

Since

2.6

QCryptoBlockFormat (Enum)

The supported full disk encryption formats

Values

qcow QCow/QCow2 built-in AES-CBC encryption. Use only for liberating data from old images.
luks LUKS encryption format. Recommended for new images

Since

2.6

QCryptoBlockOptionsBase (Object)

The common options that apply to all full disk encryption formats

Members

format: `QCryptoBlockFormat` the encryption format

Since

2.6

`QCryptoBlockOptionsQcow` (Object)

The options that apply to Qcow/Qcow2 AES-CBC encryption format

Members

key-secret: `string` (optional) the ID of a `QCryptoSecret` object providing the decryption key. Mandatory except when probing image for metadata only.

Since

2.6

`QCryptoBlockOptionsLUKS` (Object)

The options that apply to LUKS encryption format

Members

key-secret: `string` (optional) the ID of a `QCryptoSecret` object providing the decryption key. Mandatory except when probing image for metadata only.

Since

2.6

`QCryptoBlockCreateOptionsLUKS` (Object)

The options that apply to LUKS encryption format initialization

Members

cipher-alg: `QCryptoCipherAlgorithm` (optional) the cipher algorithm for data encryption. Currently defaults to 'aes-256'.

cipher-mode: `QCryptoCipherMode` (optional) the cipher mode for data encryption. Currently defaults to 'xts'.

ivgen-alg: `QCryptoIVGenAlgorithm` (optional) the initialization vector generator. Currently defaults to 'plain64'.

ivgen-hash-alg: `QCryptoHashAlgorithm` (optional) the initialization vector generator hash. Currently defaults to 'sha256'.

hash-alg: `QCryptoHashAlgorithm` (optional) the master key hash algorithm. Currently defaults to 'sha256'.

iter-time: `int` (optional) number of milliseconds to spend in PBKDF passphrase processing. Currently defaults to 2000. (since 2.8)

The members of `QCryptoBlockOptionsLUKS`

Since

2.6

`QCryptoBlockOpenOptions` (Object)

The options that are available for all encryption formats when opening an existing volume

Members

The members of `QCryptoBlockOptionsBase`

The members of `QCryptoBlockOptionsQcow` when format is "qcow"

The members of `QCryptoBlockOptionsLUKS` when format is "luks"

Since

2.6

`QCryptoBlockCreateOptions` (Object)

The options that are available for all encryption formats when initializing a new volume

Members

The members of `QCryptoBlockOptionsBase`

The members of `QCryptoBlockOptionsQcow` when format is "qcow"

The members of `QCryptoBlockCreateOptionsLUKS` when format is "luks"

Since

2.6

`QCryptoBlockInfoBase` (Object)

The common information that applies to all full disk encryption formats

Members

format: `QCryptoBlockFormat` the encryption format

Since

2.7

`QCryptoBlockInfoLUKSSlot` (Object)

Information about the LUKS block encryption key slot options

Members

active: `boolean` whether the key slot is currently in use

key-offset: `int` offset to the key material in bytes

iters: `int` (optional) number of PBKDF2 iterations for key material

stripes: `int` (optional) number of stripes for splitting key material

Since

2.7

`QCryptoBlockInfoLUKS` (Object)

Information about the LUKS block encryption options

Members

cipher-alg: `QCryptoCipherAlgorithm` the cipher algorithm for data encryption

cipher-mode: `QCryptoCipherMode` the cipher mode for data encryption

ivgen-alg: `QCryptoIVGenAlgorithm` the initialization vector generator

ivgen-hash-alg: `QCryptoHashAlgorithm` (optional) the initialization vector generator hash

hash-alg: `QCryptoHashAlgorithm` the master key hash algorithm

payload-offset: `int` offset to the payload data in bytes

master-key-iters: `int` number of PBKDF2 iterations for key material

uuid: `string` unique identifier for the volume

slots: `array of QCryptoBlockInfoLUKSSlot` information about each key slot

Since

2.7

QCryptoBlockInfo (Object)

Information about the block encryption options

Members

The members of `QCryptoBlockInfoBase`

The members of `QCryptoBlockInfoLUKS` when format is "luks"

Since

2.7

QCryptoBlockLUKSKeyslotState (Enum)

Defines state of keyslots that are affected by the update

Values

active The slots contain the given password and marked as active

inactive The slots are erased (contain garbage) and marked as inactive

Since

5.1

QCryptoBlockAmendOptionsLUKS (Object)

This struct defines the update parameters that activate/de-activate set of keyslots

Members

state: QCryptoBlockLUKSKeyslotState the desired state of the keyslots

new-secret: string (optional) The ID of a `QCryptoSecret` object providing the password to be written into added active keyslots

old-secret: string (optional) Optional (for deactivation only) If given will deactivate all keyslots that match password located in `QCryptoSecret` with this ID

iter-time: int (optional) Optional (for activation only) Number of milliseconds to spend in PBKDF passphrase processing for the newly activated keyslot. Currently defaults to 2000.

keyslot: int (optional) Optional. ID of the keyslot to activate/deactivate. For keyslot activation, keyslot should not be active already (this is unsafe to update an active keyslot), but possible if 'force' parameter is given. If keyslot is not given, first free keyslot will be written.

For keyslot deactivation, this parameter specifies the exact keyslot to deactivate

secret: string (optional) Optional. The ID of a QCryptoSecret object providing the password to use to retrieve current master key. Defaults to the same secret that was used to open the image

Since 5.1

QCryptoBlockAmendOptions (Object)

The options that are available for all encryption formats when amending encryption settings

Members

The members of QCryptoBlockOptionsBase

The members of QCryptoBlockAmendOptionsLUKS when format is "luks"

Since

5.1

Background jobs

JobType (Enum)

Type of a background job.

Values

commit block commit job type, see “block-commit”

stream block stream job type, see “block-stream”

mirror drive mirror job type, see “drive-mirror”

backup drive backup job type, see “drive-backup”

create image creation job type, see “blockdev-create” (since 3.0)

amend image options amend job type, see “x-blockdev-amend” (since 5.1)

snapshot-load snapshot load job type, see “snapshot-load” (since 6.0)

snapshot-save snapshot save job type, see “snapshot-save” (since 6.0)

snapshot-delete snapshot delete job type, see “snapshot-delete” (since 6.0)

Since

1.7

JobStatus (Enum)

Indicates the present state of a given job in its lifetime.

Values

undefined Erroneous, default state. Should not ever be visible.

created The job has been created, but not yet started.

running The job is currently running.

paused The job is running, but paused. The pause may be requested by either the QMP user or by internal processes.

ready The job is running, but is ready for the user to signal completion. This is used for long-running jobs like mirror that are designed to run indefinitely.

standby The job is ready, but paused. This is nearly identical to `paused`. The job may return to `ready` or otherwise be canceled.

waiting The job is waiting for other jobs in the transaction to converge to the waiting state. This status will likely not be visible for the last job in a transaction.

pending The job has finished its work, but has finalization steps that it needs to make prior to completing. These changes will require manual intervention via `job-finalize` if `auto-finalize` was set to false. These pending changes may still fail.

aborting The job is in the process of being aborted, and will finish with an error. The job will afterwards report that it is `concluded`. This status may not be visible to the management process.

concluded The job has finished all work. If `auto-dismiss` was set to false, the job will remain in the query list until it is dismissed via `job-dismiss`.

null The job is in the process of being dismantled. This state should not ever be visible externally.

Since

2.12

JobVerb (Enum)

Represents command verbs that can be applied to a job.

Values

cancel see `job-cancel`

pause see `job-pause`

resume see `job-resume`

set-speed see `block-job-set-speed`

complete see `job-complete`

dismiss see `job-dismiss`

finalize see `job-finalize`

Since

2.12

JOB_STATUS_CHANGE (Event)

Emitted when a job transitions to a different status.

Arguments

id: string The job identifier

status: JobStatus The new job status

Since

3.0

job-pause (Command)

Pause an active job.

This command returns immediately after marking the active job for pausing. Pausing an already paused job is an error.

The job will pause as soon as possible, which means transitioning into the PAUSED state if it was RUNNING, or into STANDBY if it was READY. The corresponding JOB_STATUS_CHANGE event will be emitted.

Cancelling a paused job automatically resumes it.

Arguments

id: string The job identifier.

Since

3.0

job-resume (Command)

Resume a paused job.

This command returns immediately after resuming a paused job. Resuming an already running job is an error.

id : The job identifier.

Arguments

id: string Not documented

Since

3.0

job-cancel (Command)

Instruct an active background job to cancel at the next opportunity. This command returns immediately after marking the active job for cancellation.

The job will cancel as soon as possible and then emit a `JOB_STATUS_CHANGE` event. Usually, the status will change to `ABORTING`, but it is possible that a job successfully completes (e.g. because it was almost done and there was no opportunity to cancel earlier than completing the job) and transitions to `PENDING` instead.

Arguments

id: string The job identifier.

Since

3.0

job-complete (Command)

Manually trigger completion of an active job in the `READY` state.

Arguments

id: string The job identifier.

Since

3.0

job-dismiss (Command)

Deletes a job that is in the `CONCLUDED` state. This command only needs to be run explicitly for jobs that don't have automatic dismiss enabled.

This command will refuse to operate on any job that has not yet reached its terminal state, `JOB_STATUS_CONCLUDED`. For jobs that make use of `JOB_READY` event, `job-cancel` or `job-complete` will still need to be used as appropriate.

Arguments

id: string The job identifier.

Since

3.0

`job-finalize` (Command)

Instructs all jobs in a transaction (or a single job if it is not part of any transaction) to finalize any graph changes and do any necessary cleanup. This command requires that all involved jobs are in the PENDING state.

For jobs in a transaction, instructing one job to finalize will force ALL jobs in the transaction to finalize, so it is only necessary to instruct a single member job to finalize.

Arguments

id: string The identifier of any job in the transaction, or of a job that is not part of any transaction.

Since

3.0

`JobInfo` (Object)

Information about a job.

Members

id: string The job identifier

type: JobType The kind of job that is being performed

status: JobStatus Current job state/status

current-progress: int Progress made until now. The unit is arbitrary and the value can only meaningfully be used for the ratio of `current-progress` to `total-progress`. The value is monotonically increasing.

total-progress: int Estimated `current-progress` value at the completion of the job. This value can arbitrarily change while the job is running, in both directions.

error: string (optional) If this field is present, the job failed; if it is still missing in the CONCLUDED state, this indicates successful completion.

The value is a human-readable error message to describe the reason for the job failure. It should not be parsed by applications.

Since

3.0

`query-jobs` (Command)

Return information about jobs.

Returns

a list with a `JobInfo` for each active job

Since

3.0

4.9.4 Socket data types

`NetworkAddressFamily` (Enum)

The network address family

Values

ipv4 IPV4 family

ipv6 IPV6 family

unix unix socket

vsock vsock family (since 2.8)

unknown otherwise

Since

2.1

`InetSocketAddressBase` (Object)

Members

host: string host part of the address

port: string port part of the address

`InetSocketAddress` (Object)

Captures a socket address or address range in the Internet namespace.

Members

numeric: boolean (optional) true if the host/port are guaranteed to be numeric, false if name resolution should be attempted. Defaults to false. (Since 2.9)

to: int (optional) If present, this is range of possible addresses, with port between `port` and `to`.

ipv4: boolean (optional) whether to accept IPv4 addresses, default try both IPv4 and IPv6

ipv6: boolean (optional) whether to accept IPv6 addresses, default try both IPv4 and IPv6

keep-alive: boolean (optional) enable keep-alive when connecting to this socket. Not supported for passive sockets. (Since 4.2)

The members of InetSocketAddressBase

Since

1.3

UnixSocketAddress (Object)

Captures a socket address in the local (“Unix socket”) namespace.

Members

path: string filesystem path to use

abstract: boolean (optional) (If: defined(CONFIG_LINUX)) if true, this is a Linux abstract socket address. `path` will be prefixed by a null byte, and optionally padded with null bytes. Defaults to false. (Since 5.1)

tight: boolean (optional) (If: defined(CONFIG_LINUX)) if false, pad an abstract socket address with enough null bytes to make it fill struct sockaddr_un member `sun_path`. Defaults to true. (Since 5.1)

Since

1.3

VsockSocketAddress (Object)

Captures a socket address in the vsock namespace.

Members

cid: string unique host identifier

port: string port

Note

string types are used to allow for possible future hostname or service resolution support.

Since

2.8

SocketAddressLegacy (Object)

Captures the address of a socket, which could also be a named file descriptor

Members

type One of `inet`, `unix`, `vsock`, `fd`

data: `InetSocketAddress` when **type** is `"inet"`

data: `UnixSocketAddress` when **type** is `"unix"`

data: `VsockSocketAddress` when **type** is `"vsock"`

data: `String` when **type** is `"fd"`

Note

This type is deprecated in favor of `SocketAddress`. The difference between `SocketAddressLegacy` and `SocketAddress` is that the latter is a flat union rather than a simple union. Flat is nicer because it avoids nesting on the wire, i.e. that form has fewer `{}`.

Since

1.3

SocketAddressType (Enum)

Available `SocketAddress` types

Values

inet Internet address

unix Unix domain socket

vsock VMCI address

fd decimal is for file descriptor number, otherwise a file descriptor name. Named file descriptors are permitted in monitor commands, in combination with the ‘getfd’ command. Decimal file descriptors are permitted at startup or other contexts where no monitor context is active.

Since

2.9

SocketAddress (Object)

Captures the address of a socket, which could also be a named file descriptor

Members

type: SocketAddressType Transport type

The members of **InetSocketAddress** when **type** is "inet"

The members of **UnixSocketAddress** when **type** is "unix"

The members of **VsockSocketAddress** when **type** is "vsock"

The members of **String** when **type** is "fd"

Since

2.9

SnapshotInfo (Object)

Members

id: string unique snapshot id

name: string user chosen name

vm-state-size: int size of the VM state

date-sec: int UTC date of the snapshot in seconds

date-nsec: int fractional part in nano seconds to be used with date-sec

vm-clock-sec: int VM clock relative to boot in seconds

vm-clock-nsec: int fractional part in nano seconds to be used with vm-clock-sec

icount: int (optional) Current instruction count. Appears when execution record/replay is enabled. Used for “time-traveling” to match the moment in the recorded execution with the snapshots. This counter may be obtained through `query-replay` command (since 5.2)

Since

1.3

ImageInfoSpecificQcow2EncryptionBase (Object)

Members

format: BlockdevQcow2EncryptionFormat The encryption format

Since

2.10

ImageInfoSpecificQCow2Encryption (Object)

Members

The members of `ImageInfoSpecificQCow2EncryptionBase`

The members of `QCryptoBlockInfoLUKS` when format is "luks"

Since

2.10

ImageInfoSpecificQCow2 (Object)

Members

compat: string compatibility level

data-file: string (optional) the filename of the external data file that is stored in the image and used as a default for opening the image (since: 4.0)

data-file-raw: boolean (optional) True if the external data file must stay valid as a standalone (read-only) raw image without looking at qcow2 metadata (since: 4.0)

extended-l2: boolean (optional) true if the image has extended L2 entries; only valid for compat \geq 1.1 (since 5.2)

lazy-refcounts: boolean (optional) on or off; only valid for compat \geq 1.1

corrupt: boolean (optional) true if the image has been marked corrupt; only valid for compat \geq 1.1 (since 2.2)

refcount-bits: int width of a refcount entry in bits (since 2.3)

encrypt: ImageInfoSpecificQCow2Encryption (optional) details about encryption parameters; only set if image is encrypted (since 2.10)

bitmaps: array of Qcow2BitmapInfo (optional) A list of qcow2 bitmap details (since 4.0)

compression-type: Qcow2CompressionType the image cluster compression method (since 5.1)

Since

1.7

ImageInfoSpecificVmdk (Object)

Members

create-type: string The create type of VMDK image

cid: int Content id of image

parent-cid: int Parent VMDK image's cid

extents: array of ImageInfo List of extent files

Since

1.7

ImageInfoSpecific (Object)

A discriminated record of image format specific information structures.

Members

type One of qcow2, vmdk, luks

data: ImageInfoSpecificQCow2 when type is "qcow2"

data: ImageInfoSpecificVmdk when type is "vmdk"

data: QCryptoBlockInfoLUKS when type is "luks"

Since

1.7

ImageInfo (Object)

Information about a QEMU image file

Members

filename: string name of the image file

format: string format of the image file

virtual-size: int maximum capacity in bytes of the image

actual-size: int (optional) actual size on disk in bytes of the image

dirty-flag: boolean (optional) true if image is not cleanly closed

cluster-size: int (optional) size of a cluster in bytes

encrypted: boolean (optional) true if the image is encrypted

compressed: boolean (optional) true if the image is compressed (Since 1.7)

backing-filename: string (optional) name of the backing file

full-backing-filename: string (optional) full path of the backing file

backing-filename-format: string (optional) the format of the backing file

snapshots: array of SnapshotInfo (optional) list of VM snapshots

backing-image: ImageInfo (optional) info of the backing image (since 1.6)

format-specific: ImageInfoSpecific (optional) structure supplying additional format-specific information (since 1.7)

Since

1.3

ImageCheck (Object)

Information about a QEMU image file check

Members

filename: string name of the image file checked

format: string format of the image file checked

check-errors: int number of unexpected errors occurred during check

image-end-offset: int (optional) offset (in bytes) where the image ends, this field is present if the driver for the image format supports it

corruptions: int (optional) number of corruptions found during the check if any

leaks: int (optional) number of leaks found during the check if any

corruptions-fixed: int (optional) number of corruptions fixed during the check if any

leaks-fixed: int (optional) number of leaks fixed during the check if any

total-clusters: int (optional) total number of clusters, this field is present if the driver for the image format supports it

allocated-clusters: int (optional) total number of allocated clusters, this field is present if the driver for the image format supports it

fragmented-clusters: int (optional) total number of fragmented clusters, this field is present if the driver for the image format supports it

compressed-clusters: int (optional) total number of compressed clusters, this field is present if the driver for the image format supports it

Since

1.4

MapEntry (Object)

Mapping information from a virtual block range to a host file range

Members

start: int virtual (guest) offset of the first byte described by this entry

length: int the number of bytes of the mapped virtual range

data: boolean reading the image will actually read data from a file (in particular, if `offset` is present this means that the sectors are not simply preallocated, but contain actual data in raw format)

zero: boolean whether the virtual blocks read as zeroes

depth: int number of layers (0 = top image, 1 = top image's backing file, ..., n - 1 = bottom image (where n is the number of images in the chain)) before reaching one for which the range is allocated

offset: int (optional) if present, the image file stores the data for this range in raw format at the given (host) offset

filename: string (optional) filename that is referred to by `offset`

Since

2.6

BlockdevCacheInfo (Object)

Cache mode information for a block device

Members

writeback: boolean true if writeback mode is enabled

direct: boolean true if the host page cache is bypassed (O_DIRECT)

no-flush: boolean true if flush requests are ignored for the device

Since

2.3

BlockDeviceInfo (Object)

Information about the backing device for a block device.

Members

file: string the filename of the backing device

node-name: string (optional) the name of the block driver node (Since 2.0)

ro: boolean true if the backing device was open read-only

drv: string the name of the block format used to open the backing device. As of 0.14 this can be: 'blkdebug', 'bochs', 'cloop', 'cow', 'dmg', 'file', 'file', 'ftp', 'ftps', 'host_cdrom', 'host_device', 'http', 'https', 'luks', 'nbd', 'parallels', 'qcow', 'qcow2', 'raw', 'vdi', 'vmdk', 'vpc', 'vvfat' 2.2: 'archipelago' added, 'cow' dropped 2.3: 'host_floppy' deprecated 2.5: 'host_floppy' dropped 2.6: 'luks' added 2.8: 'replication' added, 'tftp' dropped 2.9: 'archipelago' dropped

backing_file: string (optional) the name of the backing file (for copy-on-write)

backing_file_depth: int number of files in the backing file chain (since: 1.2)

encrypted: boolean true if the backing device is encrypted

encryption_key_missing: boolean always false

detect_zeroes: BlockdevDetectZeroesOptions detect and optimize zero writes (Since 2.1)

bps: int total throughput limit in bytes per second is specified

bps_rd: int read throughput limit in bytes per second is specified

bps_wr: int write throughput limit in bytes per second is specified

iops: int total I/O operations per second is specified

iops_rd: int read I/O operations per second is specified

iops_wr: int write I/O operations per second is specified

image: ImageInfo the info of image used (since: 1.6)

bps_max: int (optional)

total throughput limit during bursts, in bytes (Since 1.7)

bps_rd_max: int (optional)

read throughput limit during bursts, in bytes (Since 1.7)

bps_wr_max: int (optional)

write throughput limit during bursts, in bytes (Since 1.7)

iops_max: int (optional)

total I/O operations per second during bursts, in bytes (Since 1.7)

iops_rd_max: int (optional)

read I/O operations per second during bursts, in bytes (Since 1.7)

iops_wr_max: int (optional)

write I/O operations per second during bursts, in bytes (Since 1.7)

bps_max_length: int (optional)

maximum length of the **bps_max** burst period, in seconds. (Since 2.6)

bps_rd_max_length: int (optional)

maximum length of the **bps_rd_max** burst period, in seconds. (Since 2.6)

bps_wr_max_length: int (optional)

maximum length of the **bps_wr_max** burst period, in seconds. (Since 2.6)

iops_max_length: int (optional)

maximum length of the **iops** burst period, in seconds. (Since 2.6)

iops_rd_max_length: int (optional)

maximum length of the **iops_rd_max** burst period, in seconds. (Since 2.6)

iops_wr_max_length: int (optional)

maximum length of the **iops_wr_max** burst period, in seconds. (Since 2.6)

iops_size: int (optional) an I/O size in bytes (Since 1.7)

group: string (optional) throttle group name (Since 2.4)

cache: BlockdevCacheInfo the cache mode used for the block device (since: 2.3)

write_threshold: int configured write threshold for the device. 0 if disabled. (Since 2.3)

dirty-bitmaps: array of **BlockDirtyInfo** (optional) dirty bitmaps information (only present if node has one or more dirty bitmaps) (Since 4.2)

Features

deprecated Member `encryption_key_missing` is deprecated. It is always false.

Since

0.14

BlockDeviceIoStatus (Enum)

An enumeration of block device I/O status.

Values

ok The last I/O operation has succeeded

failed The last I/O operation has failed

nospace The last I/O operation has failed due to a no-space condition

Since

1.0

DirtyBitmapStatus (Enum)

An enumeration of possible states that a dirty bitmap can report to the user.

Values

frozen The bitmap is currently in-use by some operation and is immutable. If the bitmap was `active` prior to the operation, new writes by the guest are being recorded in a temporary buffer, and will not be lost. Generally, bitmaps are cleared on successful use in an operation and the temporary buffer is committed into the bitmap. On failure, the temporary buffer is merged back into the bitmap without first clearing it. Please refer to the documentation for each bitmap-using operation, See also `blockdev-backup`, `drive-backup`.

disabled The bitmap is not currently recording new writes by the guest. This is requested explicitly via `block-dirty-bitmap-disable`. It can still be cleared, deleted, or used for backup operations.

active The bitmap is actively monitoring for new writes, and can be cleared, deleted, or used for backup operations.

locked The bitmap is currently in-use by some operation and is immutable. If the bitmap was `active` prior to the operation, it is still recording new writes. If the bitmap was `disabled`, it is not recording new writes. (Since 2.12)

inconsistent This is a persistent dirty bitmap that was marked in-use on disk, and is unusable by QEMU. It can only be deleted. Please rely on the `inconsistent` field in `BlockDirtyInfo` instead, as the `status` field is deprecated. (Since 4.0)

Since

2.4

BlockDirtyInfo (Object)

Block dirty bitmap information.

Members

name: string (optional) the name of the dirty bitmap (Since 2.4)

count: int number of dirty bytes according to the dirty bitmap

granularity: int granularity of the dirty bitmap in bytes (since 1.4)

status: DirtyBitmapStatus current status of the dirty bitmap (since 2.4)

recording: boolean true if the bitmap is recording new writes from the guest. Replaces *active* and *disabled* statuses. (since 4.0)

busy: boolean true if the bitmap is in-use by some operation (NBD or jobs) and cannot be modified via QMP or used by another operation. Replaces *locked* and *frozen* statuses. (since 4.0)

persistent: boolean true if the bitmap was stored on disk, is scheduled to be stored on disk, or both. (since 4.0)

inconsistent: boolean (optional) true if this is a persistent bitmap that was improperly stored. Implies *persistent* to be true; *recording* and *busy* to be false. This bitmap cannot be used. To remove it, use `block-dirty-bitmap-remove`. (Since 4.0)

Features

deprecated Member `status` is deprecated. Use `recording` and `locked` instead.

Since

1.3

Qcow2BitmapInfoFlags (Enum)

An enumeration of flags that a bitmap can report to the user.

Values

in-use This flag is set by any process actively modifying the qcow2 file, and cleared when the updated bitmap is flushed to the qcow2 image. The presence of this flag in an offline image means that the bitmap was not saved correctly after its last usage, and may contain inconsistent data.

auto The bitmap must reflect all changes of the virtual disk by any application that would write to this qcow2 file.

Since

4.0

Qcow2BitmapInfo (Object)

Qcow2 bitmap information.

Members

name: `string` the name of the bitmap

granularity: `int` granularity of the bitmap in bytes

flags: `array of Qcow2BitmapInfoFlags` flags of the bitmap

Since

4.0

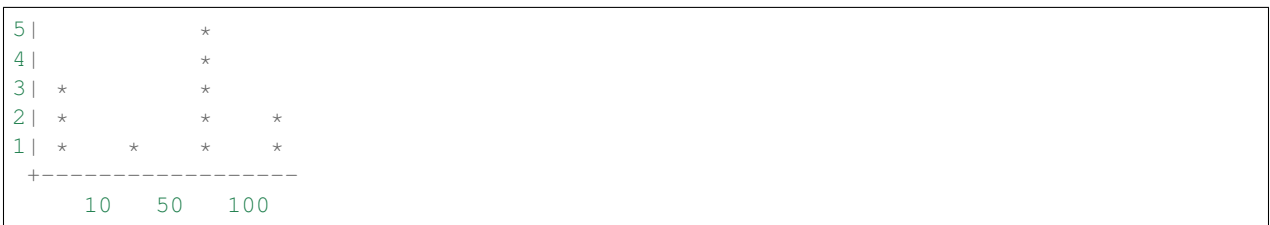
BlockLatencyHistogramInfo (Object)

Block latency histogram.

Members

boundaries: `array of int` list of interval boundary values in nanoseconds, all greater than zero and in ascending order. For example, the list [10, 50, 100] produces the following histogram intervals: [0, 10), [10, 50), [50, 100), [100, +inf).

bins: `array of int` list of io request counts corresponding to histogram intervals. `len(bins) = len(boundaries) + 1` For the example above, `bins` may be something like [3, 1, 5, 2], and corresponding histogram looks like:



Since

4.0

BlockInfo (Object)

Block device information. This structure describes a virtual device and the backing device associated with it.

Members

device: string The device name associated with the virtual device.

qdev: string (optional) The qdev ID, or if no ID is assigned, the QOM path of the block device. (since 2.10)

type: string This field is returned only for compatibility reasons, it should not be used (always returns ‘unknown’)

removable: boolean True if the device supports removable media.

locked: boolean True if the guest has locked this device from having its media removed

tray_open: boolean (optional) True if the device’s tray is open (only present if it has a tray)

dirty-bitmaps: array of BlockDirtyInfo (optional) dirty bitmaps information (only present if the driver has one or more dirty bitmaps) (Since 2.0)

io-status: BlockDeviceIoStatus (optional) BlockDeviceIoStatus. Only present if the device supports it and the VM is configured to stop on errors (supported device models: virtio-blk, IDE, SCSI except scsi-generic)

inserted: BlockDeviceInfo (optional) BlockDeviceInfo describing the device if media is present

Features

deprecated Member `dirty-bitmaps` is deprecated. Use `inserted` member `dirty-bitmaps` instead.

Since

0.14

BlockMeasureInfo (Object)

Image file size calculation information. This structure describes the size requirements for creating a new image file.

The size requirements depend on the new image file format. File size always equals virtual disk size for the ‘raw’ format, even for sparse POSIX files. Compact formats such as ‘qcow2’ represent unallocated and zero regions efficiently so file size may be smaller than virtual disk size.

The values are upper bounds that are guaranteed to fit the new image file. Subsequent modification, such as internal snapshot or further bitmap creation, may require additional space and is not covered here.

Members

required: int Size required for a new image file, in bytes, when copying just allocated guest-visible contents.

fully-allocated: int Image file size, in bytes, once data has been written to all sectors, when copying just guest-visible contents.

bitmaps: int (optional) Additional size required if all the top-level bitmap metadata in the source image were to be copied to the destination, present only when source and destination both support persistent bitmaps. (since 5.1)

Since

2.10

query-block (Command)

Get a list of BlockInfo for all virtual block devices.

Returns

a list of BlockInfo describing each virtual block device. Filter nodes that were created implicitly are skipped over.

Since

0.14

Example

```
-> { "execute": "query-block" }
<- {
  "return":[
    {
      "io-status": "ok",
      "device":"ide0-hd0",
      "locked":false,
      "removable":false,
      "inserted":{
        "ro":false,
        "drv":"qcow2",
        "encrypted":false,
        "file":"disks/test.qcow2",
        "backing_file_depth":1,
        "bps":1000000,
        "bps_rd":0,
        "bps_wr":0,
        "iops":1000000,
        "iops_rd":0,
        "iops_wr":0,
        "bps_max": 8000000,
        "bps_rd_max": 0,
        "bps_wr_max": 0,
        "iops_max": 0,
        "iops_rd_max": 0,
        "iops_wr_max": 0,
        "iops_size": 0,
        "detect_zeroes": "on",
        "write_threshold": 0,
        "image":{
          "filename":"disks/test.qcow2",
          "format":"qcow2",
          "virtual-size":2048000,
          "backing_file":"base.qcow2",
```

(continues on next page)

(continued from previous page)

```

        "full-backing-filename": "disks/base.qcow2",
        "backing-filename-format": "qcow2",
        "snapshots": [
            {
                "id": "1",
                "name": "snapshot1",
                "vm-state-size": 0,
                "date-sec": 10000200,
                "date-nsec": 12,
                "vm-clock-sec": 206,
                "vm-clock-nsec": 30
            }
        ],
        "backing-image": {
            "filename": "disks/base.qcow2",
            "format": "qcow2",
            "virtual-size": 2048000
        }
    },
    "qdev": "ide_disk",
    "type": "unknown"
},
{
    "io-status": "ok",
    "device": "ide1-cd0",
    "locked": false,
    "removable": true,
    "qdev": "/machine/unattached/device[23]",
    "tray_open": false,
    "type": "unknown"
},
{
    "device": "floppy0",
    "locked": false,
    "removable": true,
    "qdev": "/machine/unattached/device[20]",
    "type": "unknown"
},
{
    "device": "sd0",
    "locked": false,
    "removable": true,
    "type": "unknown"
}
]
}

```

BlockDeviceTimedStats (Object)

Statistics of a block device during a given interval of time.

Members

interval_length: int Interval used for calculating the statistics, in seconds.

min_rd_latency_ns: int Minimum latency of read operations in the defined interval, in nanoseconds.

min_wr_latency_ns: int Minimum latency of write operations in the defined interval, in nanoseconds.

min_flush_latency_ns: int Minimum latency of flush operations in the defined interval, in nanoseconds.

max_rd_latency_ns: int Maximum latency of read operations in the defined interval, in nanoseconds.

max_wr_latency_ns: int Maximum latency of write operations in the defined interval, in nanoseconds.

max_flush_latency_ns: int Maximum latency of flush operations in the defined interval, in nanoseconds.

avg_rd_latency_ns: int Average latency of read operations in the defined interval, in nanoseconds.

avg_wr_latency_ns: int Average latency of write operations in the defined interval, in nanoseconds.

avg_flush_latency_ns: int Average latency of flush operations in the defined interval, in nanoseconds.

avg_rd_queue_depth: number Average number of pending read operations in the defined interval.

avg_wr_queue_depth: number Average number of pending write operations in the defined interval.

Since

2.5

BlockDeviceStats (Object)

Statistics of a virtual block device or a block backing device.

Members

rd_bytes: int The number of bytes read by the device.

wr_bytes: int The number of bytes written by the device.

unmap_bytes: int The number of bytes unmapped by the device (Since 4.2)

rd_operations: int The number of read operations performed by the device.

wr_operations: int The number of write operations performed by the device.

flush_operations: int The number of cache flush operations performed by the device (since 0.15)

unmap_operations: int The number of unmap operations performed by the device (Since 4.2)

rd_total_time_ns: int Total time spent on reads in nanoseconds (since 0.15).

wr_total_time_ns: int Total time spent on writes in nanoseconds (since 0.15).

flush_total_time_ns: int Total time spent on cache flushes in nanoseconds (since 0.15).

unmap_total_time_ns: int Total time spent on unmap operations in nanoseconds (Since 4.2)

wr_highest_offset: int The offset after the greatest byte written to the device. The intended use of this information is for growable sparse files (like qcow2) that are used on top of a physical device.

rd_merged: int Number of read requests that have been merged into another request (Since 2.3).

wr_merged: int Number of write requests that have been merged into another request (Since 2.3).

unmap_merged: int Number of unmap requests that have been merged into another request (Since 4.2)

idle_time_ns: int (optional) Time since the last I/O operation, in nanoseconds. If the field is absent it means that there haven't been any operations yet (Since 2.5).

failed_rd_operations: int The number of failed read operations performed by the device (Since 2.5)

failed_wr_operations: int The number of failed write operations performed by the device (Since 2.5)

failed_flush_operations: int The number of failed flush operations performed by the device (Since 2.5)

failed_unmap_operations: int The number of failed unmap operations performed by the device (Since 4.2)

invalid_rd_operations: int
The number of invalid read operations performed by the device (Since 2.5)

invalid_wr_operations: int The number of invalid write operations performed by the device (Since 2.5)

invalid_flush_operations: int The number of invalid flush operations performed by the device (Since 2.5)

invalid_unmap_operations: int The number of invalid unmap operations performed by the device (Since 4.2)

account_invalid: boolean Whether invalid operations are included in the last access statistics (Since 2.5)

account_failed: boolean Whether failed operations are included in the latency and last access statistics (Since 2.5)

timed_stats: array of BlockDeviceTimedStats Statistics specific to the set of previously defined intervals of time (Since 2.5)

rd_latency_histogram: BlockLatencyHistogramInfo (optional) BlockLatencyHistogramInfo. (Since 4.0)

wr_latency_histogram: BlockLatencyHistogramInfo (optional) BlockLatencyHistogramInfo. (Since 4.0)

flush_latency_histogram: BlockLatencyHistogramInfo (optional) BlockLatencyHistogramInfo. (Since 4.0)

Since

0.14

BlockStatsSpecificFile (Object)

File driver statistics

Members

discard-nb-ok: int The number of successful discard operations performed by the driver.

discard-nb-failed: int The number of failed discard operations performed by the driver.

discard-bytes-ok: int The number of bytes discarded by the driver.

Since

4.2

BlockStatsSpecificNvme (Object)

NVMe driver statistics

Members

completion-errors: int The number of completion errors.

aligned-accesses: int The number of aligned accesses performed by the driver.

unaligned-accesses: int The number of unaligned accesses performed by the driver.

Since

5.2

BlockStatsSpecific (Object)

Block driver specific statistics

Members

driver: BlockdevDriver Not documented

The members of **BlockStatsSpecificFile** when **driver** is "file"

The members of **BlockStatsSpecificFile** when **driver** is "host_device"

The members of **BlockStatsSpecificNvme** when **driver** is "nvme"

Since

4.2

BlockStats (Object)

Statistics of a virtual block device or a block backing device.

Members

device: string (optional) If the stats are for a virtual block device, the name corresponding to the virtual block device.

node-name: string (optional) The node name of the device. (Since 2.3)

qdev: string (optional) The qdev ID, or if no ID is assigned, the QOM path of the block device. (since 3.0)

stats: BlockDeviceStats A BlockDeviceStats for the device.

driver-specific: BlockStatsSpecific (optional) Optional driver-specific stats. (Since 4.2)

parent: BlockStats (optional) This describes the file block device if it has one. Contains recursively the statistics of the underlying protocol (e.g. the host file for a qcow2 image). If there is no underlying protocol, this field is omitted

backing: BlockStats (optional) This describes the backing block device if it has one. (Since 2.0)

Since

0.14

query-blockstats (Command)

Query the BlockStats for all virtual block devices.

Arguments

query-nodes: boolean (optional) If true, the command will query all the block nodes that have a node name, in a list which will include “parent” information, but not “backing”. If false or omitted, the behavior is as before - query all the device backends, recursively including their “parent” and “backing”. Filter nodes that were created implicitly are skipped over in this mode. (Since 2.3)

Returns

A list of BlockStats for each virtual block devices.

Since

0.14

Example

```
-> { "execute": "query-blockstats" }
<- {
  "return": [
    {
      "device": "ide0-hd0",
      "parent": {
        "stats": {
          "wr_highest_offset": 3686448128,
          "wr_bytes": 9786368,
          "wr_operations": 751,
          "rd_bytes": 122567168,
          "rd_operations": 36772,
          "wr_total_times_ns": 313253456,
          "rd_total_times_ns": 3465673657,
          "flush_total_times_ns": 49653
        }
      }
    }
  ]
}
```

(continues on next page)

(continued from previous page)

```

        "flush_operations":61,
        "rd_merged":0,
        "wr_merged":0,
        "idle_time_ns":2953431879,
        "account_invalid":true,
        "account_failed":false
    }
},
"stats":{
    "wr_highest_offset":2821110784,
    "wr_bytes":9786368,
    "wr_operations":692,
    "rd_bytes":122739200,
    "rd_operations":36604
    "flush_operations":51,
    "wr_total_times_ns":313253456
    "rd_total_times_ns":3465673657
    "flush_total_times_ns":49653,
    "rd_merged":0,
    "wr_merged":0,
    "idle_time_ns":2953431879,
    "account_invalid":true,
    "account_failed":false
},
"qdev": "/machine/unattached/device[23]"
},
{
    "device":"ide1-cd0",
    "stats":{
        "wr_highest_offset":0,
        "wr_bytes":0,
        "wr_operations":0,
        "rd_bytes":0,
        "rd_operations":0
        "flush_operations":0,
        "wr_total_times_ns":0
        "rd_total_times_ns":0
        "flush_total_times_ns":0,
        "rd_merged":0,
        "wr_merged":0,
        "account_invalid":false,
        "account_failed":false
    },
    "qdev": "/machine/unattached/device[24]"
},
{
    "device":"floppy0",
    "stats":{
        "wr_highest_offset":0,
        "wr_bytes":0,
        "wr_operations":0,
        "rd_bytes":0,
        "rd_operations":0
        "flush_operations":0,
        "wr_total_times_ns":0
        "rd_total_times_ns":0
        "flush_total_times_ns":0,

```

(continues on next page)

(continued from previous page)

```

        "rd_merged":0,
        "wr_merged":0,
        "account_invalid":false,
        "account_failed":false
    },
    "qdev": "/machine/unattached/device[16]"
},
{
    "device":"sd0",
    "stats":{
        "wr_highest_offset":0,
        "wr_bytes":0,
        "wr_operations":0,
        "rd_bytes":0,
        "rd_operations":0,
        "flush_operations":0,
        "wr_total_times_ns":0,
        "rd_total_times_ns":0,
        "flush_total_times_ns":0,
        "rd_merged":0,
        "wr_merged":0,
        "account_invalid":false,
        "account_failed":false
    }
}
]
}

```

BlockdevOnError (Enum)

An enumeration of possible behaviors for errors on I/O operations. The exact meaning depends on whether the I/O was initiated by a guest or by a block job

Values

report for guest operations, report the error to the guest; for jobs, cancel the job

ignore ignore the error, only report a QMP event (BLOCK_IO_ERROR or BLOCK_JOB_ERROR). The backup, mirror and commit block jobs retry the failing request later and may still complete successfully. The stream block job continues to stream and will complete with an error.

enospc same as **stop** on ENOSPC, same as **report** otherwise.

stop for guest operations, stop the virtual machine; for jobs, pause the job

auto inherit the error handling policy of the backend (since: 2.7)

Since

1.3

MirrorSyncMode (Enum)

An enumeration of possible behaviors for the initial synchronization phase of storage mirroring.

Values

top copies data in the topmost image to the destination

full copies data from all images to the destination

none only copy data written from now on

incremental only copy data described by the dirty bitmap. (since: 2.4)

bitmap only copy data described by the dirty bitmap. (since: 4.2) Behavior on completion is determined by the `BitmapSyncMode`.

Since

1.3

BitmapSyncMode (Enum)

An enumeration of possible behaviors for the synchronization of a bitmap when used for data copy operations.

Values

on-success The bitmap is only synced when the operation is successful. This is the behavior always used for ‘INCREMENTAL’ backups.

never The bitmap is never synchronized with the operation, and is treated solely as a read-only manifest of blocks to copy.

always The bitmap is always synchronized with the operation, regardless of whether or not the operation was successful.

Since

4.2

MirrorCopyMode (Enum)

An enumeration whose values tell the mirror block job when to trigger writes to the target.

Values

background copy data in background only.

write-blocking when data is written to the source, write it (synchronously) to the target as well. In addition, data is copied in background just like in `background` mode.

Since

3.0

BlockJobInfo (Object)

Information about a long-running block device operation.

Members

type: string the job type ('stream' for image streaming)

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: int Estimated `offset` value at the completion of the job. This value can arbitrarily change while the job is running, in both directions.

offset: int Progress made until now. The unit is arbitrary and the value can only meaningfully be used for the ratio of `offset` to `len`. The value is monotonically increasing.

busy: boolean false if the job is known to be in a quiescent state, with no pending I/O. Since 1.3.

paused: boolean whether the job is paused or, if `busy` is true, will pause itself as soon as possible. Since 1.3.

speed: int the rate limit, bytes per second

io-status: BlockDeviceIoStatus the status of the job (since 1.3)

ready: boolean true if the job may be completed (since 2.2)

status: JobStatus Current job state/status (since 2.12)

auto-finalize: boolean Job will finalize itself when PENDING, moving to the CONCLUDED state. (since 2.12)

auto-dismiss: boolean Job will dismiss itself when CONCLUDED, moving to the NULL state and disappearing from the query list. (since 2.12)

error: string (optional) Error information if the job did not complete successfully. Not set if the job completed successfully. (since 2.12.1)

Since

1.1

query-block-jobs (Command)

Return information about long-running block device operations.

Returns

a list of `BlockJobInfo` for each active block job

Since

1.1

`block_passwd` (Command)

This command sets the password of a block device that has not been open with a password and requires one.

This command is now obsolete and will always return an error since 2.10

Arguments

device: string (optional) Not documented

node-name: string (optional) Not documented

password: string Not documented

`block_resize` (Command)

Resize a block image while a guest is running.

Either `device` or `node-name` must be set but not both.

Arguments

device: string (optional) the name of the device to get the image resized

node-name: string (optional) graph node name to get the image resized (Since 2.0)

size: int new image size in bytes

Returns

- nothing on success
- If `device` is not a valid block device, `DeviceNotFound`

Since

0.14

Example

```
-> { "execute": "block_resize",  
      "arguments": { "device": "scratch", "size": 1073741824 } }  
<- { "return": {} }
```

NewImageMode (Enum)

An enumeration that tells QEMU how to set the backing file path in a new image file.

Values

existing QEMU should look for an existing image file.

absolute-paths QEMU should create a new image with absolute paths for the backing file. If there is no backing file available, the new image will not be backed either.

Since

1.1

BlockdevSnapshotSync (Object)

Either `device` or `node-name` must be set but not both.

Members

device: string (optional) the name of the device to take a snapshot of.

node-name: string (optional) graph node name to generate the snapshot from (Since 2.0)

snapshot-file: string the target of the new overlay image. If the file exists, or if it is a device, the overlay will be created in the existing file/device. Otherwise, a new file will be created.

snapshot-node-name: string (optional) the graph node name of the new image (Since 2.0)

format: string (optional) the format of the overlay image, default is ‘qcow2’.

mode: NewImageMode (optional) whether and how QEMU should create a new image, default is ‘absolute-paths’.

BlockdevSnapshot (Object)

Members

node: string device or node name that will have a snapshot taken.

overlay: string reference to the existing block device that will become the overlay of `node`, as part of taking the snapshot. It must not have a current backing file (this can be achieved by passing “backing”: null to `blockdev-add`).

Since

2.5

BackupPerf (Object)

Optional parameters for backup. These parameters don’t affect functionality, but may significantly affect performance.

Members

use-copy-range: boolean (optional) Use copy offloading. Default false.

max-workers: int (optional) Maximum number of parallel requests for the sustained background copying process. Doesn't influence copy-before-write operations. Default 64.

max-chunk: int (optional) Maximum request length for the sustained background copying process. Doesn't influence copy-before-write operations. 0 means unlimited. If max-chunk is non-zero then it should not be less than job cluster size which is calculated as maximum of target image cluster size and 64k. Default 0.

Since

6.0

BackupCommon (Object)

Members

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device name or node-name of a root node which should be copied.

sync: MirrorSyncMode what parts of the disk image should be copied to the destination (all the disk, only the sectors allocated in the topmost image, from a dirty bitmap, or only new I/O).

speed: int (optional) the maximum speed, in bytes per second. The default is 0, for unlimited.

bitmap: string (optional) The name of a dirty bitmap to use. Must be present if sync is "bitmap" or "incremental". Can be present if sync is "full" or "top". Must not be present otherwise. (Since 2.4 (drive-backup), 3.1 (blockdev-backup))

bitmap-mode: BitmapSyncMode (optional) Specifies the type of data the bitmap should contain after the operation concludes. Must be present if a bitmap was provided, Must NOT be present otherwise. (Since 4.2)

compress: boolean (optional) true to compress data, if the target format supports it. (default: false) (since 2.8)

on-source-error: BlockdevOnError (optional) the action to take on an error on the source, default 'report'. 'stop' and 'enospc' can only be used if the block device supports io-status (see BlockInfo).

on-target-error: BlockdevOnError (optional) the action to take on an error on the target, default 'report' (no limitations, since this applies to a different block device than `device`).

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 2.12)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 2.12)

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the backup job inserts into the graph above node specified by `drive`. If this option is not given, a node name is autogenerated. (Since: 4.2)

x-perf: BackupPerf (optional) Performance options. (Since 6.0)

Note

`on-source-error` and `on-target-error` only affect background I/O. If an error occurs during a guest write request, the device's `error/werror` actions will be used.

Since

4.2

DriveBackup (Object)

Members

target: string the target of the new image. If the file exists, or if it is a device, the existing file/device will be used as the new destination. If it does not exist, a new file will be created.

format: string (optional) the format of the new destination, default is to probe if `mode` is 'existing', else the format of the source

mode: NewImageMode (optional) whether and how QEMU should create a new image, default is 'absolute-paths'.

The members of `BackupCommon`

Since

1.6

BlockdevBackup (Object)

Members

target: string the device name or node-name of the backup target node.

The members of `BackupCommon`

Since

2.3

blockdev-snapshot-sync (Command)

Takes a synchronous snapshot of a block device.

For the arguments, see the documentation of `BlockdevSnapshotSync`.

Returns

- nothing on success
- If `device` is not a valid block device, `DeviceNotFound`

Since

0.14

Example

```
-> { "execute": "blockdev-snapshot-sync",
      "arguments": { "device": "ide-hd0",
                     "snapshot-file":
                       "/some/place/my-image",
                     "format": "qcow2" } }
<- { "return": {} }
```

blockdev-snapshot (Command)

Takes a snapshot of a block device.

Take a snapshot, by installing ‘node’ as the backing image of ‘overlay’. Additionally, if ‘node’ is associated with a block device, the block device changes to using ‘overlay’ as its new active image.

For the arguments, see the documentation of BlockdevSnapshot.

Features

allow-write-only-overlay If present, the check whether this operation is safe was relaxed so that it can be used to change backing file of a destination of a blockdev-mirror. (since 5.0)

Since

2.5

Example

```
-> { "execute": "blockdev-add",
      "arguments": { "driver": "qcow2",
                     "node-name": "node1534",
                     "file": { "driver": "file",
                               "filename": "hd1.qcow2" },
                     "backing": null } }

<- { "return": {} }

-> { "execute": "blockdev-snapshot",
      "arguments": { "node": "ide-hd0",
                     "overlay": "node1534" } }

<- { "return": {} }
```

change-backing-file (Command)

Change the backing file in the image file metadata. This does not cause QEMU to reopen the image file to reparse the backing filename (it may, however, perform a reopen to change permissions from r/o -> r/w -> r/o, if needed). The new backing file string is written into the image file metadata, and the QEMU internal strings are updated.

Arguments

image-node-name: string The name of the block driver state node of the image to modify. The “device” argument is used to verify “image-node-name” is in the chain described by “device”.

device: string The device name or node-name of the root node that owns image-node-name.

backing-file: string The string to write as the backing file. This string is not validated, so care should be taken when specifying the string or the image chain may not be able to be reopened again.

Returns

- Nothing on success
- If “device” does not exist or cannot be determined, DeviceNotFound

Since

2.1

block-commit (Command)

Live commit of data from overlay image nodes into backing nodes - i.e., writes data between ‘top’ and ‘base’ into ‘base’.

If top == base, that is an error. If top has no overlays on top of it, or if it is in use by a writer, the job will not be completed by itself. The user needs to complete the job with the block-job-complete command after getting the ready event. (Since 2.0)

If the base image is smaller than top, then the base image will be resized to be the same size as top. If top is smaller than the base image, the base will not be truncated. If you want the base image size to match the size of the smaller top, you can safely truncate it yourself once the commit operation successfully completes.

Arguments

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device name or node-name of a root node

base-node: string (optional) The node name of the backing image to write data into. If not specified, this is the deepest backing image. (since: 3.1)

base: string (optional) Same as `base-node`, except that it is a file name rather than a node name. This must be the exact filename string that was used to open the node; other strings, even if addressing the same file, are not accepted

top-node: string (optional) The node name of the backing image within the image chain which contains the topmost data to be committed down. If not specified, this is the active layer. (since: 3.1)

top: string (optional) Same as `top-node`, except that it is a file name rather than a node name. This must be the exact filename string that was used to open the node; other strings, even if addressing the same file, are not accepted

backing-file: string (optional) The backing file string to write into the overlay image of ‘top’. If ‘top’ does not have an overlay image, or if ‘top’ is in use by a writer, specifying a backing file string is an error.

This filename is not validated. If a pathname string is such that it cannot be resolved by QEMU, that means that subsequent QMP or HMP commands must use node-names for the image in question, as filename lookup methods will fail.

If not specified, QEMU will automatically determine the backing file string to use, or error out if there is no obvious choice. Care should be taken when specifying the string, to specify a valid filename or protocol. (Since 2.1)

speed: int (optional) the maximum speed, in bytes per second

on-error: BlockdevOnError (optional) the action to take on an error. ‘ignore’ means that the request should be retried. (default: report; Since: 5.0)

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the commit job inserts into the graph above `top`. If this option is not given, a node name is autogenerated. (Since: 2.9)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Features

deprecated Members `base` and `top` are deprecated. Use `base-node` and `top-node` instead.

Returns

- Nothing on success
- If `device` does not exist, `DeviceNotFound`
- Any other error returns a `GenericError`.

Since

1.3

Example

```
-> { "execute": "block-commit",  
      "arguments": { "device": "virtio0",  
                     "top": "/tmp/snap1.qcow2" } }  
<- { "return": {} }
```

drive-backup (Command)

Start a point-in-time copy of a block device to a new destination. The status of ongoing drive-backup operations can be checked with `query-block-jobs` where the `BlockJobInfo.type` field has the value 'backup'. The operation can be stopped before it has completed using the `block-job-cancel` command.

Arguments

The members of DriveBackup

Returns

- nothing on success
- If `device` is not a valid block device, `GenericError`

Since

1.6

Example

```
-> { "execute": "drive-backup",  
      "arguments": { "device": "drive0",  
                     "sync": "full",  
                     "target": "backup.img" } }  
<- { "return": {} }
```

blockdev-backup (Command)

Start a point-in-time copy of a block device to a new destination. The status of ongoing blockdev-backup operations can be checked with `query-block-jobs` where the `BlockJobInfo.type` field has the value 'backup'. The operation can be stopped before it has completed using the `block-job-cancel` command.

Arguments

The members of BlockdevBackup

Returns

- nothing on success
- If device is not a valid block device, DeviceNotFound

Since

2.3

Example

```
-> { "execute": "blockdev-backup",
      "arguments": { "device": "src-id",
                     "sync": "full",
                     "target": "tgt-id" } }
<- { "return": {} }
```

query-named-block-nodes (Command)

Get the named block driver list

Arguments

flat: boolean (optional) Omit the nested data about backing image (“backing-image” key) if true. Default is false (Since 5.0)

Returns

the list of BlockDeviceInfo

Since

2.0

Example

```
-> { "execute": "query-named-block-nodes" }
<- { "return": [ { "ro":false,
                  "drv":"qcow2",
                  "encrypted":false,
                  "file":"disks/test.qcow2",
                  "node-name": "my-node",
                  "backing_file_depth":1,
                  "bps":1000000,
                  "bps_rd":0,
                  "bps_wr":0,
```

(continues on next page)

(continued from previous page)

```

    "iops":1000000,
    "iops_rd":0,
    "iops_wr":0,
    "bps_max": 8000000,
    "bps_rd_max": 0,
    "bps_wr_max": 0,
    "iops_max": 0,
    "iops_rd_max": 0,
    "iops_wr_max": 0,
    "iops_size": 0,
    "write_threshold": 0,
    "image":{
        "filename":"disks/test.qcow2",
        "format":"qcow2",
        "virtual-size":2048000,
        "backing_file":"base.qcow2",
        "full-backing-filename":"disks/base.qcow2",
        "backing-filename-format":"qcow2",
        "snapshots":[
            {
                "id": "1",
                "name": "snapshot1",
                "vm-state-size": 0,
                "date-sec": 10000200,
                "date-nsec": 12,
                "vm-clock-sec": 206,
                "vm-clock-nsec": 30
            }
        ],
        "backing-image":{
            "filename":"disks/base.qcow2",
            "format":"qcow2",
            "virtual-size":2048000
        }
    }
} } ] }

```

XDbgBlockGraphNodeType (Enum)

Values

block-backend corresponds to BlockBackend

block-job corresponds to BlockJob

block-driver corresponds to BlockDriverState

Since

4.0

XDbgBlockGraphNode (Object)

Members

id: int Block graph node identifier. This `id` is generated only for `x-debug-query-block-graph` and does not relate to any other identifiers in Qemu.

type: XDbgBlockGraphNodeType Type of graph node. Can be one of `block-backend`, `block-job` or `block-driver-state`.

name: string Human readable name of the node. Corresponds to `node-name` for `block-driver-state` nodes; is not guaranteed to be unique in the whole graph (with `block-jobs` and `block-backends`).

Since

4.0

BlockPermission (Enum)

Enum of base block permissions.

Values

consistent-read A user that has the “permission” of consistent reads is guaranteed that their view of the contents of the block device is complete and self-consistent, representing the contents of a disk at a specific point. For most block devices (including their backing files) this is true, but the property cannot be maintained in a few situations like for intermediate nodes of a commit block job.

write This permission is required to change the visible disk contents.

write-unchanged This permission (which is weaker than `BLK_PERM_WRITE`) is both enough and required for writes to the block node when the caller promises that the visible disk content doesn’t change. As the `BLK_PERM_WRITE` permission is strictly stronger, either is sufficient to perform an unchanging write.

resize This permission is required to change the size of a block node.

graph-mod This permission is required to change the node that this `BdrvChild` points to.

Since

4.0

XDbgBlockGraphEdge (Object)

Block Graph edge description for `x-debug-query-block-graph`.

Members

parent: int parent id

child: int child id

name: string name of the relation (examples are ‘file’ and ‘backing’)

perm: array of BlockPermission granted permissions for the parent operating on the child

shared-perm: array of **BlockPermission** permissions that can still be granted to other users of the child while it is still attached to this parent

Since

4.0

XDbgBlockGraph (Object)

Block Graph - list of nodes and list of edges.

Members

nodes: array of **XDbgBlockGraphNode** Not documented

edges: array of **XDbgBlockGraphEdge** Not documented

Since

4.0

x-debug-query-block-graph (Command)

Get the block graph.

Since

4.0

drive-mirror (Command)

Start mirroring a block device's writes to a new destination. `target` specifies the target of the new image. If the file exists, or if it is a device, it will be used as the new destination for writes. If it does not exist, a new file will be created. `format` specifies the format of the mirror image, default is to probe if `mode='existing'`, else the format of the source.

Arguments

The members of **DriveMirror**

Returns

- nothing on success
- If `device` is not a valid block device, **GenericError**

Since

1.3

Example

```
-> { "execute": "drive-mirror",
    "arguments": { "device": "ide-hd0",
                  "target": "/some/place/my-image",
                  "sync": "full",
                  "format": "qcow2" } }
<- { "return": {} }
```

DriveMirror (Object)

A set of parameters describing drive mirror setup.

Members

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device name or node-name of a root node whose writes should be mirrored.

target: string the target of the new image. If the file exists, or if it is a device, the existing file/device will be used as the new destination. If it does not exist, a new file will be created.

format: string (optional) the format of the new destination, default is to probe if `mode` is ‘existing’, else the format of the source

node-name: string (optional) the new block driver state node name in the graph (Since 2.1)

replaces: string (optional) with `sync=full` graph node name to be replaced by the new image when a whole image copy is done. This can be used to repair broken Quorum files. By default, `device` is replaced, although implicitly created filters on it are kept. (Since 2.1)

mode: NewImageMode (optional) whether and how QEMU should create a new image, default is ‘absolute-paths’.

speed: int (optional) the maximum speed, in bytes per second

sync: MirrorSyncMode what parts of the disk image should be copied to the destination (all the disk, only the sectors allocated in the topmost image, or only new I/O).

granularity: int (optional) granularity of the dirty bitmap, default is 64K if the image format doesn’t have clusters, 4K if the clusters are smaller than that, else the cluster size. Must be a power of 2 between 512 and 64M (since 1.4).

buf-size: int (optional) maximum amount of data in flight from source to target (since 1.4).

on-source-error: BlockdevOnError (optional) the action to take on an error on the source, default ‘report’. ‘stop’ and ‘enospc’ can only be used if the block device supports io-status (see BlockInfo).

on-target-error: BlockdevOnError (optional) the action to take on an error on the target, default ‘report’ (no limitations, since this applies to a different block device than `device`).

unmap: boolean (optional) Whether to try to unmap target sectors where source has only zero. If true, and target unallocated sectors will read as zero, target image sectors will be unmapped; otherwise, zeroes will be written. Both will result in identical contents. Default is true. (Since 2.4)

copy-mode: MirrorCopyMode (optional) when to copy data to the destination; defaults to 'background' (Since: 3.0)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Since

1.3

BlockDirtyBitmap (Object)

Members

node: string name of device/node which the bitmap is tracking

name: string name of the dirty bitmap

Since

2.4

BlockDirtyBitmapAdd (Object)

Members

node: string name of device/node which the bitmap is tracking

name: string name of the dirty bitmap (must be less than 1024 bytes)

granularity: int (optional) the bitmap granularity, default is 64k for block-dirty-bitmap-add

persistent: boolean (optional) the bitmap is persistent, i.e. it will be saved to the corresponding block device image file on its close. For now only Qcow2 disks support persistent bitmaps. Default is false for block-dirty-bitmap-add. (Since: 2.10)

disabled: boolean (optional) the bitmap is created in the disabled state, which means that it will not track drive changes. The bitmap may be enabled with `block-dirty-bitmap-enable`. Default is false. (Since: 4.0)

Since

2.4

BlockDirtyBitmapMergeSource (Alternate)

Members

local: **string** name of the bitmap, attached to the same node as target bitmap.

external: **BlockDirtyBitmap** bitmap with specified node

Since

4.1

BlockDirtyBitmapMerge (Object)

Members

node: **string** name of device/node which the `target` bitmap is tracking

target: **string** name of the destination dirty bitmap

bitmaps: **array of BlockDirtyBitmapMergeSource** name(s) of the source dirty bitmap(s) at `node` and/or fully specified `BlockDirtyBitmap` elements. The latter are supported since 4.1.

Since

4.0

block-dirty-bitmap-add (Command)

Create a dirty bitmap with a name on the node, and start tracking the writes.

Returns

- nothing on success
- If `node` is not a valid block device or node, `DeviceNotFound`
- If `name` is already taken, `GenericError` with an explanation

Since

2.4

Example

```
-> { "execute": "block-dirty-bitmap-add",
      "arguments": { "node": "drive0", "name": "bitmap0" } }
<- { "return": {} }
```

block-dirty-bitmap-remove (Command)

Stop write tracking and remove the dirty bitmap that was created with `block-dirty-bitmap-add`. If the bitmap is persistent, remove it from its storage too.

Returns

- nothing on success
- If `node` is not a valid block device or node, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation
- if `name` is frozen by an operation, `GenericError`

Since

2.4

Example

```
-> { "execute": "block-dirty-bitmap-remove",  
      "arguments": { "node": "drive0", "name": "bitmap0" } }  
<- { "return": {} }
```

block-dirty-bitmap-clear (Command)

Clear (reset) a dirty bitmap on the device, so that an incremental backup from this point in time forward will only backup clusters modified after this clear operation.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation

Since

2.4

Example

```
-> { "execute": "block-dirty-bitmap-clear",  
      "arguments": { "node": "drive0", "name": "bitmap0" } }  
<- { "return": {} }
```


block-dirty-bitmap-enable (Command)

Enables a dirty bitmap so that it will begin tracking disk changes.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation

Since

4.0

Example

```
-> { "execute": "block-dirty-bitmap-enable",  
    "arguments": { "node": "drive0", "name": "bitmap0" } }  
<- { "return": {} }
```

block-dirty-bitmap-disable (Command)

Disables a dirty bitmap so that it will stop tracking disk changes.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If `name` is not found, `GenericError` with an explanation

Since

4.0

Example

```
-> { "execute": "block-dirty-bitmap-disable",  
    "arguments": { "node": "drive0", "name": "bitmap0" } }  
<- { "return": {} }
```

block-dirty-bitmap-merge (Command)

Merge dirty bitmaps listed in `bitmaps` to the `target` dirty bitmap. Dirty bitmaps in `bitmaps` will be unchanged, except if it also appears as the `target` bitmap. Any bits already set in `target` will still be set after the merge, i.e., this operation does not clear the target. On error, `target` is unchanged.

The resulting bitmap will count as dirty any clusters that were dirty in any of the source bitmaps. This can be used to achieve backup checkpoints, or in simpler usages, to copy bitmaps.

Returns

- nothing on success
- If `node` is not a valid block device, `DeviceNotFound`
- If any bitmap in `bitmaps` or `target` is not found, `GenericError`
- If any of the bitmaps have different sizes or granularities, `GenericError`

Since

4.0

Example

```
-> { "execute": "block-dirty-bitmap-merge",  
    "arguments": { "node": "drive0", "target": "bitmap0",  
                  "bitmaps": ["bitmap1"] } }  
<- { "return": {} }
```

BlockDirtyBitmapSha256 (Object)

SHA256 hash of dirty bitmap data

Members

sha256: `string` ASCII representation of SHA256 bitmap hash

Since

2.10

x-debug-block-dirty-bitmap-sha256 (Command)

Get bitmap SHA256.

Returns

- BlockDirtyBitmapSha256 on success
- If `node` is not a valid block device, DeviceNotFound
- If `name` is not found or if hashing has failed, GenericError with an explanation

Since

2.10

blockdev-mirror (Command)

Start mirroring a block device's writes to a new destination.

Arguments

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string The device name or node-name of a root node whose writes should be mirrored.

target: string the id or node-name of the block device to mirror to. This mustn't be attached to guest.

replaces: string (optional) with sync=full graph node name to be replaced by the new image when a whole image copy is done. This can be used to repair broken Quorum files. By default, `device` is replaced, although implicitly created filters on it are kept.

speed: int (optional) the maximum speed, in bytes per second

sync: MirrorSyncMode what parts of the disk image should be copied to the destination (all the disk, only the sectors allocated in the topmost image, or only new I/O).

granularity: int (optional) granularity of the dirty bitmap, default is 64K if the image format doesn't have clusters, 4K if the clusters are smaller than that, else the cluster size. Must be a power of 2 between 512 and 64M

buf-size: int (optional) maximum amount of data in flight from source to target

on-source-error: BlockdevOnError (optional) the action to take on an error on the source, default 'report'. 'stop' and 'enospc' can only be used if the block device supports io-status (see BlockInfo).

on-target-error: BlockdevOnError (optional) the action to take on an error on the target, default 'report' (no limitations, since this applies to a different block device than `device`).

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the mirror job inserts into the graph above `device`. If this option is not given, a node name is autogenerated. (Since: 2.9)

copy-mode: MirrorCopyMode (optional) when to copy data to the destination; defaults to 'background' (Since: 3.0)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Returns

nothing on success.

Since

2.6

Example

```
-> { "execute": "blockdev-mirror",
      "arguments": { "device": "ide-hd0",
                     "target": "target0",
                     "sync": "full" } }
<- { "return": {} }
```

BlockIOThrottle (Object)

A set of parameters describing block throttling.

Members

device: string (optional) Block device name

id: string (optional) The name or QOM path of the guest device (since: 2.8)

bps: int total throughput limit in bytes per second

bps_rd: int read throughput limit in bytes per second

bps_wr: int write throughput limit in bytes per second

iops: int total I/O operations per second

iops_rd: int read I/O operations per second

iops_wr: int write I/O operations per second

bps_max: int (optional) total throughput limit during bursts, in bytes (Since 1.7)

bps_rd_max: int (optional) read throughput limit during bursts, in bytes (Since 1.7)

bps_wr_max: int (optional) write throughput limit during bursts, in bytes (Since 1.7)

iops_max: int (optional) total I/O operations per second during bursts, in bytes (Since 1.7)

iops_rd_max: int (optional) read I/O operations per second during bursts, in bytes (Since 1.7)

iops_wr_max: int (optional) write I/O operations per second during bursts, in bytes (Since 1.7)

bps_max_length: int (optional) maximum length of the `bps_max` burst period, in seconds. It must only be set if `bps_max` is set as well. Defaults to 1. (Since 2.6)

bps_rd_max_length: int (optional) maximum length of the `bps_rd_max` burst period, in seconds. It must only be set if `bps_rd_max` is set as well. Defaults to 1. (Since 2.6)

bps_wr_max_length: int (optional) maximum length of the `bps_wr_max` burst period, in seconds. It must only be set if `bps_wr_max` is set as well. Defaults to 1. (Since 2.6)

iops_max_length: int (optional) maximum length of the `iops` burst period, in seconds. It must only be set if `iops_max` is set as well. Defaults to 1. (Since 2.6)

iops_rd_max_length: int (optional) maximum length of the `iops_rd_max` burst period, in seconds. It must only be set if `iops_rd_max` is set as well. Defaults to 1. (Since 2.6)

iops_wr_max_length: int (optional) maximum length of the `iops_wr_max` burst period, in seconds. It must only be set if `iops_wr_max` is set as well. Defaults to 1. (Since 2.6)

iops_size: int (optional) an I/O size in bytes (Since 1.7)

group: string (optional) throttle group name (Since 2.4)

Features

deprecated Member `device` is deprecated. Use `id` instead.

Since

1.1

ThrottleLimits (Object)

Limit parameters for throttling. Since some limit combinations are illegal, limits should always be set in one transaction. All fields are optional. When setting limits, if a field is missing the current value is not changed.

Members

iops-total: int (optional) limit total I/O operations per second

iops-total-max: int (optional) I/O operations burst

iops-total-max-length: int (optional) length of the `iops-total-max` burst period, in seconds It must only be set if `iops-total-max` is set as well.

iops-read: int (optional) limit read operations per second

iops-read-max: int (optional) I/O operations read burst

iops-read-max-length: int (optional) length of the `iops-read-max` burst period, in seconds It must only be set if `iops-read-max` is set as well.

iops-write: int (optional) limit write operations per second

iops-write-max: int (optional) I/O operations write burst

iops-write-max-length: int (optional) length of the `iops-write-max` burst period, in seconds It must only be set if `iops-write-max` is set as well.

bps-total: int (optional) limit total bytes per second

bps-total-max: int (optional) total bytes burst

bps-total-max-length: int (optional) length of the bps-total-max burst period, in seconds. It must only be set if bps-total-max is set as well.

bps-read: int (optional) limit read bytes per second

bps-read-max: int (optional) total bytes read burst

bps-read-max-length: int (optional) length of the bps-read-max burst period, in seconds It must only be set if bps-read-max is set as well.

bps-write: int (optional) limit write bytes per second

bps-write-max: int (optional) total bytes write burst

bps-write-max-length: int (optional) length of the bps-write-max burst period, in seconds It must only be set if bps-write-max is set as well.

iops-size: int (optional) when limiting by iops max size of an I/O in bytes

Since

2.11

block-stream (Command)

Copy data from a backing file into a block device.

The block streaming operation is performed in the background until the entire backing file has been copied. This command returns immediately once streaming has started. The status of ongoing block streaming operations can be checked with query-block-jobs. The operation can be stopped before it has completed using the block-job-cancel command.

The node that receives the data is called the top image, can be located in any part of the chain (but always above the base image; see below) and can be specified using its device or node name. Earlier qemu versions only allowed ‘device’ to name the top level node; presence of the ‘base-node’ parameter during introspection can be used as a witness of the enhanced semantics of ‘device’.

If a base file is specified then sectors are not copied from that base file and its backing chain. This can be used to stream a subset of the backing file chain instead of flattening the entire image. When streaming completes the image file will have the base file as its backing file, unless that node was changed while the job was running. In that case, base’s parent’s backing (or filtered, whichever exists) child (i.e., base at the beginning of the job) will be the new backing file.

On successful completion the image file is updated to drop the backing file and the BLOCK_JOB_COMPLETED event is emitted.

In case device is a filter node, block-stream modifies the first non-filter overlay node below it to point to the new backing node instead of modifying device itself.

Arguments

job-id: string (optional) identifier for the newly-created block job. If omitted, the device name will be used. (Since 2.7)

device: string the device or node name of the top image

base: string (optional) the common backing file name. It cannot be set if base-node or bottom is also set.

base-node: string (optional) the node name of the backing file. It cannot be set if `base` or `bottom` is also set. (Since 2.8)

bottom: string (optional) the last node in the chain that should be streamed into top. It cannot be set if `base` or `base-node` is also set. It cannot be filter node. (Since 6.0)

backing-file: string (optional) The backing file string to write into the top image. This filename is not validated.

If a pathname string is such that it cannot be resolved by QEMU, that means that subsequent QMP or HMP commands must use node-names for the image in question, as filename lookup methods will fail.

If not specified, QEMU will automatically determine the backing file string to use, or error out if there is no obvious choice. Care should be taken when specifying the string, to specify a valid filename or protocol. (Since 2.1)

speed: int (optional) the maximum speed, in bytes per second

on-error: BlockdevOnError (optional) the action to take on an error (default report). ‘stop’ and ‘enospc’ can only be used if the block device supports io-status (see BlockInfo). Since 1.3.

filter-node-name: string (optional) the node name that should be assigned to the filter driver that the stream job inserts into the graph above `device`. If this option is not given, a node name is autogenerated. (Since: 6.0)

auto-finalize: boolean (optional) When false, this job will wait in a PENDING state after it has finished its work, waiting for `block-job-finalize` before making any block graph changes. When true, this job will automatically perform its abort or commit actions. Defaults to true. (Since 3.1)

auto-dismiss: boolean (optional) When false, this job will wait in a CONCLUDED state after it has completely ceased all work, and awaits `block-job-dismiss`. When true, this job will automatically disappear from the query list without user intervention. Defaults to true. (Since 3.1)

Returns

- Nothing on success.
- If device does not exist, DeviceNotFound.

Since

1.1

Example

```
-> { "execute": "block-stream",
      "arguments": { "device": "virtio0",
                     "base": "/tmp/master.qcow2" } }
<- { "return": {} }
```

block-job-set-speed (Command)

Set maximum speed for a background block operation.

This command can only be issued when there is an active block job.

Throttling can be disabled by setting the speed to 0.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

speed: int the maximum speed, in bytes per second, or 0 for unlimited. Defaults to 0.

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.1

block-job-cancel (Command)

Stop an active background block operation.

This command returns immediately after marking the active background block operation for cancellation. It is an error to call this command if no operation is in progress.

The operation will cancel as soon as possible and then emit the BLOCK_JOB_CANCELLED event. Before that happens the job is still visible when enumerated using query-block-jobs.

Note that if you issue ‘block-job-cancel’ after ‘drive-mirror’ has indicated (via the event BLOCK_JOB_READY) that the source and destination are synchronized, then the event triggered by this command changes to BLOCK_JOB_COMPLETED, to indicate that the mirroring has ended and the destination now has a point-in-time copy tied to the time of the cancellation.

For streaming, the image file retains its backing file unless the streaming operation happens to complete just as it is being cancelled. A new streaming operation can be started at a later time to finish copying all data from the backing file.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

force: boolean (optional) If true, and the job has already emitted the event BLOCK_JOB_READY, abandon the job immediately (even if it is paused) instead of waiting for the destination to complete its final synchronization (since 1.3)

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.1

block-job-pause (Command)

Pause an active background block operation.

This command returns immediately after marking the active background block operation for pausing. It is an error to call this command if no operation is in progress or if the job is already paused.

The operation will pause as soon as possible. No event is emitted when the operation is actually paused. Cancelling a paused job automatically resumes it.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.3

block-job-resume (Command)

Resume an active background block operation.

This command returns immediately after resuming a paused background block operation. It is an error to call this command if no operation is in progress or if the job is not paused.

This command also clears the error status of the job.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.3

block-job-complete (Command)

Manually trigger completion of an active background block operation. This is supported for drive mirroring, where it also switches the device to write to the target path only. The ability to complete is signaled with a BLOCK_JOB_READY event.

This command completes an active background block operation synchronously. The ordering of this command's return with the BLOCK_JOB_COMPLETED event is not defined. Note that if an I/O error occurs during the processing of this command: 1) the command itself will fail; 2) the error will be processed according to the error/werror arguments that were specified when starting the operation.

A cancelled or paused job cannot be completed.

Arguments

device: string The job identifier. This used to be a device name (hence the name of the parameter), but since QEMU 2.7 it can have other values.

Returns

- Nothing on success
- If no background operation is active on this device, DeviceNotActive

Since

1.3

block-job-dismiss (Command)

For jobs that have already concluded, remove them from the block-job-query list. This command only needs to be run for jobs which were started with QEMU 2.12+ job lifetime management semantics.

This command will refuse to operate on any job that has not yet reached its terminal state, JOB_STATUS_CONCLUDED. For jobs that make use of the BLOCK_JOB_READY event, block-job-cancel or block-job-complete will still need to be used as appropriate.

Arguments

id: string The job identifier.

Returns

Nothing on success

Since

2.12

block-job-finalize (Command)

Once a job that has `manual=true` reaches the pending state, it can be instructed to finalize any graph changes and do any necessary cleanup via this command. For jobs in a transaction, instructing one job to finalize will force ALL jobs in the transaction to finalize, so it is only necessary to instruct a single member job to finalize.

Arguments

id: string The job identifier.

Returns

Nothing on success

Since

2.12

BlockdevDiscardOptions (Enum)

Determines how to handle discard requests.

Values

ignore Ignore the request

unmap Forward as an unmap request

Since

2.9

BlockdevDetectZeroesOptions (Enum)

Describes the operation mode for the automatic conversion of plain zero writes by the OS to driver specific optimized zero write commands.

Values

off Disabled (default)

on Enabled

unmap Enabled and even try to unmap blocks if possible. This requires also that `BlockdevDiscardOptions` is set to unmap for this device.

Since

2.1

BlockdevAioOptions (Enum)

Selects the AIO backend to handle I/O requests

Values

threads Use qemu's thread pool

native Use native AIO backend (only Linux and Windows)

io_uring (If: `defined(CONFIG_LINUX_IO_URING)`) Use linux io_uring (since 5.0)

Since

2.9

BlockdevCacheOptions (Object)

Includes cache-related options for block devices

Members

direct: boolean (optional) enables use of O_DIRECT (bypass the host page cache; default: false)

no-flush: boolean (optional) ignore any flush requests for the device (default: false)

Since

2.9

BlockdevDriver (Enum)

Drivers that are supported in block device operations.

Values

throttle Since 2.11
nvme Since 2.12
copy-on-read Since 3.0
blklogwrites Since 3.0
blkreplay Since 4.2
compress Since 5.0
blkdebug Not documented
blkverify Not documented
bochs Not documented
cloop Not documented
dmg Not documented
file Not documented
ftp Not documented
ftps Not documented
gluster Not documented
host_cdrom Not documented
host_device Not documented
http Not documented
https Not documented
iscsi Not documented
luks Not documented
nbd Not documented
nfs Not documented
null-aio Not documented
null-co Not documented
parallels Not documented
preallocate Not documented
qcow Not documented
qcow2 Not documented
qed Not documented
quorum Not documented
raw Not documented
rbd Not documented
replication (If: `defined(CONFIG_REPLICATION)`) Not documented
sheepdog Not documented

ssh Not documented
vdi Not documented
vhd~~x~~ Not documented
vmdk Not documented
vpc Not documented
vvfat Not documented

Since

2.9

BlockdevOptionsFile (Object)

Driver specific block device options for the file backend.

Members

filename: **string** path to the image file
pr-manager: **string (optional)** the id for the object that will handle persistent reservations for this device (default: none, forward the commands via SG_IO; since 2.11)
aio: **BlockdevAioOptions (optional)** AIO backend (default: threads) (since: 2.8)
locking: **OnOffAuto (optional)** whether to enable file locking. If set to ‘auto’, only enable when Open File Descriptor (OFD) locking API is available (default: auto, since 2.10)
drop-cache: **boolean (optional) (If: defined (CONFIG_LINUX))** invalidate page cache during live migration. This prevents stale data on the migration destination with cache.direct=off. Currently only supported on Linux hosts. (default: on, since: 4.0)
x-check-cache-dropped: **boolean (optional)** whether to check that page cache was dropped on live migration. May cause noticeable delays if the image file is large, do not use in production. (default: off) (since: 3.0)

Features

dynamic-auto-read-only If present, enabled auto-read-only means that the driver will open the image read-only at first, dynamically reopen the image file read-write when the first writer is attached to the node and reopen read-only when the last writer is detached. This allows giving QEMU write permissions only on demand when an operation actually needs write access.

Since

2.9

BlockdevOptionsNull (Object)

Driver specific block device options for the null backend.

Members

size: int (optional) size of the device in bytes.

latency-ns: int (optional) emulated latency (in nanoseconds) in processing requests. Default to zero which completes requests immediately. (Since 2.4)

read-zeroes: boolean (optional) if true, reads from the device produce zeroes; if false, the buffer is left unchanged. (default: false; since: 4.1)

Since

2.9

BlockdevOptionsNVMe (Object)

Driver specific block device options for the NVMe backend.

Members

device: string PCI controller address of the NVMe device in format hhhh:bb:ss.f (host:bus:slot.function)

namespace: int namespace number of the device, starting from 1.

Note that the PCI `device` must have been unbound from any host kernel driver before instructing QEMU to add the `blockdev`.

Since

2.12

BlockdevOptionsVVFAT (Object)

Driver specific block device options for the vvfat protocol.

Members

dir: string directory to be exported as FAT image

fat-type: int (optional) FAT type: 12, 16 or 32

floppy: boolean (optional) whether to export a floppy image (true) or partitioned hard disk (false; default)

label: string (optional) set the volume label, limited to 11 bytes. FAT16 and FAT32 traditionally have some restrictions on labels, which are ignored by most operating systems. Defaults to “QEMU VVFAT”. (since 2.4)

rw: boolean (optional) whether to allow write operations (default: false)

Since

2.9

BlockdevOptionsGenericFormat (Object)

Driver specific block device options for image format that have no option besides their data source.

Members

file: BlockdevRef reference to or definition of the data source block device

Since

2.9

BlockdevOptionsLUKS (Object)

Driver specific block device options for LUKS.

Members

key-secret: string (optional) the ID of a QCryptoSecret object providing the decryption key (since 2.6). Mandatory except when doing a metadata-only probe of the image.

The members of BlockdevOptionsGenericFormat

Since

2.9

BlockdevOptionsGenericCOWFormat (Object)

Driver specific block device options for image format that have no option besides their data source and an optional backing file.

Members

backing: BlockdevRefOrNull (optional) reference to or definition of the backing file block device, null disables the backing file entirely. Defaults to the backing file stored the image file.

The members of BlockdevOptionsGenericFormat

Since

2.9

Qcow2OverlapCheckMode (Enum)

General overlap check modes.

Values

none Do not perform any checks

constant Perform only checks which can be done in constant time and without reading anything from disk

cached Perform only checks which can be done without reading anything from disk

all Perform all available overlap checks

Since

2.9

Qcow2OverlapCheckFlags (Object)

Structure of flags for each metadata structure. Setting a field to ‘true’ makes qemu guard that structure against unintended overwriting. The default value is chosen according to the template given.

Members

template: Qcow2OverlapCheckMode (optional) Specifies a template mode which can be adjusted using the other flags, defaults to ‘cached’

bitmap-directory: boolean (optional) since 3.0

main-header: boolean (optional) Not documented

active-11: boolean (optional) Not documented

active-12: boolean (optional) Not documented

refcount-table: boolean (optional) Not documented

refcount-block: boolean (optional) Not documented

snapshot-table: boolean (optional) Not documented

inactive-11: boolean (optional) Not documented

inactive-12: boolean (optional) Not documented

Since

2.9

Qcow2OverlapChecks (Alternate)

Specifies which metadata structures should be guarded against unintended overwriting.

Members

flags: Qcow2OverlapCheckFlags set of flags for separate specification of each metadata structure type

mode: Qcow2OverlapCheckMode named mode which chooses a specific set of flags

Since

2.9

BlockdevQcowEncryptionFormat (Enum)

Values

aes AES-CBC with plain64 initialization vectors

Since

2.10

BlockdevQcowEncryption (Object)

Members

format: `BlockdevQcowEncryptionFormat` Not documented

The members of `QCryptoBlockOptionsQcow` when **format** is "aes"

Since

2.10

BlockdevOptionsQcow (Object)

Driver specific block device options for qcow.

Members

encrypt: `BlockdevQcowEncryption (optional)` Image decryption options. Mandatory for encrypted images, except when doing a metadata-only probe of the image.

The members of `BlockdevOptionsGenericCOWFormat`

Since

2.10

BlockdevQcow2EncryptionFormat (Enum)

Values

aes AES-CBC with plain64 initialization vectors

luks Not documented

Since

2.10

BlockdevQcow2Encryption (Object)

Members

format: `BlockdevQcow2EncryptionFormat` Not documented

The members of `QCryptoBlockOptionsQcow` when **format** is "aes"

The members of `QCryptoBlockOptionsLUKS` when **format** is "luks"

Since

2.10

BlockdevOptionsPreallocate (Object)

Filter driver intended to be inserted between format and protocol node and do preallocation in protocol node on write.

Members

prealloc-align: `int (optional)` on preallocation, align file length to this number, default 1048576 (1M)

prealloc-size: `int (optional)` how much to preallocate, default 134217728 (128M)

The members of `BlockdevOptionsGenericFormat`

Since

6.0

BlockdevOptionsQcow2 (Object)

Driver specific block device options for qcow2.

Members

lazy-refcounts: `boolean (optional)` whether to enable the lazy refcounts feature (default is taken from the image file)

pass-discard-request: `boolean (optional)` whether discard requests to the qcow2 device should be forwarded to the data source

pass-discard-snapshot: `boolean (optional)` whether discard requests for the data source should be issued when a snapshot operation (e.g. deleting a snapshot) frees clusters in the qcow2 file

pass-discard-other: `boolean (optional)` whether discard requests for the data source should be issued on other occasions where a cluster gets freed

overlap-check: Qcow2OverlapChecks (optional) which overlap checks to perform for writes to the image, defaults to 'cached' (since 2.2)

cache-size: int (optional) the maximum total size of the L2 table and refcount block caches in bytes (since 2.2)

l2-cache-size: int (optional) the maximum size of the L2 table cache in bytes (since 2.2)

l2-cache-entry-size: int (optional) the size of each entry in the L2 cache in bytes. It must be a power of two between 512 and the cluster size. The default value is the cluster size (since 2.12)

refcount-cache-size: int (optional) the maximum size of the refcount block cache in bytes (since 2.2)

cache-clean-interval: int (optional) clean unused entries in the L2 and refcount caches. The interval is in seconds. The default value is 600 on supporting platforms, and 0 on other platforms. 0 disables this feature. (since 2.5)

encrypt: BlockdevQcow2Encryption (optional) Image decryption options. Mandatory for encrypted images, except when doing a metadata-only probe of the image. (since 2.10)

data-file: BlockdevRef (optional) reference to or definition of the external data file. This may only be specified for images that require an external data file. If it is not specified for such an image, the data file name is loaded from the image file. (since 4.0)

The members of `BlockdevOptionsGenericCOWFormat`

Since

2.9

SshHostKeyCheckMode (Enum)

Values

none Don't check the host key at all

hash Compare the host key with a given hash

known_hosts Check the host key against the known_hosts file

Since

2.12

SshHostKeyCheckHashType (Enum)

Values

md5 The given hash is an md5 hash

sha1 The given hash is an sha1 hash

Since

2.12

SshHostKeyHash (Object)

Members

type: **SshHostKeyCheckHashType** The hash algorithm used for the hash

hash: **string** The expected hash value

Since

2.12

SshHostKeyCheck (Object)

Members

mode: **SshHostKeyCheckMode** Not documented

The members of SshHostKeyHash when mode is "hash"

Since

2.12

BlockdevOptionsSsh (Object)

Members

server: **InetSocketAddress** host address

path: **string** path to the image on the host

user: **string (optional)** user as which to connect, defaults to current local user name

host-key-check: **SshHostKeyCheck (optional)** Defines how and what to check the host key against (default: known_hosts)

Since

2.9

BlkdebugEvent (Enum)

Trigger events supported by blkdebug.

Values

l1_shrink_write_table write zeros to the l1 table to shrink image. (since 2.11)

l1_shrink_free_l2_clusters discard the l2 tables. (since 2.11)

cor_write a write due to copy-on-read (since 2.11)

cluster_alloc_space an allocation of file space for a cluster (since 4.1)

none triggers once at creation of the blkdebug node (since 4.1)

l1_update Not documented

l1_grow_alloc_table Not documented

l1_grow_write_table Not documented

l1_grow_activate_table Not documented

l2_load Not documented

l2_update Not documented

l2_update_compressed Not documented

l2_alloc_cow_read Not documented

l2_alloc_write Not documented

read_aio Not documented

read_backing_aio Not documented

read_compressed Not documented

write_aio Not documented

write_compressed Not documented

vmstate_load Not documented

vmstate_save Not documented

cow_read Not documented

cow_write Not documented

reftable_load Not documented

reftable_grow Not documented

reftable_update Not documented

refblock_load Not documented

refblock_update Not documented

refblock_update_part Not documented

refblock_alloc Not documented

refblock_alloc_hookup Not documented

refblock_alloc_write Not documented

refblock_alloc_write_blocks Not documented

refblock_alloc_write_table Not documented

refblock_alloc_switch_table Not documented

cluster_alloc Not documented
cluster_alloc_bytes Not documented
cluster_free Not documented
flush_to_os Not documented
flush_to_disk Not documented
pwritev_rmw_head Not documented
pwritev_rmw_after_head Not documented
pwritev_rmw_tail Not documented
pwritev_rmw_after_tail Not documented
pwritev Not documented
pwritev_zero Not documented
pwritev_done Not documented
empty_image_prepare Not documented

Since

2.9

BlkdebugIOType (Enum)

Kinds of I/O that blkdebug can inject errors in.

Values

read .bdrv_co_preadv()
write .bdrv_co_pwritev()
write-zeroes .bdrv_co_pwrite_zeroes()
discard .bdrv_co_pdiscard()
flush .bdrv_co_flush_to_disk()
block-status .bdrv_co_block_status()

Since

4.1

BlkdebugInjectErrorOptions (Object)

Describes a single error injection for blkdebug.

Members

event: `BlkdebugEvent` trigger event

state: `int` (optional) the state identifier blkdebug needs to be in to actually trigger the event; defaults to “any”

iotype: `BlkdebugIOType` (optional) the type of I/O operations on which this error should be injected; defaults to “all read, write, write-zeroes, discard, and flush operations” (since: 4.1)

errno: `int` (optional) error identifier (errno) to be returned; defaults to EIO

sector: `int` (optional) specifies the sector index which has to be affected in order to actually trigger the event; defaults to “any sector”

once: `boolean` (optional) disables further events after this one has been triggered; defaults to false

immediately: `boolean` (optional) fail immediately; defaults to false

Since

2.9

`BlkdebugSetStateOptions` (Object)

Describes a single state-change event for blkdebug.

Members

event: `BlkdebugEvent` trigger event

state: `int` (optional) the current state identifier blkdebug needs to be in; defaults to “any”

new_state: `int` the state identifier blkdebug is supposed to assume if this event is triggered

Since

2.9

`BlockdevOptionsBlkdebug` (Object)

Driver specific block device options for blkdebug.

Members

image: `BlockdevRef` underlying raw block device (or image file)

config: `string` (optional) filename of the configuration file

align: `int` (optional) required alignment for requests in bytes, must be positive power of 2, or 0 for default

max-transfer: `int` (optional) maximum size for I/O transfers in bytes, must be positive multiple of `align` and of the underlying file’s request alignment (but need not be a power of 2), or 0 for default (since 2.10)

opt-write-zero: int (optional) preferred alignment for write zero requests in bytes, must be positive multiple of `align` and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

max-write-zero: int (optional) maximum size for write zero requests in bytes, must be positive multiple of `align`, of `opt-write-zero`, and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

opt-discard: int (optional) preferred alignment for discard requests in bytes, must be positive multiple of `align` and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

max-discard: int (optional) maximum size for discard requests in bytes, must be positive multiple of `align`, of `opt-discard`, and of the underlying file's request alignment (but need not be a power of 2), or 0 for default (since 2.10)

inject-error: array of BlkdebugInjectErrorOptions (optional) array of error injection descriptions

set-state: array of BlkdebugSetStateOptions (optional) array of state-change descriptions

take-child-perms: array of BlockPermission (optional) Permissions to take on `image` in addition to what is necessary anyway (which depends on how the blkdebug node is used). Defaults to none. (since 5.0)

unshare-child-perms: array of BlockPermission (optional) Permissions not to share on `image` in addition to what cannot be shared anyway (which depends on how the blkdebug node is used). Defaults to none. (since 5.0)

Since

2.9

BlockdevOptionsBlklogwrites (Object)

Driver specific block device options for blklogwrites.

Members

file: BlockdevRef block device

log: BlockdevRef block device used to log writes to `file`

log-sector-size: int (optional) sector size used in logging writes to `file`, determines granularity of offsets and sizes of writes (default: 512)

log-append: boolean (optional) append to an existing log (default: false)

log-super-update-interval: int (optional) interval of write requests after which the log super block is updated to disk (default: 4096)

Since

3.0

BlockdevOptionsBlkverify (Object)

Driver specific block device options for blkverify.

Members

test: BlockdevRef block device to be tested

raw: BlockdevRef raw image used for verification

Since

2.9

BlockdevOptionsBlkreplay (Object)

Driver specific block device options for blkreplay.

Members

image: BlockdevRef disk image which should be controlled with blkreplay

Since

4.2

QuorumReadPattern (Enum)

An enumeration of quorum read patterns.

Values

quorum read all the children and do a quorum vote on reads

fifo read only from the first child that has not failed

Since

2.9

BlockdevOptionsQuorum (Object)

Driver specific block device options for Quorum

Members

blkverify: **boolean** (optional)

true if the driver must print content mismatch set to false by default

children: **array of BlockdevRef** the children block devices to use

vote-threshold: **int** the vote limit under which a read will fail

rewrite-corrupted: **boolean** (optional) rewrite corrupted data when quorum is reached (Since 2.1)

read-pattern: **QuorumReadPattern** (optional) choose read pattern and set to quorum by default (Since 2.2)

Since

2.9

BlockdevOptionsGluster (Object)

Driver specific block device options for Gluster

Members

volume: **string** name of gluster volume where VM image resides

path: **string** absolute path to image file in gluster volume

server: **array of SocketAddress** gluster servers description

debug: **int** (optional) libgfapi log level (default '4' which is Error) (Since 2.8)

logfile: **string** (optional) libgfapi log file (default /dev/stderr) (Since 2.8)

Since

2.9

IscsiTransport (Enum)

An enumeration of libiscsi transport types

Values

tcp Not documented

iser Not documented

Since

2.9

IscsiHeaderDigest (Enum)

An enumeration of header digests supported by libiscsi

Values

crc32c Not documented

none Not documented

crc32c-none Not documented

none-crc32c Not documented

Since

2.9

BlockdevOptionsIscsi (Object)

Members

transport: IscsiTransport The iscsi transport type

portal: string The address of the iscsi portal

target: string The target iqname

lun: int (optional) LUN to connect to. Defaults to 0.

user: string (optional) User name to log in with. If omitted, no CHAP authentication is performed.

password-secret: string (optional) The ID of a QCryptoSecret object providing the password for the login.
This option is required if `user` is specified.

initiator-name: string (optional) The iqname we want to identify to the target as. If this option is not specified, an initiator name is generated automatically.

header-digest: IscsiHeaderDigest (optional) The desired header digest. Defaults to none-crc32c.

timeout: int (optional) Timeout in seconds after which a request will timeout. 0 means no timeout and is the default.

Driver specific block device options for iscsi

Since

2.9

RbdAuthMode (Enum)

Values

cephx Not documented

none Not documented

Since

3.0

BlockdevOptionsRbd (Object)

Members

pool: string Ceph pool name.

namespace: string (optional) Rados namespace name in the Ceph pool. (Since 5.0)

image: string Image name in the Ceph pool.

conf: string (optional) path to Ceph configuration file. Values in the configuration file will be overridden by options specified via QAPI.

snapshot: string (optional) Ceph snapshot name.

user: string (optional) Ceph id name.

auth-client-required: array of RbdAuthMode (optional) Acceptable authentication modes. This maps to Ceph configuration option “auth_client_required”. (Since 3.0)

key-secret: string (optional) ID of a QCryptoSecret object providing a key for cephx authentication. This maps to Ceph configuration option “key”. (Since 3.0)

server: array of InetSocketAddressBase (optional) Monitor host address and port. This maps to the “mon_host” Ceph option.

Since

2.9

BlockdevOptionsSheepdog (Object)

Driver specific block device options for sheepdog

Members

vdi: string Virtual disk image name

server: SocketAddress The Sheepdog server to connect to

snap-id: int (optional) Snapshot ID

tag: string (optional) Snapshot tag name

Only one of `snap-id` and `tag` may be present.

Since

2.9

ReplicationMode (Enum)

An enumeration of replication modes.

Values

primary Primary mode, the vm's state will be sent to secondary QEMU.

secondary Secondary mode, receive the vm's state from primary QEMU.

Since

2.9

If

`defined(CONFIG_REPLICATION)`

BlockdevOptionsReplication (Object)

Driver specific block device options for replication

Members

mode: ReplicationMode the replication mode

top-id: string (optional) In secondary mode, node name or device ID of the root node who owns the replication node chain. Must not be given in primary mode.

The members of **BlockdevOptionsGenericFormat**

Since

2.9

If

`defined(CONFIG_REPLICATION)`

NFSTransport (Enum)

An enumeration of NFS transport types

Values

inet TCP transport

Since

2.9

NFSServer (Object)

Captures the address of the socket

Members

type: **NFSTransport** transport type used for NFS (only TCP supported)

host: **string** host address for NFS server

Since

2.9

BlockdevOptionsNfs (Object)

Driver specific block device option for NFS

Members

server: **NFSServer** host address

path: **string** path of the image on the host

user: **int (optional)** UID value to use when talking to the server (defaults to 65534 on Windows and getuid() on unix)

group: **int (optional)** GID value to use when talking to the server (defaults to 65534 on Windows and getgid() in unix)

tcp-syn-count: **int (optional)** number of SYNs during the session establishment (defaults to libnfs default)

readahead-size: **int (optional)** set the readahead size in bytes (defaults to libnfs default)

page-cache-size: **int (optional)** set the pagecache size in bytes (defaults to libnfs default)

debug: **int (optional)** set the NFS debug level (max 2) (defaults to libnfs default)

Since

2.9

BlockdevOptionsCurlBase (Object)

Driver specific block device options shared by all protocols supported by the curl backend.

Members

url: **string** URL of the image file

readahead: **int** (**optional**) Size of the read-ahead cache; must be a multiple of 512 (defaults to 256 kB)

timeout: **int** (**optional**) Timeout for connections, in seconds (defaults to 5)

username: **string** (**optional**) Username for authentication (defaults to none)

password-secret: **string** (**optional**) ID of a QCryptoSecret object providing a password for authentication (defaults to no password)

proxy-username: **string** (**optional**) Username for proxy authentication (defaults to none)

proxy-password-secret: **string** (**optional**) ID of a QCryptoSecret object providing a password for proxy authentication (defaults to no password)

Since

2.9

BlockdevOptionsCurlHttp (Object)

Driver specific block device options for HTTP connections over the curl backend. URLs must start with “[http://](#)”.

Members

cookie: **string** (**optional**) List of cookies to set; format is “name1=content1; name2=content2;” as explained by CURLOPT_COOKIE(3). Defaults to no cookies.

cookie-secret: **string** (**optional**) ID of a QCryptoSecret object providing the cookie data in a secure way. See `cookie` for the format. (since 2.10)

The members of `BlockdevOptionsCurlBase`

Since

2.9

BlockdevOptionsCurlHttps (Object)

Driver specific block device options for HTTPS connections over the curl backend. URLs must start with “[https://](#)”.

Members

cookie: **string** (**optional**) List of cookies to set; format is “name1=content1; name2=content2;” as explained by CURLOPT_COOKIE(3). Defaults to no cookies.

sslverify: **boolean** (**optional**) Whether to verify the SSL certificate’s validity (defaults to true)

cookie-secret: **string** (**optional**) ID of a QCryptoSecret object providing the cookie data in a secure way. See `cookie` for the format. (since 2.10)

The members of `BlockdevOptionsCurlBase`

Since

2.9

`BlockdevOptionsCurlFtp` (Object)

Driver specific block device options for FTP connections over the curl backend. URLs must start with “`ftp://`”.

Members

The members of `BlockdevOptionsCurlBase`

Since

2.9

`BlockdevOptionsCurlFtps` (Object)

Driver specific block device options for FTPS connections over the curl backend. URLs must start with “`ftps://`”.

Members

sslverify: boolean (optional) Whether to verify the SSL certificate’s validity (defaults to true)

The members of `BlockdevOptionsCurlBase`

Since

2.9

`BlockdevOptionsNbd` (Object)

Driver specific block device options for NBD.

Members

server: SocketAddress NBD server address

export: string (optional) export name

tls-creds: string (optional) TLS credentials ID

x-dirty-bitmap: string (optional) A metadata context name such as “`qemu:dirty-bitmap:NAME`” or “`qemu:allocation-depth`” to query in place of the traditional “`base:allocation`” block status (see `NBD_OPT_LIST_META_CONTEXT` in the NBD protocol; and yes, naming this option x-context would have made more sense) (since 3.0)

reconnect-delay: int (optional) On an unexpected disconnect, the nbd client tries to connect again until succeeding or encountering a serious error. During the first `reconnect-delay` seconds, all requests are paused and will be rerun on a successful reconnect. After that time, any delayed requests and all future requests before a successful reconnect will immediately fail. Default 0 (Since 4.2)

Since

2.9

BlockdevOptionsRaw (Object)

Driver specific block device options for the raw driver.

Members

offset: int (optional) position where the block device starts

size: int (optional) the assumed size of the device

The members of `BlockdevOptionsGenericFormat`

Since

2.9

BlockdevOptionsThrottle (Object)

Driver specific block device options for the throttle driver

Members

throttle-group: string the name of the throttle-group object to use. It must already exist.

file: BlockdevRef reference to or definition of the data source block device

Since

2.11

BlockdevOptionsCor (Object)

Driver specific block device options for the copy-on-read driver.

Members

bottom: string (optional) The name of a non-filter node (allocation-bearing layer) that limits the COR operations in the backing chain (inclusive), so that no data below this node will be copied by this filter. If option is absent, the limit is not applied, so that data from all backing layers may be copied.

The members of `BlockdevOptionsGenericFormat`

Since

6.0

BlockdevOptions (Object)

Options for creating a block device. Many options are available for all block devices, independent of the block driver:

Members

driver: BlockdevDriver block driver name

node-name: string (optional) the node name of the new node (Since 2.0). This option is required on the top level of `blockdev-add`. Valid node names start with an alphabetic character and may contain only alphanumeric characters, '-', '.' and '_'. Their maximum length is 31 characters.

discard: BlockdevDiscardOptions (optional) discard-related options (default: ignore)

cache: BlockdevCacheOptions (optional) cache-related options

read-only: boolean (optional) whether the block device should be read-only (default: false). Note that some block drivers support only read-only access, either generally or in certain configurations. In this case, the default value does not work and the option must be specified explicitly.

auto-read-only: boolean (optional) if true and `read-only` is false, QEMU may automatically decide not to open the image read-write as requested, but fall back to read-only instead (and switch between the modes later), e.g. depending on whether the image file is writable or whether a writing user is attached to the node (default: false, since 3.1)

detect-zeroes: BlockdevDetectZeroesOptions (optional) detect and optimize zero writes (Since 2.1) (default: off)

force-share: boolean (optional) force share all permission on added nodes. Requires `read-only=true`. (Since 2.10)

The members of `BlockdevOptionsBlkdebug` when `driver` is "blkdebug"

The members of `BlockdevOptionsBlklogwrites` when `driver` is "blklogwrites"

The members of `BlockdevOptionsBlkverify` when `driver` is "blkverify"

The members of `BlockdevOptionsBlkreplay` when `driver` is "blkreplay"

The members of `BlockdevOptionsGenericFormat` when `driver` is "bochs"

The members of `BlockdevOptionsGenericFormat` when `driver` is "cloop"

The members of `BlockdevOptionsGenericFormat` when `driver` is "compress"

The members of `BlockdevOptionsCor` when `driver` is "copy-on-read"

The members of `BlockdevOptionsGenericFormat` when driver is "dmg"

The members of `BlockdevOptionsFile` when driver is "file"

The members of `BlockdevOptionsCurlFtp` when driver is "ftp"

The members of `BlockdevOptionsCurlFtps` when driver is "ftps"

The members of `BlockdevOptionsGluster` when driver is "gluster"

The members of `BlockdevOptionsFile` when driver is "host_cdrom"

The members of `BlockdevOptionsFile` when driver is "host_device"

The members of `BlockdevOptionsCurlHttp` when driver is "http"

The members of `BlockdevOptionsCurlHttps` when driver is "https"

The members of `BlockdevOptionsIscsi` when driver is "iscsi"

The members of `BlockdevOptionsLUKS` when driver is "luks"

The members of `BlockdevOptionsNbd` when driver is "nbd"

The members of `BlockdevOptionsNfs` when driver is "nfs"

The members of `BlockdevOptionsNull` when driver is "null-aio"

The members of `BlockdevOptionsNull` when driver is "null-co"

The members of `BlockdevOptionsNVMe` when driver is "nvme"

The members of `BlockdevOptionsGenericFormat` when driver is "parallels"

The members of `BlockdevOptionsPreallocate` when driver is "preallocate"

The members of `BlockdevOptionsQcow2` when driver is "qcow2"

The members of `BlockdevOptionsQcow` when driver is "qcow"

The members of `BlockdevOptionsGenericCOWFormat` when driver is "qed"

The members of `BlockdevOptionsQuorum` when driver is "quorum"

The members of `BlockdevOptionsRaw` when driver is "raw"

The members of `BlockdevOptionsRbd` when driver is "rbd"

The members of `BlockdevOptionsReplication` when driver is "replication" (If: defined(CONFIG_REPLICAT

The members of `BlockdevOptionsSheepdog` when driver is "sheepdog"

The members of `BlockdevOptionsSsh` when driver is "ssh"

The members of `BlockdevOptionsThrottle` when driver is "throttle"

The members of `BlockdevOptionsGenericFormat` when driver is "vdi"

The members of `BlockdevOptionsGenericFormat` when driver is "vhdx"

The members of `BlockdevOptionsGenericCOWFormat` when driver is "vmdk"

The members of `BlockdevOptionsGenericFormat` when driver is "vpc"

The members of `BlockdevOptionsVVFAT` when driver is "vvfat"

Remaining options are determined by the block driver.

Since

2.9

BlockdevRef (Alternate)

Reference to a block device.

Members

definition: BlockdevOptions defines a new block device inline

reference: string references the ID of an existing block device

Since

2.9

BlockdevRefOrNull (Alternate)

Reference to a block device.

Members

definition: BlockdevOptions defines a new block device inline

reference: string references the ID of an existing block device. An empty string means that no block device should be referenced. Deprecated; use null instead.

null: null No block device should be referenced (since 2.10)

Since

2.9

blockdev-add (Command)

Creates a new block device.

Arguments

The members of **BlockdevOptions**

Since

2.9

Example

```
1.
-> { "execute": "blockdev-add",
    "arguments": {
        "driver": "qcow2",
        "node-name": "test1",
        "file": {
            "driver": "file",
            "filename": "test.qcow2"
        }
    }
}
<- { "return": {} }

2.
-> { "execute": "blockdev-add",
    "arguments": {
        "driver": "qcow2",
        "node-name": "node0",
        "discard": "unmap",
        "cache": {
            "direct": true
        },
        "file": {
            "driver": "file",
            "filename": "/tmp/test.qcow2"
        },
        "backing": {
            "driver": "raw",
            "file": {
                "driver": "file",
                "filename": "/dev/fdset/4"
            }
        }
    }
}
<- { "return": {} }
```

x-blockdev-reopen (Command)

Reopens a block device using the given set of options. Any option not specified will be reset to its default value regardless of its previous status. If an option cannot be changed or a particular driver does not support reopening then the command will return an error.

The top-level `node-name` option (from `BlockdevOptions`) must be specified and is used to select the block device to be reopened. Other `node-name` options must be either omitted or set to the current name of the appropriate node. This command won't change any node name and any attempt to do it will result in an error.

In the case of options that refer to child nodes, the behavior of this command depends on the value:

- 1) A set of options (`BlockdevOptions`): the child is reopened with the specified set of options.
- 2) A reference to the current child: the child is reopened using its existing set of options.
- 3) A reference to a different node: the current child is replaced with the specified one.

- 4) NULL: the current child (if any) is detached.

Options (1) and (2) are supported in all cases, but at the moment only `backing` allows replacing or detaching an existing child.

Unlike with `blockdev-add`, the `backing` option must always be present unless the node being reopened does not have a backing file and its image does not have a default backing file name as part of its metadata.

Arguments

The members of `BlockdevOptions`

Since

4.0

`blockdev-del` (Command)

Deletes a block device that has been added using `blockdev-add`. The command will fail if the node is attached to a device or is otherwise being used.

Arguments

node-name: `string` Name of the graph node to delete.

Since

2.9

Example

```
-> { "execute": "blockdev-add",
      "arguments": {
        "driver": "qcow2",
        "node-name": "node0",
        "file": {
          "driver": "file",
          "filename": "test.qcow2"
        }
      }
    }
<- { "return": {} }

-> { "execute": "blockdev-del",
      "arguments": { "node-name": "node0" }
    }
<- { "return": {} }
```

BlockdevCreateOptionsFile (Object)

Driver specific image creation options for file.

Members

filename: string Filename for the new image file

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, falloc (if defined CONFIG_POSIX_FALLOCATE), full (if defined CONFIG_POSIX))

nocow: boolean (optional) Turn off copy-on-write (valid only on btrfs; default: off)

extent-size-hint: int (optional) Extent size hint to add to the image file; 0 for not adding an extent size hint (default: 1 MB, since 5.1)

Since

2.12

BlockdevCreateOptionsGluster (Object)

Driver specific image creation options for gluster.

Members

location: BlockdevOptionsGluster Where to store the new image file

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, falloc (if defined CONFIG_GLUSTERFS_FALLOCATE), full (if defined CONFIG_GLUSTERFS_ZEROFILL))

Since

2.12

BlockdevCreateOptionsLUKS (Object)

Driver specific image creation options for LUKS.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (since: 4.2) (default: off; allowed values: off, metadata, falloc, full)

The members of `QCryptoBlockCreateOptionsLUKS`

Since

2.12

`BlockdevCreateOptionsNfs` (Object)

Driver specific image creation options for NFS.

Members

location: `BlockdevOptionsNfs` Where to store the new image file

size: `int` Size of the virtual disk in bytes

Since

2.12

`BlockdevCreateOptionsParallels` (Object)

Driver specific image creation options for parallels.

Members

file: `BlockdevRef` Node to create the image format on

size: `int` Size of the virtual disk in bytes

cluster-size: `int` (optional) Cluster size in bytes (default: 1 MB)

Since

2.12

`BlockdevCreateOptionsQcow` (Object)

Driver specific image creation options for qcow.

Members

file: `BlockdevRef` Node to create the image format on

size: `int` Size of the virtual disk in bytes

backing-file: `string` (optional) File name of the backing file if a backing file should be used

encrypt: `QCryptoBlockCreateOptions` (optional) Encryption options if the image should be encrypted

Since

2.12

BlockdevQcow2Version (Enum)

Values

v2 The original QCOW2 format as introduced in qemu 0.10 (version 2)

v3 The extended QCOW2 format as introduced in qemu 1.1 (version 3)

Since

2.12

Qcow2CompressionType (Enum)

Compression type used in qcow2 image file

Values

zlib zlib compression, see <<http://zlib.net/>>

zstd (If: defined(CONFIG_ZSTD)) zstd compression, see <<http://github.com/facebook/zstd>>

Since

5.1

BlockdevCreateOptionsQcow2 (Object)

Driver specific image creation options for qcow2.

Members

file: BlockdevRef Node to create the image format on

data-file: BlockdevRef (optional) Node to use as an external data file in which all guest data is stored so that only metadata remains in the qcow2 file (since: 4.0)

data-file-raw: boolean (optional) True if the external data file must stay valid as a standalone (read-only) raw image without looking at qcow2 metadata (default: false; since: 4.0)

extended-l2: boolean (optional) True to make the image have extended L2 entries (default: false; since 5.2)

size: int Size of the virtual disk in bytes

version: BlockdevQcow2Version (optional) Compatibility level (default: v3)

backing-file: string (optional) File name of the backing file if a backing file should be used

backing-fmt: BlockdevDriver (optional) Name of the block driver to use for the backing file

encrypt: QCryptoBlockCreateOptions (optional) Encryption options if the image should be encrypted

cluster-size: int (optional) qcow2 cluster size in bytes (default: 65536)

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, falloc, full, metadata)

lazy-refcounts: boolean (optional) True if refcounts may be updated lazily (default: off)

refcount-bits: int (optional) Width of reference counts in bits (default: 16)

compression-type: Qcow2CompressionType (optional) The image cluster compression method (default: zlib, since 5.1)

Since

2.12

BlockdevCreateOptionsQed (Object)

Driver specific image creation options for qed.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

backing-file: string (optional) File name of the backing file if a backing file should be used

backing-fmt: BlockdevDriver (optional) Name of the block driver to use for the backing file

cluster-size: int (optional) Cluster size in bytes (default: 65536)

table-size: int (optional) L1/L2 table size (in clusters)

Since

2.12

BlockdevCreateOptionsRbd (Object)

Driver specific image creation options for rbd/Ceph.

Members

location: BlockdevOptionsRbd Where to store the new image file. This location cannot point to a snapshot.

size: int Size of the virtual disk in bytes

cluster-size: int (optional) RBD object size

Since

2.12

BlockdevVmdkSubformat (Enum)

Subformat options for VMDK images

Values

monolithicSparse Single file image with sparse cluster allocation

monolithicFlat Single flat data image and a descriptor file

twoGbMaxExtentSparse Data is split into 2GB (per virtual LBA) sparse extent files, in addition to a descriptor file

twoGbMaxExtentFlat Data is split into 2GB (per virtual LBA) flat extent files, in addition to a descriptor file

streamOptimized Single file image sparse cluster allocation, optimized for streaming over network.

Since

4.0

BlockdevVmdkAdapterType (Enum)

Adapter type info for VMDK images

Values

ide Not documented

buslogic Not documented

lsilogic Not documented

legacyESX Not documented

Since

4.0

BlockdevCreateOptionsVmdk (Object)

Driver specific image creation options for VMDK.

Members

file: BlockdevRef Where to store the new image file. This refers to the image file for monolithicSparse and streamOptimized format, or the descriptor file for other formats.

size: int Size of the virtual disk in bytes

extents: array of BlockdevRef (optional) Where to store the data extents. Required for monolithicFlat, twoGbMaxExtentSparse and twoGbMaxExtentFlat formats. For monolithicFlat, only one entry is required; for twoGbMaxExtent* formats, the number of entries required is calculated as $\text{extent_number} = \text{virtual_size} / 2\text{GB}$. Providing more extents than will be used is an error.

subformat: BlockdevVmdkSubformat (optional) The subformat of the VMDK image. Default: “monolithic-Sparse”.

backing-file: string (optional) The path of backing file. Default: no backing file is used.

adapter-type: BlockdevVmdkAdapterType (optional) The adapter type used to fill in the descriptor. Default: ide.

hwversion: string (optional) Hardware version. The meaningful options are “4” or “6”. Default: “4”.

zeroed-grain: boolean (optional) Whether to enable zeroed-grain feature for sparse subformats. Default: false.

Since

4.0

SheepdogRedundancyType (Enum)

Values

full Create a fully replicated vdi with x copies

erasure-coded Create an erasure coded vdi with x data strips and y parity strips

Since

2.12

SheepdogRedundancyFull (Object)

Members

copies: int Number of copies to use (between 1 and 31)

Since

2.12

SheepdogRedundancyErasureCoded (Object)

Members

data-strips: int Number of data strips to use (one of {2,4,8,16})

parity-strips: int Number of parity strips to use (between 1 and 15)

Since

2.12

SheepdogRedundancy (Object)

Members

type: SheepdogRedundancyType Not documented

The members of `SheepdogRedundancyFull` when `type` is "full"

The members of `SheepdogRedundancyErasureCoded` when `type` is "erasure-coded"

Since

2.12

BlockdevCreateOptionsSheepdog (Object)

Driver specific image creation options for Sheepdog.

Members

location: BlockdevOptionsSheepdog Where to store the new image file

size: int Size of the virtual disk in bytes

backing-file: string (optional) File name of a base image

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, full)

redundancy: SheepdogRedundancy (optional) Redundancy of the image

object-size: int (optional) Object size of the image

Since

2.12

BlockdevCreateOptionsSsh (Object)

Driver specific image creation options for SSH.

Members

location: BlockdevOptionsSsh Where to store the new image file

size: int Size of the virtual disk in bytes

Since

2.12

BlockdevCreateOptionsVdi (Object)

Driver specific image creation options for VDI.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

preallocation: PreallocMode (optional) Preallocation mode for the new image (default: off; allowed values: off, metadata)

Since

2.12

BlockdevVhdxSubformat (Enum)**Values**

dynamic Growing image file

fixed Preallocated fixed-size image file

Since

2.12

BlockdevCreateOptionsVhdx (Object)

Driver specific image creation options for vhdx.

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

log-size: int (optional) Log size in bytes, must be a multiple of 1 MB (default: 1 MB)

block-size: int (optional) Block size in bytes, must be a multiple of 1 MB and not larger than 256 MB (default: automatically choose a block size depending on the image size)

subformat: BlockdevVhdxSubformat (optional) vhdx subformat (default: dynamic)

block-state-zero: boolean (optional) Force use of payload blocks of type 'ZERO'. Non-standard, but default. Do not set to 'off' when using 'qemu-img convert' with subformat=dynamic.

Since

2.12

BlockdevVpcSubformat (Enum)

Values

dynamic Growing image file

fixed Preallocated fixed-size image file

Since

2.12

BlockdevCreateOptionsVpc (Object)

Driver specific image creation options for vpc (VHD).

Members

file: BlockdevRef Node to create the image format on

size: int Size of the virtual disk in bytes

subformat: BlockdevVpcSubformat (optional) vhdx subformat (default: dynamic)

force-size: boolean (optional) Force use of the exact byte size instead of rounding to the next size that can be represented in CHS geometry (default: false)

Since

2.12

BlockdevCreateOptions (Object)

Options for creating an image format on a given node.

Members

driver: BlockdevDriver block driver to create the image format

The members of **BlockdevCreateOptionsFile** when driver is "file"

The members of **BlockdevCreateOptionsGluster** when driver is "gluster"

The members of **BlockdevCreateOptionsLUKS** when driver is "luks"

The members of **BlockdevCreateOptionsNfs** when driver is "nfs"

The members of **BlockdevCreateOptionsParallels** when driver is "parallels"

The members of **BlockdevCreateOptionsQcow** when driver is "qcow"

The members of **BlockdevCreateOptionsQcow2** when driver is "qcow2"

The members of **BlockdevCreateOptionsQed** when driver is "qed"

The members of **BlockdevCreateOptionsRbd** when driver is "rbd"

The members of **BlockdevCreateOptionsSheepdog** when driver is "sheepdog"

The members of **BlockdevCreateOptionsSsh** when driver is "ssh"

The members of **BlockdevCreateOptionsVdi** when driver is "vdi"

The members of **BlockdevCreateOptionsVhdx** when driver is "vhdx"

The members of **BlockdevCreateOptionsVmdk** when driver is "vmdk"

The members of **BlockdevCreateOptionsVpc** when driver is "vpc"

Since

2.12

blockdev-create (Command)

Starts a job to create an image format on a given node. The job is automatically finalized, but a manual job-dismiss is required.

Arguments

job-id: string Identifier for the newly created job.

options: BlockdevCreateOptions Options for the image creation.

Since

3.0

BlockdevAmendOptionsLUKS (Object)

Driver specific image amend options for LUKS.

Members

The members of `QCryptoBlockAmendOptionsLUKS`

Since

5.1

BlockdevAmendOptionsQcow2 (Object)

Driver specific image amend options for qcow2. For now, only encryption options can be amended
`encrypt` Encryption options to be amended

Members

encrypt: `QCryptoBlockAmendOptions` (optional) Not documented

Since

5.1

BlockdevAmendOptions (Object)

Options for amending an image format

Members

driver: `BlockdevDriver` Block driver of the node to amend.

The members of `BlockdevAmendOptionsLUKS` when `driver` is "luks"

The members of `BlockdevAmendOptionsQcow2` when `driver` is "qcow2"

Since

5.1

x-blockdev-amend (Command)

Starts a job to amend format specific options of an existing open block device The job is automatically finalized, but a manual job-dismiss is required.

Arguments

job-id: string Identifier for the newly created job.

node-name: string Name of the block node to work on

options: BlockdevAmendOptions Options (driver specific)

force: boolean (optional) Allow unsafe operations, format specific For luks that allows erase of the last active keyslot (permanent loss of data), and replacement of an active keyslot (possible loss of data if IO error happens)

Since

5.1

BlockErrorAction (Enum)

An enumeration of action that has been taken when a DISK I/O occurs

Values

ignore error has been ignored

report error has been reported to the device

stop error caused VM to be stopped

Since

2.1

BLOCK_IMAGE_CORRUPTED (Event)

Emitted when a disk image is being marked corrupt. The image can be identified by its device or node name. The ‘device’ field is always present for compatibility reasons, but it can be empty (“”) if the image does not have a device name associated.

Arguments

device: string device name. This is always present for compatibility reasons, but it can be empty (“”) if the image does not have a device name associated.

node-name: string (optional) node name (Since: 2.4)

msg: string informative message for human consumption, such as the kind of corruption being detected. It should not be parsed by machine as it is not guaranteed to be stable

offset: int (optional) if the corruption resulted from an image access, this is the host’s access offset into the image

size: int (optional) if the corruption resulted from an image access, this is the access size

fatal: boolean if set, the image is marked corrupt and therefore unusable after this event and must be repaired (Since 2.2; before, every BLOCK_IMAGE_CORRUPTED event was fatal)

Note

If action is “stop”, a STOP event will eventually follow the BLOCK_IO_ERROR event.

Example

```
<- { "event": "BLOCK_IMAGE_CORRUPTED",
      "data": { "device": "ide0-hd0", "node-name": "node0",
                 "msg": "Prevented active L1 table overwrite", "offset": 196608,
                 "size": 65536 },
      "timestamp": { "seconds": 1378126126, "microseconds": 966463 } }
```

Since

1.7

BLOCK_IO_ERROR (Event)

Emitted when a disk I/O error occurs

Arguments

device: string device name. This is always present for compatibility reasons, but it can be empty (“”) if the image does not have a device name associated.

node-name: string (optional) node name. Note that errors may be reported for the root node that is directly attached to a guest device rather than for the node where the error occurred. The node name is not present if the drive is empty. (Since: 2.8)

operation: IoOperationType I/O operation

action: BlockErrorAction action that has been taken

nospace: boolean (optional) true if I/O error was caused due to a no-space condition. This key is only present if query-block’s io-status is present, please see query-block documentation for more information (since: 2.2)

reason: string human readable string describing the error cause. (This field is a debugging aid for humans, it should not be parsed by applications) (since: 2.2)

Note

If action is “stop”, a STOP event will eventually follow the BLOCK_IO_ERROR event

Since

0.13

Example

```
<- { "event": "BLOCK_IO_ERROR",
      "data": { "device": "ide0-hd1",
                  "node-name": "#block212",
                  "operation": "write",
                  "action": "stop" },
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

BLOCK_JOB_COMPLETED (Event)

Emitted when a block job has completed

Arguments

type: JobType job type

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: int maximum progress value

offset: int current progress value. On success this is equal to len. On failure this is less than len

speed: int rate limit, bytes per second

error: string (optional) error message. Only present on failure. This field contains a human-readable error message. There are no semantics other than that streaming has failed and clients should not try to interpret the error string

Since

1.1

Example

```
<- { "event": "BLOCK_JOB_COMPLETED",
      "data": { "type": "stream", "device": "virtio-disk0",
                  "len": 10737418240, "offset": 10737418240,
                  "speed": 0 },
      "timestamp": { "seconds": 1267061043, "microseconds": 959568 } }
```

BLOCK_JOB_CANCELLED (Event)

Emitted when a block job has been cancelled

Arguments

type: JobType job type

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: int maximum progress value

offset: int current progress value. On success this is equal to len. On failure this is less than len

speed: int rate limit, bytes per second

Since

1.1

Example

```
<- { "event": "BLOCK_JOB_CANCELLED",
      "data": { "type": "stream", "device": "virtio-disk0",
                  "len": 10737418240, "offset": 134217728,
                  "speed": 0 },
      "timestamp": { "seconds": 1267061043, "microseconds": 959568 } }
```

BLOCK_JOB_ERROR (Event)

Emitted when a block job encounters an error

Arguments

device: string The job identifier. Originally the device name but other values are allowed since QEMU 2.7

operation: IoOperationType I/O operation

action: BlockErrorAction action that has been taken

Since

1.3

Example

```
<- { "event": "BLOCK_JOB_ERROR",
      "data": { "device": "ide0-hd1",
                  "operation": "write",
                  "action": "stop" },
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

BLOCK_JOB_READY (Event)

Emitted when a block job is ready to complete

Arguments

type: `JobType` job type

device: `string` The job identifier. Originally the device name but other values are allowed since QEMU 2.7

len: `int` maximum progress value

offset: `int` current progress value. On success this is equal to len. On failure this is less than len

speed: `int` rate limit, bytes per second

Note

The “ready to complete” status is always reset by a `BLOCK_JOB_ERROR` event

Since

1.3

Example

```
<- { "event": "BLOCK_JOB_READY",
      "data": { "device": "drive0", "type": "mirror", "speed": 0,
                 "len": 2097152, "offset": 2097152 }
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

BLOCK_JOB_PENDING (Event)

Emitted when a block job is awaiting explicit authorization to finalize graph changes via `block-job-finalize`. If this job is part of a transaction, it will not emit this event until the transaction has converged first.

Arguments

type: `JobType` job type

id: `string` The job identifier.

Since

2.12

Example

```
<- { "event": "BLOCK_JOB_WAITING",
      "data": { "device": "drive0", "type": "mirror" },
      "timestamp": { "seconds": 1265044230, "microseconds": 450486 } }
```

PreallocMode (Enum)

Preallocation mode of QEMU image file

Values

off no preallocation

metadata preallocate only for metadata

falloc like `full` preallocation but allocate disk space by `posix_fallocate()` rather than writing data.

full preallocate all data by writing it to the device to ensure disk space is really available. This data may or may not be zero, depending on the image format and storage. `full` preallocation also sets up metadata correctly.

Since

2.2

BLOCK_WRITE_THRESHOLD (Event)

Emitted when writes on block device reaches or exceeds the configured write threshold. For thin-provisioned devices, this means the device should be extended to avoid pausing for disk exhaustion. The event is one shot. Once triggered, it needs to be re-registered with another `block-set-write-threshold` command.

Arguments

node-name: string graph node name on which the threshold was exceeded.

amount-exceeded: int amount of data which exceeded the threshold, in bytes.

write-threshold: int last configured threshold, in bytes.

Since

2.3

block-set-write-threshold (Command)

Change the write threshold for a block drive. An event will be delivered if a write to this block drive crosses the configured threshold. The threshold is an offset, thus must be non-negative. Default is no write threshold. Setting the threshold to zero disables it.

This is useful to transparently resize thin-provisioned drives without the guest OS noticing.

Arguments

node-name: string graph node name on which the threshold must be set.

write-threshold: int configured threshold for the block device, bytes. Use 0 to disable the threshold.

Since

2.3

Example

```
-> { "execute": "block-set-write-threshold",
      "arguments": { "node-name": "mydev",
                     "write-threshold": 17179869184 } }
<- { "return": {} }
```

x-blockdev-change (Command)

Dynamically reconfigure the block driver state graph. It can be used to add, remove, insert or replace a graph node. Currently only the Quorum driver implements this feature to add or remove its child. This is useful to fix a broken quorum child.

If `node` is specified, it will be inserted under `parent`. `child` may not be specified in this case. If both `parent` and `child` are specified but `node` is not, `child` will be detached from `parent`.

Arguments

parent: string the id or name of the parent node.

child: string (optional) the name of a child under the given parent node.

node: string (optional) the name of the node that will be added.

Note

this command is experimental, and its API is not stable. It does not support all kinds of operations, all kinds of children, nor all block drivers.

FIXME Removing children from a quorum node means introducing gaps in the child indices. This cannot be represented in the 'children' list of `BlockdevOptionsQuorum`, as returned by `.bdrv_refresh_filename()`.

Warning: The data in a new quorum child MUST be consistent with that of the rest of the array.

Since

2.7

Example

```
1. Add a new node to a quorum
-> { "execute": "blockdev-add",
      "arguments": {
        "driver": "raw",
        "node-name": "new_node",
```

(continues on next page)

(continued from previous page)

```
        "file": { "driver": "file",
                  "filename": "test.raw" } } }
<- { "return": {} }
-> { "execute": "x-blockdev-change",
    "arguments": { "parent": "disk1",
                  "node": "new_node" } }
<- { "return": {} }

2. Delete a quorum's node
-> { "execute": "x-blockdev-change",
    "arguments": { "parent": "disk1",
                  "child": "children.1" } }
<- { "return": {} }
```

x-blockdev-set-iothread (Command)

Move node and its children into the `iothread`. If `iothread` is null then move node and its children into the main loop.

The node must not be attached to a BlockBackend.

Arguments

node-name: string the name of the block driver node

iothread: StrOrNull the name of the IOThread object or null for the main loop

force: boolean (optional) true if the node and its children should be moved when a BlockBackend is already attached

Note

this command is experimental and intended for test cases that need control over IOThreads only.

Since

2.12

Example

```
1. Move a node into an IOThread
-> { "execute": "x-blockdev-set-iothread",
    "arguments": { "node-name": "disk1",
                  "iothread": "iothread0" } }
<- { "return": {} }

2. Move a node into the main loop
-> { "execute": "x-blockdev-set-iothread",
    "arguments": { "node-name": "disk1",
```

(continues on next page)

(continued from previous page)

```

        "iothread": null } }
<- { "return": {} }

```

QuorumOpType (Enum)

An enumeration of the quorum operation types

Values

read read operation

write write operation

flush flush operation

Since

2.6

QUORUM_FAILURE (Event)

Emitted by the Quorum block driver if it fails to establish a quorum

Arguments

reference: string device name if defined else node name

sector-num: int number of the first sector of the failed read operation

sectors-count: int failed read operation sector count

Note

This event is rate-limited.

Since

2.0

Example

```

<- { "event": "QUORUM_FAILURE",
      "data": { "reference": "usrl", "sector-num": 345435, "sectors-count": 5 },
      "timestamp": { "seconds": 1344522075, "microseconds": 745528 } }

```

QUORUM_REPORT_BAD (Event)

Emitted to report a corruption of a Quorum file

Arguments

type: QuorumOpType quorum operation type (Since 2.6)

error: string (optional) error message. Only present on failure. This field contains a human-readable error message. There are no semantics other than that the block layer reported an error and clients should not try to interpret the error string.

node-name: string the graph node name of the block driver state

sector-num: int number of the first sector of the failed read operation

sectors-count: int failed read operation sector count

Note

This event is rate-limited.

Since

2.0

Example

```
1. Read operation
{ "event": "QUORUM_REPORT_BAD",
  "data": { "node-name": "node0", "sector-num": 345435, "sectors-count": 5,
            "type": "read" },
  "timestamp": { "seconds": 1344522075, "microseconds": 745528 } }

2. Flush operation
{ "event": "QUORUM_REPORT_BAD",
  "data": { "node-name": "node0", "sector-num": 0, "sectors-count": 2097120,
            "type": "flush", "error": "Broken pipe" },
  "timestamp": { "seconds": 1456406829, "microseconds": 291763 } }
```

BlockdevSnapshotInternal (Object)

Members

device: string the device name or node-name of a root node to generate the snapshot from

name: string the name of the internal snapshot to be created

Notes

In transaction, if `name` is empty, or any snapshot matching `name` exists, the operation will fail. Only some image formats support it, for example, `qcow2`, `rbd`, and `sheepdog`.

Since

1.7

`blockdev-snapshot-internal-sync` (Command)

Synchronously take an internal snapshot of a block device, when the format of the image used supports it. If the name is an empty string, or a snapshot with name already exists, the operation will fail.

For the arguments, see the documentation of `BlockdevSnapshotInternal`.

Returns

- nothing on success
- If `device` is not a valid block device, `GenericError`
- If any snapshot matching `name` exists, or `name` is empty, `GenericError`
- If the format of the image used does not support it, `BlockFormatFeatureNotSupported`

Since

1.7

Example

```
-> { "execute": "blockdev-snapshot-internal-sync",
    "arguments": { "device": "ide-hd0",
                  "name": "snapshot0" }
  }
<- { "return": {} }
```

`blockdev-snapshot-delete-internal-sync` (Command)

Synchronously delete an internal snapshot of a block device, when the format of the image used support it. The snapshot is identified by name or id or both. One of the name or id is required. Return `SnapshotInfo` for the successfully deleted snapshot.

Arguments

device: string the device name or node-name of a root node to delete the snapshot from

id: string (optional) optional the snapshot's ID to be deleted

name: **string (optional)** optional the snapshot's name to be deleted

Returns

- SnapshotInfo on success
- If `device` is not a valid block device, `GenericError`
- If snapshot not found, `GenericError`
- If the format of the image used does not support it, `BlockFormatFeatureNotSupported`
- If `id` and `name` are both not specified, `GenericError`

Since

1.7

Example

```
-> { "execute": "blockdev-snapshot-delete-internal-sync",
    "arguments": { "device": "ide-hd0",
                  "name": "snapshot0" }
}
<- { "return": {
    "id": "1",
    "name": "snapshot0",
    "vm-state-size": 0,
    "date-sec": 1000012,
    "date-nsec": 10,
    "vm-clock-sec": 100,
    "vm-clock-nsec": 20,
    "icount": 220414
  }
}
```

Block device exports

NbdServerOptions (Object)

Keep this type consistent with the `nbd-server-start` arguments. The only intended difference is using `SocketAddress` instead of `SocketAddressLegacy`.

Members

addr: **SocketAddress** Address on which to listen.

tls-creds: **string (optional)** ID of the TLS credentials object (since 2.6).

tls-authz: **string (optional)** ID of the QAuthZ authorization object used to validate the client's x509 distinguished name. This object is only resolved at time of use, so can be deleted and recreated on the fly while the NBD server is active. If missing, it will default to denying access (since 4.0).

max-connections: int (optional) The maximum number of connections to allow at the same time, 0 for unlimited. (since 5.2; default: 0)

Since

4.2

nbd-server-start (Command)

Start an NBD server listening on the given host and port. Block devices can then be exported using `nbd-server-add`. The NBD server will present them as named exports; for example, another QEMU instance could refer to them as “nbd:HOST:PORT:exportname=NAME”.

Keep this type consistent with the `NbdServerOptions` type. The only intended difference is using `SocketAddressLegacy` instead of `SocketAddress`.

Arguments

addr: SocketAddressLegacy Address on which to listen.

tls-creds: string (optional) ID of the TLS credentials object (since 2.6).

tls-authz: string (optional) ID of the QAuthZ authorization object used to validate the client’s x509 distinguished name. This object is only resolved at time of use, so can be deleted and recreated on the fly while the NBD server is active. If missing, it will default to denying access (since 4.0).

max-connections: int (optional) The maximum number of connections to allow at the same time, 0 for unlimited. (since 5.2; default: 0)

Returns

error if the server is already running.

Since

1.3

BlockExportOptionsNbdBase (Object)

An NBD block export (common options shared between `nbd-server-add` and the NBD branch of `block-export-add`).

Members

name: string (optional) Export name. If unspecified, the `device` parameter is used as the export name. (Since 2.12)

description: string (optional) Free-form description of the export, up to 4096 bytes. (Since 5.0)

Since

5.0

BlockExportOptionsNbd (Object)

An NBD block export (distinct options used in the NBD branch of block-export-add).

Members

bitmaps: array of string (optional) Also export each of the named dirty bitmaps reachable from `device`, so the NBD client can use `NBD_OPT_SET_META_CONTEXT` with the metadata context name “qemu:dirty-bitmap:BITMAP” to inspect each bitmap.

allocation-depth: boolean (optional) Also export the allocation depth map for `device`, so the NBD client can use `NBD_OPT_SET_META_CONTEXT` with the metadata context name “qemu:allocation-depth” to inspect allocation details. (since 5.2)

The members of **BlockExportOptionsNbdBase**

Since

5.2

BlockExportOptionsVhostUserBlk (Object)

A vhost-user-blk block export.

Members

addr: SocketAddress The vhost-user socket on which to listen. Both ‘unix’ and ‘fd’ `SocketAddress` types are supported. Passed fds must be UNIX domain sockets.

logical-block-size: int (optional) Logical block size in bytes. Defaults to 512 bytes.

num-queues: int (optional) Number of request virtqueues. Must be greater than 0. Defaults to 1.

Since

5.2

BlockExportOptionsFuse (Object)

Options for exporting a block graph node on some (file) mountpoint as a raw image.

Members

mountpoint: string Path on which to export the block device via FUSE. This must point to an existing regular file.

growable: boolean (optional) Whether writes beyond the EOF should grow the block node accordingly. (default: false)

Since

6.0

If

`defined(CONFIG_FUSE)`

NbdServerAddOptions (Object)

An NBD block export, per legacy `nbd-server-add` command.

Members

device: string The device name or node name of the node to be exported

writable: boolean (optional) Whether clients should be able to write to the device via the NBD connection (default false).

bitmap: string (optional) Also export a single dirty bitmap reachable from `device`, so the NBD client can use `NBD_OPT_SET_META_CONTEXT` with the metadata context name “`qemu:dirty-bitmap:BITMAP`” to inspect the bitmap (since 4.0).

The members of `BlockExportOptionsNbdBase`

Since

5.0

nbd-server-add (Command)

Export a block node to QEMU’s embedded NBD server.

The export name will be used as the id for the resulting block export.

Arguments

The members of `NbdServerAddOptions`

Features

deprecated This command is deprecated. Use `block-export-add` instead.

Returns

error if the server is not running, or export with the same name already exists.

Since

1.3

BlockExportRemoveMode (Enum)

Mode for removing a block export.

Values

safe Remove export if there are no existing connections, fail otherwise.

hard Drop all connections immediately and remove export.

Potential additional modes to be added in the future:

hide: Just hide export from new clients, leave existing connections as is. Remove export after all clients are disconnected.

soft: Hide export from new clients, answer with ESHUTDOWN for all further requests from existing clients.

Since

2.12

nbd-server-remove (Command)

Remove NBD export by name.

Arguments

name: string Block export id.

mode: BlockExportRemoveMode (optional) Mode of command operation. See `BlockExportRemoveMode` description. Default is 'safe'.

Features

deprecated This command is deprecated. Use `block-export-del` instead.

Returns

error if

- the server is not running
- export is not found
- mode is ‘safe’ and there are existing connections

Since

2.12

`nbd-server-stop` (Command)

Stop QEMU’s embedded NBD server, and unregister all devices previously added via `nbd-server-add`.

Since

1.3

`BlockExportType` (Enum)

An enumeration of block export types

Values

`nbd` NBD export

`vhost-user-blk` vhost-user-blk export (since 5.2)

`fuse` (If: `defined(CONFIG_FUSE)`) FUSE export (since: 6.0)

Since

4.2

`BlockExportOptions` (Object)

Describes a block export, i.e. how single node should be exported on an external interface.

Members

id: `string` A unique identifier for the block export (across all export types)

node-name: `string` The node name of the block node to be exported (since: 5.2)

writable: `boolean` (optional) True if clients should be able to write to the export (default false)

writethrough: boolean (optional) If true, caches are flushed after every write request to the export before completion is signalled. (since: 5.2; default: false)

iothread: string (optional) The name of the iothread object where the export will run. The default is to use the thread currently associated with the block node. (since: 5.2)

fixed-iothread: boolean (optional) True prevents the block node from being moved to another thread while the export is active. If true and `iothread` is given, export creation fails if the block node cannot be moved to the iothread. The default is false. (since: 5.2)

type: BlockExportType Not documented

The members of `BlockExportOptionsNbd` when type is "nbd"

The members of `BlockExportOptionsVhostUserBlk` when type is "vhost-user-blk"

The members of `BlockExportOptionsFuse` when type is "fuse" (If: `defined(CONFIG_FUSE)`)

Since

4.2

`block-export-add` (Command)

Creates a new block export.

Arguments

The members of `BlockExportOptions`

Since

5.2

`block-export-del` (Command)

Request to remove a block export. This drops the user's reference to the export, but the export may still stay around after this command returns until the shutdown of the export has completed.

Arguments

id: string Block export id.

mode: BlockExportRemoveMode (optional) Mode of command operation. See `BlockExportRemoveMode` description. Default is 'safe'.

Returns

Error if the export is not found or mode is 'safe' and the export is still in use (e.g. by existing client connections)

Since

5.2

BLOCK_EXPORT_DELETED (Event)

Emitted when a block export is removed and its id can be reused.

Arguments

id: **string** Block export id.

Since

5.2

BlockExportInfo (Object)

Information about a single block export.

Members

id: **string** The unique identifier for the block export

type: **BlockExportType** The block export type

node-name: **string** The node name of the block node that is exported

shutting-down: **boolean** True if the export is shutting down (e.g. after a block-export-del command, but before the shutdown has completed)

Since

5.2

query-block-exports (Command)**Returns**

A list of BlockExportInfo describing all block exports

Since

5.2

4.9.5 Character devices

ChardevInfo (Object)

Information about a character device.

Members

label: **string** the label of the character device

filename: **string** the filename of the character device

frontend-open: **boolean** shows whether the frontend device attached to this backend (eg. with the `chardev=...` option) is in open or closed state (since 2.1)

Notes

`filename` is encoded using the QEMU command line character device encoding. See the QEMU man page for details.

Since

0.14

query-chardev (Command)

Returns information about current character devices.

Returns

a list of ChardevInfo

Since

0.14

Example

```
-> { "execute": "query-chardev" }
<- {
  "return": [
    {
      "label": "charchannel0",
      "filename": "unix:/var/lib/libvirt/qemu/seabios.rhel6.agent,server",
      "frontend-open": false
    },
    {
      "label": "charmonitor",
```

(continues on next page)

(continued from previous page)

```

        "filename": "unix:/var/lib/libvirt/qemu/seabios.rhel6.monitor,server",
        "frontend-open": true
    },
    {
        "label": "charserial0",
        "filename": "pty:/dev/pts/2",
        "frontend-open": true
    }
]
}

```

ChardevBackendInfo (Object)

Information about a character device backend

Members

name: string The backend name

Since

2.0

query-chardev-backends (Command)

Returns information about character device backends.

Returns

a list of ChardevBackendInfo

Since

2.0

Example

```

-> { "execute": "query-chardev-backends" }
<- {
    "return": [
        {
            "name": "udp"
        },
        {
            "name": "tcp"
        },
    ],
}

```

(continues on next page)

(continued from previous page)

```
{
  "name": "unix"
},
{
  "name": "spiceport"
}
]
```

DataFormat (Enum)

An enumeration of data format.

Values

utf8 Data is a UTF-8 string (RFC 3629)

base64 Data is Base64 encoded binary (RFC 3548)

Since

1.4

ringbuf-write (Command)

Write to a ring buffer character device.

Arguments

device: string the ring buffer character device name

data: string data to write

format: DataFormat (optional) data encoding (default 'utf8').

- base64: data must be base64 encoded text. Its binary decoding gets written.
- utf8: data's UTF-8 encoding is written
- data itself is always Unicode regardless of format, like any other string.

Returns

Nothing on success

Since

1.4

Example

```
-> { "execute": "ringbuf-write",
      "arguments": { "device": "foo",
                     "data": "abcdefgh",
                     "format": "utf8" } }
<- { "return": {} }
```

ringbuf-read (Command)

Read from a ring buffer character device.

Arguments

device: string the ring buffer character device name

size: int how many bytes to read at most

format: DataFormat (optional) data encoding (default 'utf8').

- base64: the data read is returned in base64 encoding.
- utf8: the data read is interpreted as UTF-8. Bug: can screw up when the buffer contains invalid UTF-8 sequences, NUL characters, after the ring buffer lost data, and when reading stops because the size limit is reached.
- The return value is always Unicode regardless of format, like any other string.

Returns

data read from the device

Since

1.4

Example

```
-> { "execute": "ringbuf-read",
      "arguments": { "device": "foo",
                     "size": 1000,
                     "format": "utf8" } }
<- { "return": "abcdefgh" }
```

ChardevCommon (Object)

Configuration shared across all chardev backends

Members

logfile: string (optional) The name of a logfile to save output

logappend: boolean (optional) true to append instead of truncate (default to false to truncate)

Since

2.6

ChardevFile (Object)

Configuration info for file chardevs.

Members

in: string (optional) The name of the input file

out: string The name of the output file

append: boolean (optional) Open the file in append mode (default false to truncate) (Since 2.6)

The members of ChardevCommon

Since

1.4

ChardevHostdev (Object)

Configuration info for device and pipe chardevs.

Members

device: string The name of the special file for the device, i.e. /dev/ttyS0 on Unix or COM1: on Windows

The members of ChardevCommon

Since

1.4

ChardevSocket (Object)

Configuration info for (stream) socket chardevs.

Members

addr: SocketAddressLegacy socket address to listen on (server=true) or connect to (server=false)

tls-creds: string (optional) the ID of the TLS credentials object (since 2.6)

tls-authz: string (optional) the ID of the QAuthZ authorization object against which the client's x509 distinguished name will be validated. This object is only resolved at time of use, so can be deleted and recreated on the fly while the chardev server is active. If missing, it will default to denying access (since 4.0)

server: boolean (optional) create server socket (default: true)

wait: boolean (optional) wait for incoming connection on server sockets (default: false). Silently ignored with server: false. This use is deprecated.

nodelay: boolean (optional) set TCP_NODELAY socket option (default: false)

telnet: boolean (optional) enable telnet protocol on server sockets (default: false)

tn3270: boolean (optional) enable tn3270 protocol on server sockets (default: false) (Since: 2.10)

websocket: boolean (optional) enable websocket protocol on server sockets (default: false) (Since: 3.1)

reconnect: int (optional) For a client socket, if a socket is disconnected, then attempt a reconnect after the given number of seconds. Setting this to zero disables this function. (default: 0) (Since: 2.2)

The members of ChardevCommon

Since

1.4

ChardevUdp (Object)

Configuration info for datagram socket chardevs.

Members

remote: SocketAddressLegacy remote address

local: SocketAddressLegacy (optional) local address

The members of ChardevCommon

Since

1.5

ChardevMux (Object)

Configuration info for mux chardevs.

Members

chardev: **string** name of the base chardev.

The members of `ChardevCommon`

Since

1.5

`ChardevStdio` (Object)

Configuration info for stdio chardevs.

Members

signal: **boolean (optional)** Allow signals (such as SIGINT triggered by ^C) be delivered to qemu. Default: true.

The members of `ChardevCommon`

Since

1.5

`ChardevSpiceChannel` (Object)

Configuration info for spice vm channel chardevs.

Members

type: **string** kind of channel (for example vdagent).

The members of `ChardevCommon`

Since

1.5

If

`defined(CONFIG_SPICE)`

`ChardevSpicePort` (Object)

Configuration info for spice port chardevs.

Members

fqdn: string name of the channel (see docs/spice-port-fqdn.txt)

The members of `ChardevCommon`

Since

1.5

If

`defined(CONFIG_SPICE)`

ChardevVC (Object)

Configuration info for virtual console chardevs.

Members

width: int (optional) console width, in pixels

height: int (optional) console height, in pixels

cols: int (optional) console width, in chars

rows: int (optional) console height, in chars

The members of `ChardevCommon`

Since

1.5

ChardevRingbuf (Object)

Configuration info for ring buffer chardevs.

Members

size: int (optional) ring buffer size, must be power of two, default is 65536

The members of `ChardevCommon`

Since

1.5

ChardevBackend (Object)

Configuration info for the new chardev backend.

Members

type One of file, serial, parallel, pipe, socket, udp, pty, null, mux, msmouse, wctablet, braille, testdev, stdio, console, spicevmc, spiceport, vc, ringbuf, memory

data: ChardevFile when type is "file"

data: ChardevHostdev when type is "serial"

data: ChardevHostdev when type is "parallel"

data: ChardevHostdev when type is "pipe"

data: ChardevSocket when type is "socket"

data: ChardevUdp when type is "udp"

data: ChardevCommon when type is "pty"

data: ChardevCommon when type is "null"

data: ChardevMux when type is "mux"

data: ChardevCommon when type is "msmouse"

data: ChardevCommon when type is "wctablet"

data: ChardevCommon when type is "braille"

data: ChardevCommon when type is "testdev"

data: ChardevStdio when type is "stdio"

data: ChardevCommon when type is "console"

data: ChardevSpiceChannel when type is "spicevmc" (If: defined(CONFIG_SPICE))

data: ChardevSpicePort when type is "spiceport" (If: defined(CONFIG_SPICE))

data: ChardevVC when type is "vc"

data: ChardevRingbuf when type is "ringbuf"

data: ChardevRingbuf when type is "memory"

Since

1.4 (testdev since 2.2, wctablet since 2.9)

ChardevReturn (Object)

Return info about the chardev backend just created.

Members

pty: **string** (optional) name of the slave pseudoterminal device, present if and only if a chardev of type 'pty' was created

Since

1.4

chardev-add (Command)

Add a character device backend

Arguments

id: **string** the chardev's ID, must be unique

backend: **ChardevBackend** backend type and parameters

Returns

ChardevReturn.

Since

1.4

Example

```
-> { "execute" : "chardev-add",
      "arguments" : { "id" : "foo",
                      "backend" : { "type" : "null", "data" : {} } } }
<- { "return": {} }

-> { "execute" : "chardev-add",
      "arguments" : { "id" : "bar",
                      "backend" : { "type" : "file",
                                    "data" : { "out" : "/tmp/bar.log" } } } }
<- { "return": {} }

-> { "execute" : "chardev-add",
      "arguments" : { "id" : "baz",
                      "backend" : { "type" : "pty", "data" : {} } } }
<- { "return": { "pty" : "/dev/pty/42" } }
```

chardev-change (Command)

Change a character device backend

Arguments

id: **string** the chardev's ID, must exist

backend: **ChardevBackend** new backend type and parameters

Returns

ChardevReturn.

Since

2.10

Example

```
-> { "execute" : "chardev-change",
      "arguments" : { "id" : "baz",
                      "backend" : { "type" : "pty", "data" : {} } } }
<- { "return": { "pty" : "/dev/pty/42" } }

-> {"execute" : "chardev-change",
    "arguments" : {
      "id" : "charchannel2",
      "backend" : {
        "type" : "socket",
        "data" : {
          "addr" : {
            "type" : "unix" ,
            "data" : {
              "path" : "/tmp/charchannel2.socket"
            }
          },
          "server" : true,
          "wait" : false }}}
<- {"return": {}}
```

chardev-remove (Command)

Remove a character device backend

Arguments

id: **string** the chardev's ID, must exist and not be in use

Returns

Nothing on success

Since

1.4

Example

```
-> { "execute": "chardev-remove", "arguments": { "id" : "foo" } }  
<- { "return": {} }
```

chardev-send-break (Command)

Send a break to a character device

Arguments

id: string the chardev's ID, must exist

Returns

Nothing on success

Since

2.10

Example

```
-> { "execute": "chardev-send-break", "arguments": { "id" : "foo" } }  
<- { "return": {} }
```

VSERPORT_CHANGE (Event)

Emitted when the guest opens or closes a virtio-serial port.

Arguments

id: string device identifier of the virtio-serial port

open: boolean true if the guest has opened the virtio-serial port

Note

This event is rate-limited.

Since

2.1

Example

```
<- { "event": "VSERPORT_CHANGE",  
      "data": { "id": "channel0", "open": true },  
      "timestamp": { "seconds": 1401385907, "microseconds": 422329 } }
```

4.9.6 QMP monitor control

qmp_capabilities (Command)

Enable QMP capabilities.

Arguments:

Arguments

enable: array of QMPCapability (optional) An optional list of QMPCapability values to enable. The client must not enable any capability that is not mentioned in the QMP greeting message. If the field is not provided, it means no QMP capabilities will be enabled. (since 2.12)

Example

```
-> { "execute": "qmp_capabilities",  
      "arguments": { "enable": [ "oob" ] } }  
<- { "return": {} }
```

Notes

This command is valid exactly when first connecting: it must be issued before any other command will be accepted, and will fail once the monitor is accepting other commands. (see [qemu docs/interop/qmp-spec.txt](#))

The QMP client needs to explicitly enable QMP capabilities, otherwise all the QMP capabilities will be turned off by default.

Since

0.13

QMPCapability (Enum)

Enumeration of capabilities to be advertised during initial client connection, used for agreeing on particular QMP extension behaviors.

Values

oob QMP ability to support out-of-band requests. (Please refer to qmp-spec.txt for more information on OOB)

Since

2.12

VersionTriple (Object)

A three-part version number.

Members

major: int The major version number.

minor: int The minor version number.

micro: int The micro version number.

Since

2.4

VersionInfo (Object)

A description of QEMU's version.

Members

qemu: VersionTriple The version of QEMU. By current convention, a micro version of 50 signifies a development branch. A micro version greater than or equal to 90 signifies a release candidate for the next minor version. A micro version of less than 50 signifies a stable release.

package: string QEMU will always set this field to an empty string. Downstream versions of QEMU should set this to a non-empty string. The exact format depends on the downstream however it is highly recommended that a unique name is used.

Since

0.14

query-version (Command)

Returns the current version of QEMU.

Returns

A `VersionInfo` object describing the current version of QEMU.

Since

0.14

Example

```
-> { "execute": "query-version" }
<- {
  "return": {
    "qemu": {
      "major": 0,
      "minor": 11,
      "micro": 5
    },
    "package": ""
  }
}
```

CommandInfo (Object)

Information about a QMP command

Members

name: string The command name

Since

0.14

query-commands (Command)

Return a list of supported QMP commands by this server

Returns

A list of `CommandInfo` for all supported commands

Since

0.14

Example

```
-> { "execute": "query-commands" }
<- {
  "return": [
    {
      "name": "query-balloon"
    },
    {
      "name": "system_powerdown"
    }
  ]
}
```

Note

This example has been shortened as the real response is too long.

EventInfo (Object)

Information about a QMP event

Members

name: string The event name

Since

1.2

query-events (Command)

Return information on QMP events.

Features

deprecated This command is deprecated, because its output doesn't reflect compile-time configuration. Use 'query-qmp-schema' instead.

Returns

A list of EventInfo.

Since

1.2

Example

```
-> { "execute": "query-events" }
<- {
  "return": [
    {
      "name": "SHUTDOWN"
    },
    {
      "name": "RESET"
    }
  ]
}
```

Note

This example has been shortened as the real response is too long.

quit (Command)

This command will cause the QEMU process to exit gracefully. While every attempt is made to send the QMP response before terminating, this is not guaranteed. When using this interface, a premature EOF would not be unexpected.

Since

0.14

Example

```
-> { "execute": "quit" }
<- { "return": {} }
```

MonitorMode (Enum)

An enumeration of monitor modes.

Values

readline HMP monitor (human-oriented command line interface)

control QMP monitor (JSON-based machine interface)

Since

5.0

MonitorOptions (Object)

Options to be used for adding a new monitor.

Members

id: `string` (optional) Name of the monitor

mode: `MonitorMode` (optional) Selects the monitor mode (default: readline in the system emulator, control in qemu-storage-daemon)

pretty: `boolean` (optional) Enables pretty printing (QMP only)

chardev: `string` Name of a character device to expose the monitor on

Since

5.0

4.9.7 QMP introspection

query-qmp-schema (Command)

Command `query-qmp-schema` exposes the QMP wire ABI as an array of `SchemaInfo`. This lets QMP clients figure out what commands and events are available in this QEMU, and their parameters and results.

However, the `SchemaInfo` can't reflect all the rules and restrictions that apply to QMP. It's interface introspection (figuring out what's there), not interface specification. The specification is in the QAPI schema.

Furthermore, while we strive to keep the QMP wire format backwards-compatible across qemu versions, the introspection output is not guaranteed to have the same stability. For example, one version of qemu may list an object member as an optional non-variant, while another lists the same member only through the object's variants; or the type of a member may change from a generic string into a specific enum or from one specific type into an alternate that includes the original type alongside something else.

Returns

array of `SchemaInfo`, where each element describes an entity in the ABI: command, event, type, ...

The order of the various `SchemaInfo` is unspecified; however, all names are guaranteed to be unique (no name will be duplicated with different meta-types).

Note

the QAPI schema is also used to help define *internal* interfaces, by defining QAPI types. These are not part of the QMP wire ABI, and therefore not returned by this command.

Since

2.5

SchemaMetaType (Enum)

This is a `SchemaInfo`'s meta type, i.e. the kind of entity it describes.

Values

builtin a predefined type such as 'int' or 'bool'.

enum an enumeration type

array an array type

object an object type (struct or union)

alternate an alternate type

command a QMP command

event a QMP event

Since

2.5

SchemaInfo (Object)

Members

name: string the entity's name, inherited from `base`. The `SchemaInfo` is always referenced by this name. Commands and events have the name defined in the QAPI schema. Unlike command and event names, type names are not part of the wire ABI. Consequently, type names are meaningless strings here, although they are still guaranteed unique regardless of `meta-type`.

meta-type: SchemaMetaType the entity's meta type, inherited from `base`.

features: array of string (optional) names of features associated with the entity, in no particular order. (since 4.1 for object types, 4.2 for commands, 5.0 for the rest)

The members of `SchemaInfoBuiltin` when `meta-type` is "builtin"

The members of `SchemaInfoEnum` when `meta-type` is "enum"

The members of `SchemaInfoArray` when `meta-type` is "array"

The members of `SchemaInfoObject` when `meta-type` is "object"

The members of `SchemaInfoAlternate` when `meta-type` is "alternate"

The members of `SchemaInfoCommand` when `meta-type` is "command"

The members of `SchemaInfoEvent` when `meta-type` is "event"

Additional members depend on the value of `meta-type`.

Since

2.5

SchemaInfoBuiltin (Object)

Additional SchemaInfo members for meta-type 'builtin'.

Members

json-type: JSONType the JSON type used for this type on the wire.

Since

2.5

JSONType (Enum)

The four primitive and two structured types according to RFC 8259 section 1, plus 'int' (split off 'number'), plus the obvious top type 'value'.

Values

string Not documented

number Not documented

int Not documented

boolean Not documented

null Not documented

object Not documented

array Not documented

value Not documented

Since

2.5

SchemaInfoEnum (Object)

Additional SchemaInfo members for meta-type 'enum'.

Members

values: array of string the enumeration type's values, in no particular order.

Values of this type are JSON string on the wire.

Since

2.5

SchemaInfoArray (Object)

Additional SchemaInfo members for meta-type ‘array’.

Members

element-type: string the array type’s element type.

Values of this type are JSON array on the wire.

Since

2.5

SchemaInfoObject (Object)

Additional SchemaInfo members for meta-type ‘object’.

Members

members: array of SchemaInfoObjectMember the object type’s (non-variant) members, in no particular order.

tag: string (optional) the name of the member serving as type tag. An element of `members` with this name must exist.

variants: array of SchemaInfoObjectVariant (optional) variant members, i.e. additional members that depend on the type tag’s value. Present exactly when `tag` is present. The variants are in no particular order, and may even differ from the order of the values of the enum type of the `tag`.

Values of this type are JSON object on the wire.

Since

2.5

SchemaInfoObjectMember (Object)

An object member.

Members

name: string the member's name, as defined in the QAPI schema.

type: string the name of the member's type.

default: value (optional) default when used as command parameter. If absent, the parameter is mandatory. If present, the value must be null. The parameter is optional, and behavior when it's missing is not specified here. Future extension: if present and non-null, the parameter is optional, and defaults to this value.

features: array of string (optional) names of features associated with the member, in no particular order. (since 5.0)

Since

2.5

SchemaInfoObjectVariant (Object)

The variant members for a value of the type tag.

Members

case: string a value of the type tag.

type: string the name of the object type that provides the variant members when the type tag has value `case`.

Since

2.5

SchemaInfoAlternate (Object)

Additional SchemaInfo members for meta-type 'alternate'.

Members

members: array of SchemaInfoAlternateMember the alternate type's members, in no particular order. The members' wire encoding is distinct, see docs/devel/qapi-code-gen.txt section Alternate types.

On the wire, this can be any of the members.

Since

2.5

SchemaInfoAlternateMember (Object)

An alternate member.

Members

type: **string** the name of the member's type.

Since

2.5

SchemaInfoCommand (Object)

Additional SchemaInfo members for meta-type 'command'.

Members

arg-type: **string** the name of the object type that provides the command's parameters.

ret-type: **string** the name of the command's result type.

allow-oob: **boolean (optional)** whether the command allows out-of-band execution, defaults to false (Since: 2.12)

TODO

`success-response` (currently irrelevant, because it's QGA, not QMP)

Since

2.5

SchemaInfoEvent (Object)

Additional SchemaInfo members for meta-type 'event'.

Members

arg-type: **string** the name of the object type that provides the event's parameters.

Since

2.5

4.9.8 QEMU Object Model (QOM)

ObjectPropertyInfo (Object)

Members

name: string the name of the property

type: string the type of the property. This will typically come in one of four forms:

- 1) A primitive type such as 'u8', 'u16', 'bool', 'str', or 'double'. These types are mapped to the appropriate JSON type.
- 2) A child type in the form 'child<subtype>' where subtype is a qdev device type name. Child properties create the composition tree.
- 3) A link type in the form 'link<subtype>' where subtype is a qdev device type name. Link properties form the device model graph.

description: string (optional) if specified, the description of the property.

default-value: value (optional) the default value, if any (since 5.0)

Since

1.2

qom-list (Command)

This command will list any properties of an object given a path in the object model.

Arguments

path: string the path within the object model. See `qom-get` for a description of this parameter.

Returns

a list of `ObjectPropertyInfo` that describe the properties of the object.

Since

1.2

Example

```
-> { "execute": "qom-list",
      "arguments": { "path": "/chardevs" } }
<- { "return": [ { "name": "type", "type": "string" },
                  { "name": "parallel0", "type": "child<chardev-vc>" },
                  { "name": "serial0", "type": "child<chardev-vc>" },
                  { "name": "mon0", "type": "child<chardev-stdio>" } ] }
```

qom-get (Command)

This command will get a property from a object model path and return the value.

Arguments

path: string The path within the object model. There are two forms of supported paths—absolute and partial paths.

Absolute paths are derived from the root object and can follow `child<>` or `link<>` properties. Since they can follow `link<>` properties, they can be arbitrarily long. Absolute paths look like absolute filenames and are prefixed with a leading slash.

Partial paths look like relative filenames. They do not begin with a prefix. The matching rules for partial paths are subtle but designed to make specifying objects easy. At each level of the composition tree, the partial path is matched as an absolute path. The first match is not returned. At least two matches are searched for. A successful result is only returned if only one match is found. If more than one match is found, a flag is return to indicate that the match was ambiguous.

property: string The property name to read

Returns

The property value. The type depends on the property type. `child<>` and `link<>` properties are returned as `#str` pathnames. All integer property types (`u8`, `u16`, etc) are returned as `#int`.

Since

1.2

Example

```
1. Use absolute path
-> { "execute": "qom-get",
    "arguments": { "path": "/machine/unattached/device[0]",
                  "property": "hotplugged" } }
<- { "return": false }

2. Use partial path
-> { "execute": "qom-get",
    "arguments": { "path": "unattached/sysbus",
                  "property": "type" } }
<- { "return": "System" }
```

qom-set (Command)

This command will set a property from a object model path.

Arguments

path: string see `qom-get` for a description of this parameter

property: string the property name to set

value: value a value whose type is appropriate for the property type. See `qom-get` for a description of type mapping.

Since

1.2

Example

```

-> { "execute": "qom-set",
      "arguments": { "path": "/machine",
                     "property": "graphics",
                     "value": false } }
<- { "return": {} }

```

ObjectTypeInfo (Object)

This structure describes a search result from `qom-list-types`

Members

name: string the type name found in the search

abstract: boolean (optional) the type is abstract and can't be directly instantiated. Omitted if false. (since 2.10)

parent: string (optional) Name of parent type, if any (since 2.10)

Since

1.1

qom-list-types (Command)

This command will return a list of types given search parameters

Arguments

implements: string (optional) if specified, only return types that implement this type name

abstract: boolean (optional) if true, include abstract types in the results

Returns

a list of `ObjectTypeInfo` or an empty list if no results are found

Since

1.1

`qom-list-properties` (Command)

List properties associated with a QOM object.

Arguments

typename: `string` the type name of an object

Note

objects can create properties at runtime, for example to describe links between different devices and/or objects. These properties are not included in the output of this command.

Returns

a list of `ObjectPropertyInfo` describing object properties

Since

2.12

`object-add` (Command)

Create a QOM object.

Arguments

qom-type: `string` the class name for the object to be created

id: `string` the name of the new object

props: `value` (optional) a dictionary of properties to be passed to the backend. Deprecated since 5.0, specify the properties on the top level instead. It is an error to specify the same option both on the top level and in `props`.

Additional arguments depend on `qom-type` and are passed to the backend unchanged.

Returns

Nothing on success Error if `qom-type` is not a valid class name

Since

2.0

Example

```
-> { "execute": "object-add",  
      "arguments": { "qom-type": "rng-random", "id": "rng1",  
                     "filename": "/dev/hwrng" } }  
<- { "return": {} }
```

object-del (Command)

Remove a QOM object.

Arguments

id: `string` the name of the QOM object to remove

Returns

Nothing on success Error if `id` is not a valid id for a QOM object

Since

2.0

Example

```
-> { "execute": "object-del", "arguments": { "id": "rng1" } }  
<- { "return": {} }
```

4.9.9 Transactions

Abort (Object)

This action can be used to test transaction failure.

Since

1.6

ActionCompletionMode (Enum)

An enumeration of Transactional completion modes.

Values

individual Do not attempt to cancel any other Actions if any Actions fail after the Transaction request succeeds. All Actions that can complete successfully will do so without waiting on others. This is the default.

grouped If any Action fails after the Transaction succeeds, cancel all Actions. Actions do not complete until all Actions are ready to complete. May be rejected by Actions that do not support this completion mode.

Since

2.5

TransactionAction (Object)

A discriminated record of operations that can be performed with `transaction`. Action type can be:

- `abort`: since 1.6
- `block-dirty-bitmap-add`: since 2.5
- `block-dirty-bitmap-remove`: since 4.2
- `block-dirty-bitmap-clear`: since 2.5
- `block-dirty-bitmap-enable`: since 4.0
- `block-dirty-bitmap-disable`: since 4.0
- `block-dirty-bitmap-merge`: since 4.0
- `blockdev-backup`: since 2.3
- `blockdev-snapshot`: since 2.5
- `blockdev-snapshot-internal-sync`: since 1.7
- `blockdev-snapshot-sync`: since 1.1
- `drive-backup`: since 1.6

Members

type One of `abort`, `block-dirty-bitmap-add`, `block-dirty-bitmap-remove`, `block-dirty-bitmap-clear`, `block-dirty-bitmap-enable`, `block-dirty-bitmap-disable`, `block-dirty-bitmap-merge`, `blockdev-backup`, `blockdev-snapshot`, `blockdev-snapshot-internal-sync`, `blockdev-snapshot-sync`, `drive-backup`

data: Abort when type is "abort"

data: BlockDirtyBitmapAdd when type is "block-dirty-bitmap-add"

data: BlockDirtyBitmap when type is "block-dirty-bitmap-remove"

data: BlockDirtyBitmap when type is "block-dirty-bitmap-clear"

data: BlockDirtyBitmap when type is "block-dirty-bitmap-enable"

data: BlockDirtyBitmap when type is "block-dirty-bitmap-disable"

data: BlockDirtyBitmapMerge when type is "block-dirty-bitmap-merge"

data: BlockdevBackup when type is "blockdev-backup"

data: BlockdevSnapshot when type is "blockdev-snapshot"

data: BlockdevSnapshotInternal when type is "blockdev-snapshot-internal-sync"

data: BlockdevSnapshotSync when type is "blockdev-snapshot-sync"

data: DriveBackup when type is "drive-backup"

Since

1.1

TransactionProperties (Object)

Optional arguments to modify the behavior of a Transaction.

Members

completion-mode: ActionCompletionMode (optional) Controls how jobs launched asynchronously by Actions will complete or fail as a group. See `ActionCompletionMode` for details.

Since

2.5

transaction (Command)

Executes a number of transactionable QMP commands atomically. If any operation fails, then the entire set of actions will be abandoned and the appropriate error returned.

For external snapshots, the dictionary contains the device, the file to use for the new snapshot, and the format. The default format, if not specified, is qcow2.

Each new snapshot defaults to being created by QEMU (wiping any contents if the file already exists), but it is also possible to reuse an externally-created file. In the latter case, you should ensure that the new image file has the same contents as the current one; QEMU cannot perform any meaningful check. Typically this is achieved by using the current image file as the backing file for the new image.

On failure, the original disks pre-snapshot attempt will be used.

For internal snapshots, the dictionary contains the device and the snapshot's name. If an internal snapshot matching name already exists, the request will be rejected. Only some image formats support it, for example, qcow2, rbd, and sheepdog.

On failure, qemu will try delete the newly created internal snapshot in the transaction. When an I/O error occurs during deletion, the user needs to fix it later with `qemu-img` or other command.

Arguments

actions: array of TransactionAction List of TransactionAction; information needed for the respective operations.

properties: TransactionProperties (optional) structure of additional options to control the execution of the transaction. See TransactionProperties for additional detail.

Returns

nothing on success

Errors depend on the operations of the transaction

Note

The transaction aborts on the first failure. Therefore, there will be information on only one failed operation returned in an error condition, and subsequent actions will not have been attempted.

Since

1.1

Example

```
-> { "execute": "transaction",  
    "arguments": { "actions": [  
        { "type": "blockdev-snapshot-sync", "data" : { "device": "ide-hd0",  
                                                    "snapshot-file": "/some/place/my-image",  
                                                    "format": "qcow2" } },  
        { "type": "blockdev-snapshot-sync", "data" : { "node-name": "myfile",  
                                                    "snapshot-file": "/some/place/my-image2",  
                                                    "snapshot-node-name": "node3432",  
                                                    "mode": "existing",  
                                                    "format": "qcow2" } },  
        { "type": "blockdev-snapshot-sync", "data" : { "device": "ide-hd1",  
                                                    "snapshot-file": "/some/place/my-image2",  
                                                    "mode": "existing",  
                                                    "format": "qcow2" } },  
        { "type": "blockdev-snapshot-internal-sync", "data" : {  
                                                    "device": "ide-hd2",  
                                                    "name": "snapshot0" } } ] } }  
<- { "return": {} }
```

4.10 Vhost-user Protocol

Copyright 2014 Virtual Open Systems Sarl.

Copyright 2019 Intel Corporation

Licence This work is licensed under the terms of the GNU GPL, version 2 or later. See the COPYING file in the top-level directory.

Table of Contents

- *Vhost-user Protocol*
 - *Introduction*
 - *Message Specification*
 - * *Header*
 - * *Payload*
 - *A single 64-bit integer*
 - *A vring state description*
 - *A vring address description*
 - *Memory regions description*
 - *Single memory region description*
 - *Log description*
 - *An IOTLB message*
 - *Virtio device config space*
 - *Vring area description*
 - *Inflight description*
 - * *C structure*
 - *Communication*
 - * *Starting and stopping rings*
 - * *Multiple queue support*
 - * *Migration*
 - * *Memory access*
 - * *IOMMU support*
 - * *Slave communication*
 - * *Inflight I/O tracking*
 - * *In-band notifications*
 - * *Protocol features*
 - * *Master message types*
 - * *Slave message types*
 - * *VHOST_USER_PROTOCOL_F_REPLY_ACK*
 - *Backend program conventions*
 - * *vhost-user-input*
 - * *vhost-user-gpu*

* *vhost-user-blk*

4.10.1 Introduction

This protocol is aiming to complement the `ioctl` interface used to control the vhost implementation in the Linux kernel. It implements the control plane needed to establish virtqueue sharing with a user space process on the same host. It uses communication over a Unix domain socket to share file descriptors in the ancillary data of the message.

The protocol defines 2 sides of the communication, *master* and *slave*. *Master* is the application that shares its virtqueues, in our case QEMU. *Slave* is the consumer of the virtqueues.

In the current implementation QEMU is the *master*, and the *slave* is the external process consuming the virtio queues, for example a software Ethernet switch running in user space, such as Snabbswitch, or a block device backend processing read & write to a virtual disk. In order to facilitate interoperability between various backend implementations, it is recommended to follow the [Backend program conventions](#).

Master and *slave* can be either a client (i.e. connecting) or server (listening) in the socket communication.

4.10.2 Message Specification

Note: All numbers are in the machine native byte order.

A vhost-user message consists of 3 header fields and a payload.

request	flags	size	payload
---------	-------	------	---------

Header

request 32-bit type of the request

flags 32-bit bit field

- Lower 2 bits are the version (currently 0x01)
- Bit 2 is the reply flag - needs to be sent on each reply from the slave
- Bit 3 is the need_reply flag - see [REPLY_ACK](#) for details.

size 32-bit size of the payload

Payload

Depending on the request type, **payload** can be:

A single 64-bit integer

u64

u64 a 64-bit unsigned integer

A vring state description

index	num
-------	-----

index a 32-bit index

num a 32-bit number

A vring address description

index	flags	size	descriptor	used	available	log
-------	-------	------	------------	------	-----------	-----

index a 32-bit vring index

flags a 32-bit vring flags

descriptor a 64-bit ring address of the vring descriptor table

used a 64-bit ring address of the vring used ring

available a 64-bit ring address of the vring available ring

log a 64-bit guest address for logging

Note that a ring address is an IOVA if `VIRTIO_F_IOMMU_PLATFORM` has been negotiated. Otherwise it is a user address.

Memory regions description

num regions	padding	region0	...	region7
-------------	---------	---------	-----	---------

num regions a 32-bit number of regions

padding 32-bit

A region is:

guest address	size	user address	mmap offset
---------------	------	--------------	-------------

guest address a 64-bit guest address of the region

size a 64-bit size

user address a 64-bit user address

mmap offset 64-bit offset where region starts in the mapped memory

Single memory region description

padding	guest address	size	user address	mmap offset
---------	---------------	------	--------------	-------------

padding 64-bit

guest address a 64-bit guest address of the region

size a 64-bit size

user address a 64-bit user address

mmap offset 64-bit offset where region starts in the mapped memory

Log description

log size	log offset
----------	------------

log size size of area used for logging

log offset offset from start of supplied file descriptor where logging starts (i.e. where guest address 0 would be logged)

An IOTLB message

iova	size	user address	permissions flags	type
------	------	--------------	-------------------	------

iova a 64-bit I/O virtual address programmed by the guest

size a 64-bit size

user address a 64-bit user address

permissions flags an 8-bit value: - 0: No access - 1: Read access - 2: Write access - 3: Read/Write access

type an 8-bit IOTLB message type: - 1: IOTLB miss - 2: IOTLB update - 3: IOTLB invalidate - 4: IOTLB access fail

Virtio device config space

offset	size	flags	payload
--------	------	-------	---------

offset a 32-bit offset of virtio device's configuration space

size a 32-bit configuration space access size in bytes

flags a 32-bit value: - 0: Vhost master messages used for writeable fields - 1: Vhost master messages used for live migration

payload Size bytes array holding the contents of the virtio device's configuration space

Vring area description

u64	size	offset
-----	------	--------

u64 a 64-bit integer contains vring index and flags

size a 64-bit size of this area

offset a 64-bit offset of this area from the start of the supplied file descriptor

Inflight description

mmap size	mmap offset	num queues	queue size
-----------	-------------	------------	------------

mmap size a 64-bit size of area to track inflight I/O

mmap offset a 64-bit offset of this area from the start of the supplied file descriptor

num queues a 16-bit number of virtqueues

queue size a 16-bit size of virtqueues

C structure

In QEMU the vhost-user message is implemented with the following struct:

```
typedef struct VhostUserMsg {
    VhostUserRequest request;
    uint32_t flags;
    uint32_t size;
    union {
        uint64_t u64;
        struct vhost_vring_state state;
        struct vhost_vring_addr addr;
        VhostUserMemory memory;
        VhostUserLog log;
        struct vhost_iotlb_msg iotlb;
        VhostUserConfig config;
        VhostUserVringArea area;
        VhostUserInflight inflight;
    };
} QEMU_PACKED VhostUserMsg;
```

4.10.3 Communication

The protocol for vhost-user is based on the existing implementation of vhost for the Linux Kernel. Most messages that can be sent via the Unix domain socket implementing vhost-user have an equivalent ioctl to the kernel implementation.

The communication consists of *master* sending message requests and *slave* sending message replies. Most of the requests don't require replies. Here is a list of the ones that do:

- VHOST_USER_GET_FEATURES
- VHOST_USER_GET_PROTOCOL_FEATURES
- VHOST_USER_GET_VRING_BASE
- VHOST_USER_SET_LOG_BASE (if VHOST_USER_PROTOCOL_F_LOG_SHMFD)
- VHOST_USER_GET_INFLIGHT_FD (if VHOST_USER_PROTOCOL_F_INFLIGHT_SHMFD)

See also:

REPLY_ACK The section on REPLY_ACK protocol extension.

There are several messages that the master sends with file descriptors passed in the ancillary data:

- VHOST_USER_SET_MEM_TABLE
- VHOST_USER_SET_LOG_BASE (if VHOST_USER_PROTOCOL_F_LOG_SHMFD)

- `VHOST_USER_SET_LOG_FD`
- `VHOST_USER_SET_VRING_KICK`
- `VHOST_USER_SET_VRING_CALL`
- `VHOST_USER_SET_VRING_ERR`
- `VHOST_USER_SET_SLAVE_REQ_FD`
- `VHOST_USER_SET_INFLIGHT_FD` (if `VHOST_USER_PROTOCOL_F_INFLIGHT_SHMFD`)

If *master* is unable to send the full message or receives a wrong reply it will close the connection. An optional reconnection mechanism can be implemented.

If *slave* detects some error such as incompatible features, it may also close the connection. This should only happen in exceptional circumstances.

Any protocol extensions are gated by protocol feature bits, which allows full backwards compatibility on both master and slave. As older slaves don't support negotiating protocol features, a feature bit was dedicated for this purpose:

```
#define VHOST_USER_F_PROTOCOL_FEATURES 30
```

Starting and stopping rings

Client must only process each ring when it is started.

Client must only pass data between the ring and the backend, when the ring is enabled.

If ring is started but disabled, client must process the ring without talking to the backend.

For example, for a networking device, in the disabled state client must not supply any new RX packets, but must process and discard any TX packets.

If `VHOST_USER_F_PROTOCOL_FEATURES` has not been negotiated, the ring is initialized in an enabled state.

If `VHOST_USER_F_PROTOCOL_FEATURES` has been negotiated, the ring is initialized in a disabled state. Client must not pass data to/from the backend until ring is enabled by `VHOST_USER_SET_VRING_ENABLE` with parameter 1, or after it has been disabled by `VHOST_USER_SET_VRING_ENABLE` with parameter 0.

Each ring is initialized in a stopped state, client must not process it until ring is started, or after it has been stopped.

Client must start ring upon receiving a kick (that is, detecting that file descriptor is readable) on the descriptor specified by `VHOST_USER_SET_VRING_KICK` or receiving the in-band message `VHOST_USER_VRING_KICK` if negotiated, and stop ring upon receiving `VHOST_USER_GET_VRING_BASE`.

While processing the rings (whether they are enabled or not), client must support changing some configuration aspects on the fly.

Multiple queue support

Many devices have a fixed number of virtqueues. In this case the master already knows the number of available virtqueues without communicating with the slave.

Some devices do not have a fixed number of virtqueues. Instead the maximum number of virtqueues is chosen by the slave. The number can depend on host resource availability or slave implementation details. Such devices are called multiple queue devices.

Multiple queue support allows the slave to advertise the maximum number of queues. This is treated as a protocol extension, hence the slave has to implement protocol features first. The multiple queues feature is supported only when the protocol feature `VHOST_USER_PROTOCOL_F_MQ` (bit 0) is set.

The max number of queues the slave supports can be queried with message `VHOST_USER_GET_QUEUE_NUM`. Master should stop when the number of requested queues is bigger than that.

As all queues share one connection, the master uses a unique index for each queue in the sent message to identify a specified queue.

The master enables queues by sending message `VHOST_USER_SET_VRING_ENABLE`. `vhost-user-net` has historically automatically enabled the first queue pair.

Slaves should always implement the `VHOST_USER_PROTOCOL_F_MQ` protocol feature, even for devices with a fixed number of virtqueues, since it is simple to implement and offers a degree of introspection.

Masters must not rely on the `VHOST_USER_PROTOCOL_F_MQ` protocol feature for devices with a fixed number of virtqueues. Only true multiqueue devices require this protocol feature.

Migration

During live migration, the master may need to track the modifications the slave makes to the memory mapped regions. The client should mark the dirty pages in a log. Once it complies to this logging, it may declare the `VHOST_F_LOG_ALL` vhost feature.

To start/stop logging of data/used ring writes, server may send messages `VHOST_USER_SET_FEATURES` with `VHOST_F_LOG_ALL` and `VHOST_USER_SET_VRING_ADDR` with `VHOST_VRING_F_LOG` in ring's flags set to 1/0, respectively.

All the modifications to memory pointed by vring “descriptor” should be marked. Modifications to “used” vring should be marked if `VHOST_VRING_F_LOG` is part of ring's flags.

Dirty pages are of size:

```
#define VHOST_LOG_PAGE 0x1000
```

The log memory fd is provided in the ancillary data of `VHOST_USER_SET_LOG_BASE` message when the slave has `VHOST_USER_PROTOCOL_F_LOG_SHMFD` protocol feature.

The size of the log is supplied as part of `VhostUserMsg` which should be large enough to cover all known guest addresses. Log starts at the supplied offset in the supplied file descriptor. The log covers from address 0 to the maximum of guest regions. In pseudo-code, to mark page at `addr` as dirty:

```
page = addr / VHOST_LOG_PAGE
log[page / 8] |= 1 << page % 8
```

Where `addr` is the guest physical address.

Use atomic operations, as the log may be concurrently manipulated.

Note that when logging modifications to the used ring (when `VHOST_VRING_F_LOG` is set for this ring), `log_guest_addr` should be used to calculate the log offset: the write to first byte of the used ring is logged at this offset from log start. Also note that this value might be outside the legal guest physical address range (i.e. does not have to be covered by the `VhostUserMemory` table), but the bit offset of the last byte of the ring must fall within the size supplied by `VhostUserLog`.

`VHOST_USER_SET_LOG_FD` is an optional message with an eventfd in ancillary data, it may be used to inform the master that the log has been modified.

Once the source has finished migration, rings will be stopped by the source. No further update must be done before rings are restarted.

In postcopy migration the slave is started before all the memory has been received from the source host, and care must be taken to avoid accessing pages that have yet to be received. The slave opens a ‘userfault’-fd and registers the

memory with it; this fd is then passed back over to the master. The master services requests on the userfaultfd for pages that are accessed and when the page is available it performs `WAKE ioctl`'s on the userfaultfd to wake the stalled slave. The client indicates support for this via the `VHOST_USER_PROTOCOL_F_PAGEFAULT` feature.

Memory access

The master sends a list of vhost memory regions to the slave using the `VHOST_USER_SET_MEM_TABLE` message. Each region has two base addresses: a guest address and a user address.

Messages contain guest addresses and/or user addresses to reference locations within the shared memory. The mapping of these addresses works as follows.

User addresses map to the vhost memory region containing that user address.

When the `VIRTIO_F_IOMMU_PLATFORM` feature has not been negotiated:

- Guest addresses map to the vhost memory region containing that guest address.

When the `VIRTIO_F_IOMMU_PLATFORM` feature has been negotiated:

- Guest addresses are also called I/O virtual addresses (IOVAs). They are translated to user addresses via the IOTLB.
- The vhost memory region guest address is not used.

IOMMU support

When the `VIRTIO_F_IOMMU_PLATFORM` feature has been negotiated, the master sends IOTLB entries update & invalidation by sending `VHOST_USER_IOTLB_MSG` requests to the slave with a `struct vhost_iotlb_msg` as payload. For update events, the `iotlb` payload has to be filled with the update message type (2), the I/O virtual address, the size, the user virtual address, and the permissions flags. Addresses and size must be within vhost memory regions set via the `VHOST_USER_SET_MEM_TABLE` request. For invalidation events, the `iotlb` payload has to be filled with the invalidation message type (3), the I/O virtual address and the size. On success, the slave is expected to reply with a zero payload, non-zero otherwise.

The slave relies on the slave communication channel (see *Slave communication* section below) to send IOTLB miss and access failure events, by sending `VHOST_USER_SLAVE_IOTLB_MSG` requests to the master with a `struct vhost_iotlb_msg` as payload. For miss events, the `iotlb` payload has to be filled with the miss message type (1), the I/O virtual address and the permissions flags. For access failure event, the `iotlb` payload has to be filled with the access failure message type (4), the I/O virtual address and the permissions flags. For synchronization purpose, the slave may rely on the reply-ack feature, so the master may send a reply when operation is completed if the reply-ack feature is negotiated and slaves requests a reply. For miss events, completed operation means either master sent an update message containing the IOTLB entry containing requested address and permission, or master sent nothing if the IOTLB miss message is invalid (invalid IOVA or permission).

The master isn't expected to take the initiative to send IOTLB update messages, as the slave sends IOTLB miss messages for the guest virtual memory areas it needs to access.

Slave communication

An optional communication channel is provided if the slave declares `VHOST_USER_PROTOCOL_F_SLAVE_REQ` protocol feature, to allow the slave to make requests to the master.

The fd is provided via `VHOST_USER_SET_SLAVE_REQ_FD` ancillary data.

A slave may then send `VHOST_USER_SLAVE_*` messages to the master using this fd communication channel.

If `VHOST_USER_PROTOCOL_F_SLAVE_SEND_FD` protocol feature is negotiated, slave can send file descriptors (at most 8 descriptors in each message) to master via ancillary data using this fd communication channel.

Inflight I/O tracking

To support reconnecting after restart or crash, slave may need to resubmit inflight I/Os. If virtqueue is processed in order, we can easily achieve that by getting the inflight descriptors from descriptor table (split virtqueue) or descriptor ring (packed virtqueue). However, it can't work when we process descriptors out-of-order because some entries which store the information of inflight descriptors in available ring (split virtqueue) or descriptor ring (packed virtqueue) might be overridden by new entries. To solve this problem, slave need to allocate an extra buffer to store this information of inflight descriptors and share it with master for persistent. `VHOST_USER_GET_INFLIGHT_FD` and `VHOST_USER_SET_INFLIGHT_FD` are used to transfer this buffer between master and slave. And the format of this buffer is described below:

queue0 region	queue1 region	...	queueN region
---------------	---------------	-----	---------------

N is the number of available virtqueues. Slave could get it from `num queues` field of `VhostUserInflight`.

For split virtqueue, queue region can be implemented as:

```
typedef struct DescStateSplit {
    /* Indicate whether this descriptor is inflight or not.
     * Only available for head-descriptor. */
    uint8_t inflight;

    /* Padding */
    uint8_t padding[5];

    /* Maintain a list for the last batch of used descriptors.
     * Only available when batching is used for submitting */
    uint16_t next;

    /* Used to preserve the order of fetching available descriptors.
     * Only available for head-descriptor. */
    uint64_t counter;
} DescStateSplit;

typedef struct QueueRegionSplit {
    /* The feature flags of this region. Now it's initialized to 0. */
    uint64_t features;

    /* The version of this region. It's 1 currently.
     * Zero value indicates an uninitialized buffer */
    uint16_t version;

    /* The size of DescStateSplit array. It's equal to the virtqueue
     * size. Slave could get it from queue size field of VhostUserInflight. */
    uint16_t desc_num;

    /* The head of list that track the last batch of used descriptors. */
    uint16_t last_batch_head;

    /* Store the idx value of used ring */
    uint16_t used_idx;

    /* Used to track the state of each descriptor in descriptor table */

```

(continues on next page)

(continued from previous page)

```
DescStateSplit desc[];
} QueueRegionSplit;
```

To track inflight I/O, the queue region should be processed as follows:

When receiving available buffers from the driver:

1. Get the next available head-descriptor index from available ring, `i`
2. Set `desc[i].counter` to the value of global counter
3. Increase global counter by 1
4. Set `desc[i].inflight` to 1

When supplying used buffers to the driver:

1. Get corresponding used head-descriptor index, `i`
2. Set `desc[i].next` to `last_batch_head`
3. Set `last_batch_head` to `i`
4. Steps 1,2,3 may be performed repeatedly if batching is possible
5. Increase the `idx` value of used ring by the size of the batch
6. Set the `inflight` field of each `DescStateSplit` entry in the batch to 0
7. Set `used_idx` to the `idx` value of used ring

When reconnecting:

1. If the value of `used_idx` does not match the `idx` value of used ring (means the `inflight` field of `DescStateSplit` entries in last batch may be incorrect),
 - a. Subtract the value of `used_idx` from the `idx` value of used ring to get last batch size of `DescStateSplit` entries
 - b. Set the `inflight` field of each `DescStateSplit` entry to 0 in last batch list which starts from `last_batch_head`
 - c. Set `used_idx` to the `idx` value of used ring
2. Resubmit inflight `DescStateSplit` entries in order of their counter value

For packed virtqueue, queue region can be implemented as:

```
typedef struct DescStatePacked {
    /* Indicate whether this descriptor is inflight or not.
     * Only available for head-descriptor. */
    uint8_t inflight;

    /* Padding */
    uint8_t padding;

    /* Link to the next free entry */
    uint16_t next;

    /* Link to the last entry of descriptor list.
     * Only available for head-descriptor. */
    uint16_t last;

    /* The length of descriptor list.

```

(continues on next page)

(continued from previous page)

```

    * Only available for head-descriptor. */
    uint16_t num;

    /* Used to preserve the order of fetching available descriptors.
     * Only available for head-descriptor. */
    uint64_t counter;

    /* The buffer id */
    uint16_t id;

    /* The descriptor flags */
    uint16_t flags;

    /* The buffer length */
    uint32_t len;

    /* The buffer address */
    uint64_t addr;
} DescStatePacked;

typedef struct QueueRegionPacked {
    /* The feature flags of this region. Now it's initialized to 0. */
    uint64_t features;

    /* The version of this region. It's 1 currently.
     * Zero value indicates an uninitialized buffer */
    uint16_t version;

    /* The size of DescStatePacked array. It's equal to the virtqueue
     * size. Slave could get it from queue size field of VhostUserInflight. */
    uint16_t desc_num;

    /* The head of free DescStatePacked entry list */
    uint16_t free_head;

    /* The old head of free DescStatePacked entry list */
    uint16_t old_free_head;

    /* The used index of descriptor ring */
    uint16_t used_idx;

    /* The old used index of descriptor ring */
    uint16_t old_used_idx;

    /* Device ring wrap counter */
    uint8_t used_wrap_counter;

    /* The old device ring wrap counter */
    uint8_t old_used_wrap_counter;

    /* Padding */
    uint8_t padding[7];

    /* Used to track the state of each descriptor fetched from descriptor ring */
    DescStatePacked desc[];
} QueueRegionPacked;

```

To track inflight I/O, the queue region should be processed as follows:

When receiving available buffers from the driver:

1. Get the next available descriptor entry from descriptor ring, `d`
2. If `d` is head descriptor,
 - a. Set `desc[old_free_head].num` to 0
 - b. Set `desc[old_free_head].counter` to the value of global counter
 - c. Increase global counter by 1
 - d. Set `desc[old_free_head].inflight` to 1
3. If `d` is last descriptor, set `desc[old_free_head].last` to `free_head`
4. Increase `desc[old_free_head].num` by 1
5. Set `desc[free_head].addr`, `desc[free_head].len`, `desc[free_head].flags`, `desc[free_head].id` to `d.addr`, `d.len`, `d.flags`, `d.id`
6. Set `free_head` to `desc[free_head].next`
7. If `d` is last descriptor, set `old_free_head` to `free_head`

When supplying used buffers to the driver:

1. Get corresponding used head-descriptor entry from descriptor ring, `d`
2. Get corresponding `DescStatePacked` entry, `e`
3. Set `desc[e.last].next` to `free_head`
4. Set `free_head` to the index of `e`
5. Steps 1,2,3,4 may be performed repeatedly if batching is possible
6. Increase `used_idx` by the size of the batch and update `used_wrap_counter` if needed
7. Update `d.flags`
8. Set the `inflight` field of each head `DescStatePacked` entry in the batch to 0
9. Set `old_free_head`, `old_used_idx`, `old_used_wrap_counter` to `free_head`, `used_idx`, `used_wrap_counter`

When reconnecting:

1. If `used_idx` does not match `old_used_idx` (means the `inflight` field of `DescStatePacked` entries in last batch may be incorrect),
 - a. Get the next descriptor ring entry through `old_used_idx`, `d`
 - b. Use `old_used_wrap_counter` to calculate the available flags
 - c. If `d.flags` is not equal to the calculated flags value (means slave has submitted the buffer to guest driver before crash, so it has to commit the in-progress update), set `old_free_head`, `old_used_idx`, `old_used_wrap_counter` to `free_head`, `used_idx`, `used_wrap_counter`
2. Set `free_head`, `used_idx`, `used_wrap_counter` to `old_free_head`, `old_used_idx`, `old_used_wrap_counter` (roll back any in-progress update)
3. Set the `inflight` field of each `DescStatePacked` entry in free list to 0
4. Resubmit inflight `DescStatePacked` entries in order of their counter value

In-band notifications

In some limited situations (e.g. for simulation) it is desirable to have the kick, call and error (if used) signals done via in-band messages instead of asynchronous eventfd notifications. This can be done by negotiating the `VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS` protocol feature.

Note that due to the fact that too many messages on the sockets can cause the sending application(s) to block, it is not advised to use this feature unless absolutely necessary. It is also considered an error to negotiate this feature without also negotiating `VHOST_USER_PROTOCOL_F_SLAVE_REQ` and `VHOST_USER_PROTOCOL_F_REPLY_ACK`, the former is necessary for getting a message channel from the slave to the master, while the latter needs to be used with the in-band notification messages to block until they are processed, both to avoid blocking later and for proper processing (at least in the simulation use case.) As it has no other way of signalling this error, the slave should close the connection as a response to a `VHOST_USER_SET_PROTOCOL_FEATURES` message that sets the in-band notifications feature flag without the other two.

Protocol features

<code>#define VHOST_USER_PROTOCOL_F_MQ</code>	<code>0</code>
<code>#define VHOST_USER_PROTOCOL_F_LOG_SHMFD</code>	<code>1</code>
<code>#define VHOST_USER_PROTOCOL_F_RARP</code>	<code>2</code>
<code>#define VHOST_USER_PROTOCOL_F_REPLY_ACK</code>	<code>3</code>
<code>#define VHOST_USER_PROTOCOL_F_MTU</code>	<code>4</code>
<code>#define VHOST_USER_PROTOCOL_F_SLAVE_REQ</code>	<code>5</code>
<code>#define VHOST_USER_PROTOCOL_F_CROSS_ENDIAN</code>	<code>6</code>
<code>#define VHOST_USER_PROTOCOL_F_CRYPTO_SESSION</code>	<code>7</code>
<code>#define VHOST_USER_PROTOCOL_F_PAGEFAULT</code>	<code>8</code>
<code>#define VHOST_USER_PROTOCOL_F_CONFIG</code>	<code>9</code>
<code>#define VHOST_USER_PROTOCOL_F_SLAVE_SEND_FD</code>	<code>10</code>
<code>#define VHOST_USER_PROTOCOL_F_HOST_NOTIFIER</code>	<code>11</code>
<code>#define VHOST_USER_PROTOCOL_F_INFLIGHT_SHMFD</code>	<code>12</code>
<code>#define VHOST_USER_PROTOCOL_F_RESET_DEVICE</code>	<code>13</code>
<code>#define VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS</code>	<code>14</code>
<code>#define VHOST_USER_PROTOCOL_F_CONFIGURE_MEM_SLOTS</code>	<code>15</code>
<code>#define VHOST_USER_PROTOCOL_F_STATUS</code>	<code>16</code>

Master message types

`VHOST_USER_GET_FEATURES`

id 1

equivalent ioctl `VHOST_GET_FEATURES`

master payload N/A

slave payload `u64`

Get from the underlying vhost implementation the features bitmask. Feature bit `VHOST_USER_F_PROTOCOL_FEATURES` signals slave support for `VHOST_USER_GET_PROTOCOL_FEATURES` and `VHOST_USER_SET_PROTOCOL_FEATURES`.

`VHOST_USER_SET_FEATURES`

id 2

equivalent ioctl `VHOST_SET_FEATURES`

master payload `u64`

Enable features in the underlying vhost implementation using a bitmask. Feature bit `VHOST_USER_F_PROTOCOL_FEATURES` signals slave support for `VHOST_USER_GET_PROTOCOL_FEATURES` and `VHOST_USER_SET_PROTOCOL_FEATURES`.

VHOST_USER_GET_PROTOCOL_FEATURES

id 15

equivalent ioctl `VHOST_GET_FEATURES`

master payload N/A

slave payload `u64`

Get the protocol feature bitmask from the underlying vhost implementation. Only legal if feature bit `VHOST_USER_F_PROTOCOL_FEATURES` is present in `VHOST_USER_GET_FEATURES`.

Note: Slave that reported `VHOST_USER_F_PROTOCOL_FEATURES` must support this message even before `VHOST_USER_SET_FEATURES` was called.

VHOST_USER_SET_PROTOCOL_FEATURES

id 16

equivalent ioctl `VHOST_SET_FEATURES`

master payload `u64`

Enable protocol features in the underlying vhost implementation.

Only legal if feature bit `VHOST_USER_F_PROTOCOL_FEATURES` is present in `VHOST_USER_GET_FEATURES`.

Note: Slave that reported `VHOST_USER_F_PROTOCOL_FEATURES` must support this message even before `VHOST_USER_SET_FEATURES` was called.

VHOST_USER_SET_OWNER

id 3

equivalent ioctl `VHOST_SET_OWNER`

master payload N/A

Issued when a new connection is established. It sets the current *master* as an owner of the session. This can be used on the *slave* as a “session start” flag.

VHOST_USER_RESET_OWNER

id 4

master payload N/A

Deprecated

This is no longer used. Used to be sent to request disabling all rings, but some clients interpreted it to also discard connection state (this interpretation would lead to bugs). It is recommended that clients either ignore this message, or use it to disable all rings.

VHOST_USER_SET_MEM_TABLE

id 5**equivalent ioctl** VHOST_SET_MEM_TABLE**master payload** memory regions description**slave payload** (postcopy only) memory regions description

Sets the memory map regions on the slave so it can translate the vring addresses. In the ancillary data there is an array of file descriptors for each memory mapped region. The size and ordering of the fds matches the number and ordering of memory regions.

When VHOST_USER_POSTCOPY_LISTEN has been received, SET_MEM_TABLE replies with the bases of the memory mapped regions to the master. The slave must have mmap'd the regions but not yet accessed them and should not yet generate a userfault event.

Note: NEED_REPLY_MASK is not set in this case. QEMU will then reply back to the list of mappings with an empty VHOST_USER_SET_MEM_TABLE as an acknowledgement; only upon reception of this message may the guest start accessing the memory and generating faults.

VHOST_USER_SET_LOG_BASE**id 6****equivalent ioctl** VHOST_SET_LOG_BASE**master payload** u64**slave payload** N/A

Sets logging shared memory space.

When slave has VHOST_USER_PROTOCOL_F_LOG_SHMFD protocol feature, the log memory fd is provided in the ancillary data of VHOST_USER_SET_LOG_BASE message, the size and offset of shared memory area provided in the message.

VHOST_USER_SET_LOG_FD**id 7****equivalent ioctl** VHOST_SET_LOG_FD**master payload** N/A

Sets the logging file descriptor, which is passed as ancillary data.

VHOST_USER_SET_VRING_NUM**id 8****equivalent ioctl** VHOST_SET_VRING_NUM**master payload** vring state description

Set the size of the queue.

VHOST_USER_SET_VRING_ADDR**id 9****equivalent ioctl** VHOST_SET_VRING_ADDR**master payload** vring address description**slave payload** N/A

Sets the addresses of the different aspects of the vring.

VHOST_USER_SET_VRING_BASE

id 10

equivalent ioctl VHOST_SET_VRING_BASE

master payload vring state description

Sets the base offset in the available vring.

VHOST_USER_GET_VRING_BASE

id 11

equivalent ioctl VHOST_USER_GET_VRING_BASE

master payload vring state description

slave payload vring state description

Get the available vring base offset.

VHOST_USER_SET_VRING_KICK

id 12

equivalent ioctl VHOST_SET_VRING_KICK

master payload `u64`

Set the event file descriptor for adding buffers to the vring. It is passed in the ancillary data.

Bits (0-7) of the payload contain the vring index. Bit 8 is the invalid FD flag. This flag is set when there is no file descriptor in the ancillary data. This signals that polling should be used instead of waiting for the kick. Note that if the protocol feature `VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS` has been negotiated this message isn't necessary as the ring is also started on the `VHOST_USER_VRING_KICK` message, it may however still be used to set an event file descriptor (which will be preferred over the message) or to enable polling.

VHOST_USER_SET_VRING_CALL

id 13

equivalent ioctl VHOST_SET_VRING_CALL

master payload `u64`

Set the event file descriptor to signal when buffers are used. It is passed in the ancillary data.

Bits (0-7) of the payload contain the vring index. Bit 8 is the invalid FD flag. This flag is set when there is no file descriptor in the ancillary data. This signals that polling will be used instead of waiting for the call. Note that if the protocol features `VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS` and `VHOST_USER_PROTOCOL_F_SLAVE_REQ` have been negotiated this message isn't necessary as the `VHOST_USER_SLAVE_VRING_CALL` message can be used, it may however still be used to set an event file descriptor or to enable polling.

VHOST_USER_SET_VRING_ERR

id 14

equivalent ioctl VHOST_SET_VRING_ERR

master payload `u64`

Set the event file descriptor to signal when error occurs. It is passed in the ancillary data.

Bits (0-7) of the payload contain the vring index. Bit 8 is the invalid FD flag. This flag is set when there is no file descriptor in the ancillary data. Note that if the protocol features `VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS` and `VHOST_USER_PROTOCOL_F_SLAVE_REQ` have been negotiated this message isn't necessary as the `VHOST_USER_SLAVE_VRING_ERR` message can be used, it may however still be used to set an event file descriptor (which will be preferred over the message).

VHOST_USER_GET_QUEUE_NUM

id 17

equivalent ioctl N/A

master payload N/A

slave payload u64

Query how many queues the backend supports.

This request should be sent only when `VHOST_USER_PROTOCOL_F_MQ` is set in queried protocol features by `VHOST_USER_GET_PROTOCOL_FEATURES`.

VHOST_USER_SET_VRING_ENABLE

id 18

equivalent ioctl N/A

master payload vring state description

Signal slave to enable or disable corresponding vring.

This request should be sent only when `VHOST_USER_F_PROTOCOL_FEATURES` has been negotiated.

VHOST_USER_SEND_RARP

id 19

equivalent ioctl N/A

master payload u64

Ask vhost user backend to broadcast a fake RARP to notify the migration is terminated for guest that does not support `GUEST_ANNOUNCE`.

Only legal if feature bit `VHOST_USER_F_PROTOCOL_FEATURES` is present in `VHOST_USER_GET_FEATURES` and protocol feature bit `VHOST_USER_PROTOCOL_F_RARP` is present in `VHOST_USER_GET_PROTOCOL_FEATURES`. The first 6 bytes of the payload contain the mac address of the guest to allow the vhost user backend to construct and broadcast the fake RARP.

VHOST_USER_NET_SET_MTU

id 20

equivalent ioctl N/A

master payload u64

Set host MTU value exposed to the guest.

This request should be sent only when `VIRTIO_NET_F_MTU` feature has been successfully negotiated, `VHOST_USER_F_PROTOCOL_FEATURES` is present in `VHOST_USER_GET_FEATURES` and protocol feature bit `VHOST_USER_PROTOCOL_F_NET_MTU` is present in `VHOST_USER_GET_PROTOCOL_FEATURES`.

If `VHOST_USER_PROTOCOL_F_REPLY_ACK` is negotiated, slave must respond with zero in case the specified MTU is valid, or non-zero otherwise.

VHOST_USER_SET_SLAVE_REQ_FD

id 21

equivalent ioctl N/A

master payload N/A

Set the socket file descriptor for slave initiated requests. It is passed in the ancillary data.

This request should be sent only when `VHOST_USER_F_PROTOCOL_FEATURES` has been negotiated, and protocol feature bit `VHOST_USER_PROTOCOL_F_SLAVE_REQ` bit is present in `VHOST_USER_GET_PROTOCOL_FEATURES`. If `VHOST_USER_PROTOCOL_F_REPLY_ACK` is negotiated, slave must respond with zero for success, non-zero otherwise.

VHOST_USER_IOTLB_MSG

id 22

equivalent ioctl N/A (equivalent to `VHOST_IOTLB_MSG` message type)

master payload `struct vhost_iotlb_msg`

slave payload `u64`

Send IOTLB messages with `struct vhost_iotlb_msg` as payload.

Master sends such requests to update and invalidate entries in the device IOTLB. The slave has to acknowledge the request with sending zero as `u64` payload for success, non-zero otherwise.

This request should be send only when `VIRTIO_F_IOMMU_PLATFORM` feature has been successfully negotiated.

VHOST_USER_SET_VRING_ENDIAN

id 23

equivalent ioctl `VHOST_SET_VRING_ENDIAN`

master payload vring state description

Set the endianness of a VQ for legacy devices. Little-endian is indicated with `state.num` set to 0 and big-endian is indicated with `state.num` set to 1. Other values are invalid.

This request should be sent only when `VHOST_USER_PROTOCOL_F_CROSS_ENDIAN` has been negotiated. Backends that negotiated this feature should handle both endiannesses and expect this message once (per VQ) during device configuration (ie. before the master starts the VQ).

VHOST_USER_GET_CONFIG

id 24

equivalent ioctl N/A

master payload virtio device config space

slave payload virtio device config space

When `VHOST_USER_PROTOCOL_F_CONFIG` is negotiated, this message is submitted by the vhost-user master to fetch the contents of the virtio device configuration space, vhost-user slave's payload size **MUST** match master's request, vhost-user slave uses zero length of payload to indicate an error to vhost-user master. The vhost-user master may cache the contents to avoid repeated `VHOST_USER_GET_CONFIG` calls.

VHOST_USER_SET_CONFIG

id 25**equivalent ioctl** N/A**master payload** virtio device config space**slave payload** N/A

When `VHOST_USER_PROTOCOL_F_CONFIG` is negotiated, this message is submitted by the vhost-user master when the Guest changes the virtio device configuration space and also can be used for live migration on the destination host. The vhost-user slave must check the flags field, and slaves **MUST NOT** accept `SET_CONFIG` for read-only configuration space fields unless the live migration bit is set.

VHOST_USER_CREATE_CRYPTO_SESSION**id** 26**equivalent ioctl** N/A**master payload** crypto session description**slave payload** crypto session description

Create a session for crypto operation. The server side must return the session id, 0 or positive for success, negative for failure. This request should be sent only when `VHOST_USER_PROTOCOL_F_CRYPTO_SESSION` feature has been successfully negotiated. It's a required feature for crypto devices.

VHOST_USER_CLOSE_CRYPTO_SESSION**id** 27**equivalent ioctl** N/A**master payload** u64

Close a session for crypto operation which was previously created by `VHOST_USER_CREATE_CRYPTO_SESSION`.

This request should be sent only when `VHOST_USER_PROTOCOL_F_CRYPTO_SESSION` feature has been successfully negotiated. It's a required feature for crypto devices.

VHOST_USER_POSTCOPY_ADVISE**id** 28**master payload** N/A**slave payload** userfault fd

When `VHOST_USER_PROTOCOL_F_PAGEFAULT` is supported, the master advises slave that a migration with postcopy enabled is underway, the slave must open a userfaultfd for later use. Note that at this stage the migration is still in precopy mode.

VHOST_USER_POSTCOPY_LISTEN**id** 29**master payload** N/A

Master advises slave that a transition to postcopy mode has happened. The slave must ensure that shared memory is registered with userfaultfd to cause faulting of non-present pages.

This is always sent sometime after a `VHOST_USER_POSTCOPY_ADVISE`, and thus only when `VHOST_USER_PROTOCOL_F_PAGEFAULT` is supported.

VHOST_USER_POSTCOPY_END**id** 30

slave payload u64

Master advises that postcopy migration has now completed. The slave must disable the userfaultfd. The response is an acknowledgement only.

When `VHOST_USER_PROTOCOL_F_PAGEFAULT` is supported, this message is sent at the end of the migration, after `VHOST_USER_POSTCOPY_LISTEN` was previously sent.

The value returned is an error indication; 0 is success.

VHOST_USER_GET_INFLIGHT_FD

id 31

equivalent ioctl N/A

master payload inflight description

When `VHOST_USER_PROTOCOL_F_INFLIGHT_SHMFD` protocol feature has been successfully negotiated, this message is submitted by master to get a shared buffer from slave. The shared buffer will be used to track inflight I/O by slave. QEMU should retrieve a new one when vm reset.

VHOST_USER_SET_INFLIGHT_FD

id 32

equivalent ioctl N/A

master payload inflight description

When `VHOST_USER_PROTOCOL_F_INFLIGHT_SHMFD` protocol feature has been successfully negotiated, this message is submitted by master to send the shared inflight buffer back to slave so that slave could get inflight I/O after a crash or restart.

VHOST_USER_GPU_SET_SOCKET

id 33

equivalent ioctl N/A

master payload N/A

Sets the GPU protocol socket file descriptor, which is passed as ancillary data. The GPU protocol is used to inform the master of rendering state and updates. See `vhost-user-gpu.rst` for details.

VHOST_USER_RESET_DEVICE

id 34

equivalent ioctl N/A

master payload N/A

slave payload N/A

Ask the vhost user backend to disable all rings and reset all internal device state to the initial state, ready to be reinitialized. The backend retains ownership of the device throughout the reset operation.

Only valid if the `VHOST_USER_PROTOCOL_F_RESET_DEVICE` protocol feature is set by the backend.

VHOST_USER_VRING_KICK

id 35

equivalent ioctl N/A

slave payload vring state description

master payload N/A

When the `VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS` protocol feature has been successfully negotiated, this message may be submitted by the master to indicate that a buffer was added to the vring instead of signalling it using the vring's kick file descriptor or having the slave rely on polling.

The `state.num` field is currently reserved and must be set to 0.

VHOST_USER_GET_MAX_MEM_SLOTS

id 36

equivalent ioctl N/A

slave payload u64

When the `VHOST_USER_PROTOCOL_F_CONFIGURE_MEM_SLOTS` protocol feature has been successfully negotiated, this message is submitted by master to the slave. The slave should return the message with a u64 payload containing the maximum number of memory slots for QEMU to expose to the guest. The value returned by the backend will be capped at the maximum number of ram slots which can be supported by the target platform.

VHOST_USER_ADD_MEM_REG

id 37

equivalent ioctl N/A

slave payload single memory region description

When the `VHOST_USER_PROTOCOL_F_CONFIGURE_MEM_SLOTS` protocol feature has been successfully negotiated, this message is submitted by the master to the slave. The message payload contains a memory region descriptor struct, describing a region of guest memory which the slave device must map in. When the `VHOST_USER_PROTOCOL_F_CONFIGURE_MEM_SLOTS` protocol feature has been successfully negotiated, along with the `VHOST_USER_REM_MEM_REG` message, this message is used to set and update the memory tables of the slave device.

VHOST_USER_REM_MEM_REG

id 38

equivalent ioctl N/A

slave payload single memory region description

When the `VHOST_USER_PROTOCOL_F_CONFIGURE_MEM_SLOTS` protocol feature has been successfully negotiated, this message is submitted by the master to the slave. The message payload contains a memory region descriptor struct, describing a region of guest memory which the slave device must unmap. When the `VHOST_USER_PROTOCOL_F_CONFIGURE_MEM_SLOTS` protocol feature has been successfully negotiated, along with the `VHOST_USER_ADD_MEM_REG` message, this message is used to set and update the memory tables of the slave device.

VHOST_USER_SET_STATUS

id 39

equivalent ioctl `VHOST_VDPA_SET_STATUS`

slave payload N/A

master payload u64

When the `VHOST_USER_PROTOCOL_F_STATUS` protocol feature has been successfully negotiated, this message is submitted by the master to notify the backend with updated device status as defined in the Virtio specification.

VHOST_USER_GET_STATUS

id 40

equivalent ioctl VHOST_VDPA_GET_STATUS

slave payload u64

master payload N/A

When the VHOST_USER_PROTOCOL_F_STATUS protocol feature has been successfully negotiated, this message is submitted by the master to query the backend for its device status as defined in the Virtio specification.

Slave message types

VHOST_USER_SLAVE_IOTLB_MSG

id 1

equivalent ioctl N/A (equivalent to VHOST_IOTLB_MSG message type)

slave payload struct vhost_iotlb_msg

master payload N/A

Send IOTLB messages with struct vhost_iotlb_msg as payload. Slave sends such requests to notify of an IOTLB miss, or an IOTLB access failure. If VHOST_USER_PROTOCOL_F_REPLY_ACK is negotiated, and slave set the VHOST_USER_NEED_REPLY flag, master must respond with zero when operation is successfully completed, or non-zero otherwise. This request should be send only when VIRTIO_F_IOMMU_PLATFORM feature has been successfully negotiated.

VHOST_USER_SLAVE_CONFIG_CHANGE_MSG

id 2

equivalent ioctl N/A

slave payload N/A

master payload N/A

When VHOST_USER_PROTOCOL_F_CONFIG is negotiated, vhost-user slave sends such messages to notify that the virtio device's configuration space has changed, for those host devices which can support such feature, host driver can send VHOST_USER_GET_CONFIG message to slave to get the latest content. If VHOST_USER_PROTOCOL_F_REPLY_ACK is negotiated, and slave set the VHOST_USER_NEED_REPLY flag, master must respond with zero when operation is successfully completed, or non-zero otherwise.

VHOST_USER_SLAVE_VRING_HOST_NOTIFIER_MSG

id 3

equivalent ioctl N/A

slave payload vring area description

master payload N/A

Sets host notifier for a specified queue. The queue index is contained in the u64 field of the vring area description. The host notifier is described by the file descriptor (typically it's a VFIO device fd) which is passed as ancillary data and the size (which is mmap size and should be the same as host page size) and offset (which is mmap offset) carried in the vring area description. QEMU can mmap the file descriptor based on the size and offset to get a memory range. Registering a host notifier means mapping this memory range to the VM as the specified queue's notify MMIO region. Slave sends this request to tell QEMU to de-register the existing notifier if any and register the new notifier if the request is sent with a file descriptor.

This request should be sent only when `VHOST_USER_PROTOCOL_F_HOST_NOTIFIER` protocol feature has been successfully negotiated.

VHOST_USER_SLAVE_VRING_CALL

id 4

equivalent ioctl N/A

slave payload vring state description

master payload N/A

When the `VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS` protocol feature has been successfully negotiated, this message may be submitted by the slave to indicate that a buffer was used from the vring instead of signalling this using the vring's call file descriptor or having the master relying on polling.

The `state.num` field is currently reserved and must be set to 0.

VHOST_USER_SLAVE_VRING_ERR

id 5

equivalent ioctl N/A

slave payload vring state description

master payload N/A

When the `VHOST_USER_PROTOCOL_F_INBAND_NOTIFICATIONS` protocol feature has been successfully negotiated, this message may be submitted by the slave to indicate that an error occurred on the specific vring, instead of signalling the error file descriptor set by the master via `VHOST_USER_SET_VRING_ERR`.

The `state.num` field is currently reserved and must be set to 0.

VHOST_USER_PROTOCOL_F_REPLY_ACK

The original vhost-user specification only demands replies for certain commands. This differs from the vhost protocol implementation where commands are sent over an `ioctl()` call and block until the client has completed.

With this protocol extension negotiated, the sender (QEMU) can set the `need_reply` [Bit 3] flag to any command. This indicates that the client **MUST** respond with a `Payload VhostUserMsg` indicating success or failure. The payload should be set to zero on success or non-zero on failure, unless the message already has an explicit reply body.

The response payload gives QEMU a deterministic indication of the result of the command. Today, QEMU is expected to terminate the main vhost-user loop upon receiving such errors. In future, qemu could be taught to be more resilient for selective requests.

For the message types that already solicit a reply from the client, the presence of `VHOST_USER_PROTOCOL_F_REPLY_ACK` or `need_reply` bit being set brings no behavioural change. (See the *Communication* section for details.)

4.10.4 Backend program conventions

vhost-user backends can provide various devices & services and may need to be configured manually depending on the use case. However, it is a good idea to follow the conventions listed here when possible. Users, QEMU or libvirt, can then rely on some common behaviour to avoid heterogeneous configuration and management of the backend programs and facilitate interoperability.

Each backend installed on a host system should come with at least one JSON file that conforms to the `vhost-user.json` schema. Each file informs the management applications about the backend type, and binary location. In addition, it

defines rules for management apps for picking the highest priority backend when multiple match the search criteria (see @VhostUserBackend documentation in the schema file).

If the backend is not capable of enabling a requested feature on the host (such as 3D acceleration with virgl), or the initialization failed, the backend should fail to start early and exit with a status `!= 0`. It may also print a message to `stderr` for further details.

The backend program must not daemonize itself, but it may be daemonized by the management layer. It may also have a restricted access to the system.

File descriptors 0, 1 and 2 will exist, and have regular `stdin/stdout/stderr` usage (they may have been redirected to `/dev/null` by the management layer, or to a log handler).

The backend program must end (as quickly and cleanly as possible) when the `SIGTERM` signal is received. Eventually, it may receive `SIGKILL` by the management layer after a few seconds.

The following command line options have an expected behaviour. They are mandatory, unless explicitly said differently:

- socket-path=PATH** This option specify the location of the vhost-user Unix domain socket. It is incompatible with `-fd`.
- fd=FDNUM** When this argument is given, the backend program is started with the vhost-user socket as file descriptor `FDNUM`. It is incompatible with `--socket-path`.
- print-capabilities** Output to `stdout` the backend capabilities in JSON format, and then exit successfully. Other options and arguments should be ignored, and the backend program should not perform its normal function. The capabilities can be reported dynamically depending on the host capabilities.

The JSON output is described in the `vhost-user.json` schema, by `@VHostUserBackendCapabilities`. Example:

```
{
  "type": "foo",
  "features": [
    "feature-a",
    "feature-b"
  ]
}
```

vhost-user-input

Command line options:

- evdev-path=PATH** Specify the linux input device.
(optional)
- no-grab** Do no request exclusive access to the input device.
(optional)

vhost-user-gpu

Command line options:

- render-node=PATH** Specify the GPU DRM render node.
(optional)

--virgl Enable virgl rendering support.
(optional)

vhost-user-blk

Command line options:

--blk-file=PATH Specify block device or file path.
(optional)

--read-only Enable read-only.
(optional)

4.11 Vhost-user-gpu Protocol

Licence This work is licensed under the terms of the GNU GPL, version 2 or later. See the COPYING file in the top-level directory.

Table of Contents

- *Vhost-user-gpu Protocol*
 - *Introduction*
 - *Wire format*
 - * *Header*
 - * *Payload types*
 - *VhostUserGpuCursorPos*
 - *VhostUserGpuCursorUpdate*
 - *VhostUserGpuScanout*
 - *VhostUserGpuUpdate*
 - *VhostUserGpuDMABUFScanout*
 - * *C structure*
 - * *Protocol features*
 - *Communication*
 - * *Message types*

4.11.1 Introduction

The vhost-user-gpu protocol is aiming at sharing the rendering result of a virtio-gpu, done from a vhost-user slave process to a vhost-user master process (such as QEMU). It bears a resemblance to a display server protocol, if you consider QEMU as the display server and the slave as the client, but in a very limited way. Typically, it will work by setting a scanout/display configuration, before sending flush events for the display updates. It will also update the cursor shape and position.

The protocol is sent over a UNIX domain stream socket, since it uses socket ancillary data to share opened file descriptors (DMABUF fds or shared memory). The socket is usually obtained via `VHOST_USER_GPU_SET_SOCKET`.

Requests are sent by the *slave*, and the optional replies by the *master*.

4.11.2 Wire format

Unless specified differently, numbers are in the machine native byte order.

A vhost-user-gpu message (request and reply) consists of 3 header fields and a payload.

request	flags	size	payload
---------	-------	------	---------

Header

request `u32`, type of the request

flags `u32`, 32-bit bit field:

- Bit 2 is the reply flag - needs to be set on each reply

size `u32`, size of the payload

Payload types

Depending on the request type, **payload** can be:

VhostUserGpuCursorPos

scanout-id	x	y
------------	---	---

scanout-id `u32`, the scanout where the cursor is located

x/y `u32`, the cursor position

VhostUserGpuCursorUpdate

pos	hot_x	hot_y	cursor
-----	-------	-------	--------

pos a `VhostUserGpuCursorPos`, the cursor location

hot_x/hot_y `u32`, the cursor hot location

cursor [`u32; 64 * 64`], 64x64 RGBA cursor data (PIXMAN_a8r8g8b8 format)

VhostUserGpuScanout

scanout-id	w	h
------------	---	---

scanout-id `u32`, the scanout configuration to set

w/h `u32`, the scanout width/height size

VhostUserGpuUpdate

scanout-id	x	y	w	h	data
------------	---	---	---	---	------

scanout-id u32, the scanout content to update

x/y/w/h u32, region of the update

data RGB data (PIXMAN_x8r8g8b8 format)

VhostUserGpuDMABUFScanout

scanout-id	x	y	w	h	fdw	fwh	stride	flags	fourcc
------------	---	---	---	---	-----	-----	--------	-------	--------

scanout-id u32, the scanout configuration to set

x/y u32, the location of the scanout within the DMABUF

w/h u32, the scanout width/height size

fdw/fdh/stride/flags u32, the DMABUF width/height/stride/flags

fourcc i32, the DMABUF fourcc

C structure

In QEMU the vhost-user-gpu message is implemented with the following struct:

```
typedef struct VhostUserGpuMsg {
    uint32_t request; /* VhostUserGpuRequest */
    uint32_t flags;
    uint32_t size; /* the following payload size */
    union {
        VhostUserGpuCursorPos cursor_pos;
        VhostUserGpuCursorUpdate cursor_update;
        VhostUserGpuScanout scanout;
        VhostUserGpuUpdate update;
        VhostUserGpuDMABUFScanout dmabuf_scanout;
        struct virtio_gpu_resp_display_info display_info;
        uint64_t u64;
    } payload;
} QEMU_PACKED VhostUserGpuMsg;
```

Protocol features

None yet.

As the protocol may need to evolve, new messages and communication changes are negotiated thanks to preliminary `VHOST_USER_GPU_GET_PROTOCOL_FEATURES` and `VHOST_USER_GPU_SET_PROTOCOL_FEATURES` requests.

4.11.3 Communication

Message types

VHOST_USER_GPU_GET_PROTOCOL_FEATURES

id 1

request payload N/A

reply payload u64

Get the supported protocol features bitmask.

VHOST_USER_GPU_SET_PROTOCOL_FEATURES

id 2

request payload u64

reply payload N/A

Enable protocol features using a bitmask.

VHOST_USER_GPU_GET_DISPLAY_INFO

id 3

request payload N/A

reply payload struct virtio_gpu_resp_display_info (from virtio specification)

Get the preferred display configuration.

VHOST_USER_GPU_CURSOR_POS

id 4

request payload VhostUserGpuCursorPos

reply payload N/A

Set/show the cursor position.

VHOST_USER_GPU_CURSOR_POS_HIDE

id 5

request payload VhostUserGpuCursorPos

reply payload N/A

Set/hide the cursor.

VHOST_USER_GPU_CURSOR_UPDATE

id 6

request payload VhostUserGpuCursorUpdate

reply payload N/A

Update the cursor shape and location.

VHOST_USER_GPU_SCANOUT

id 7

request payload VhostUserGpuScanout

reply payload N/A

Set the scanout resolution. To disable a scanout, the dimensions width/height are set to 0.

VHOST_USER_GPU_UPDATE

id 8

request payload VhostUserGpuUpdate

reply payload N/A

Update the scanout content. The data payload contains the graphical bits. The display should be flushed and presented.

VHOST_USER_GPU_DMABUF_SCANOUT

id 9

request payload VhostUserGpuDMABUFScanout

reply payload N/A

Set the scanout resolution/configuration, and share a DMABUF file descriptor for the scanout content, which is passed as ancillary data. To disable a scanout, the dimensions width/height are set to 0, there is no file descriptor passed.

VHOST_USER_GPU_DMABUF_UPDATE

id 10

request payload VhostUserGpuUpdate

reply payload empty payload

The display should be flushed and presented according to updated region from VhostUserGpuUpdate.

Note: there is no data payload, since the scanout is shared thanks to DMABUF, that must have been set previously with VHOST_USER_GPU_DMABUF_SCANOUT.

4.12 Vhost-vdpa Protocol

4.12.1 Introduction

vDPA(Virtual data path acceleration) device is a device that uses a datapath which complies with the virtio specifications with vendor specific control path. vDPA devices can be both physically located on the hardware or emulated by software.

This document describes the vDPA support in qemu

Here is the kernel commit here <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?id=4c8cf31885f69e86be0b5b9e6677a26797365e1d>

TODO : More information will add later

QEMU System Emulation Guest Hardware Specifications

Contents:

5.1 POWER9 XIVE interrupt controller

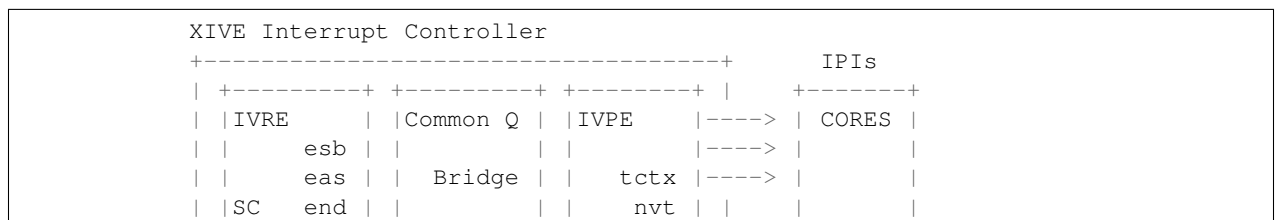
The POWER9 processor comes with a new interrupt controller architecture, called XIVE as “eXternal Interrupt Virtualization Engine”.

Compared to the previous architecture, the main characteristics of XIVE are to support a larger number of interrupt sources and to deliver interrupts directly to virtual processors without hypervisor assistance. This removes the context switches required for the delivery process.

5.1.1 XIVE architecture

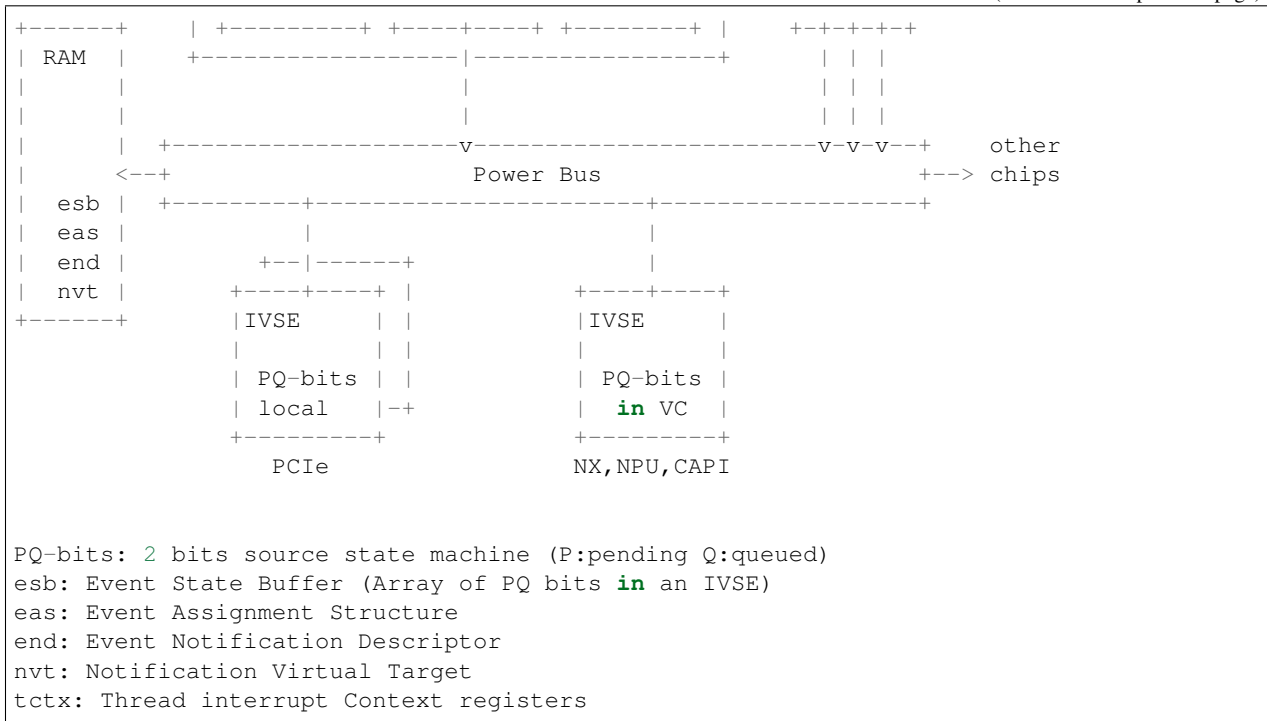
The XIVE IC is composed of three sub-engines, each taking care of a processing layer of external interrupts:

- Interrupt Virtualization Source Engine (IVSE), or Source Controller (SC). These are found in PCI PHBs, in the Processor Service Interface (PSI) host bridge Controller, but also inside the main controller for the core IPIs and other sub-chips (NX, CAP, NPU) of the chip/processor. They are configured to feed the IVRE with events.
- Interrupt Virtualization Routing Engine (IVRE) or Virtualization Controller (VC). It handles event coalescing and perform interrupt routing by matching an event source number with an Event Notification Descriptor (END).
- Interrupt Virtualization Presentation Engine (IVPE) or Presentation Controller (PC). It maintains the interrupt context state of each thread and handles the delivery of the external interrupt to the thread.



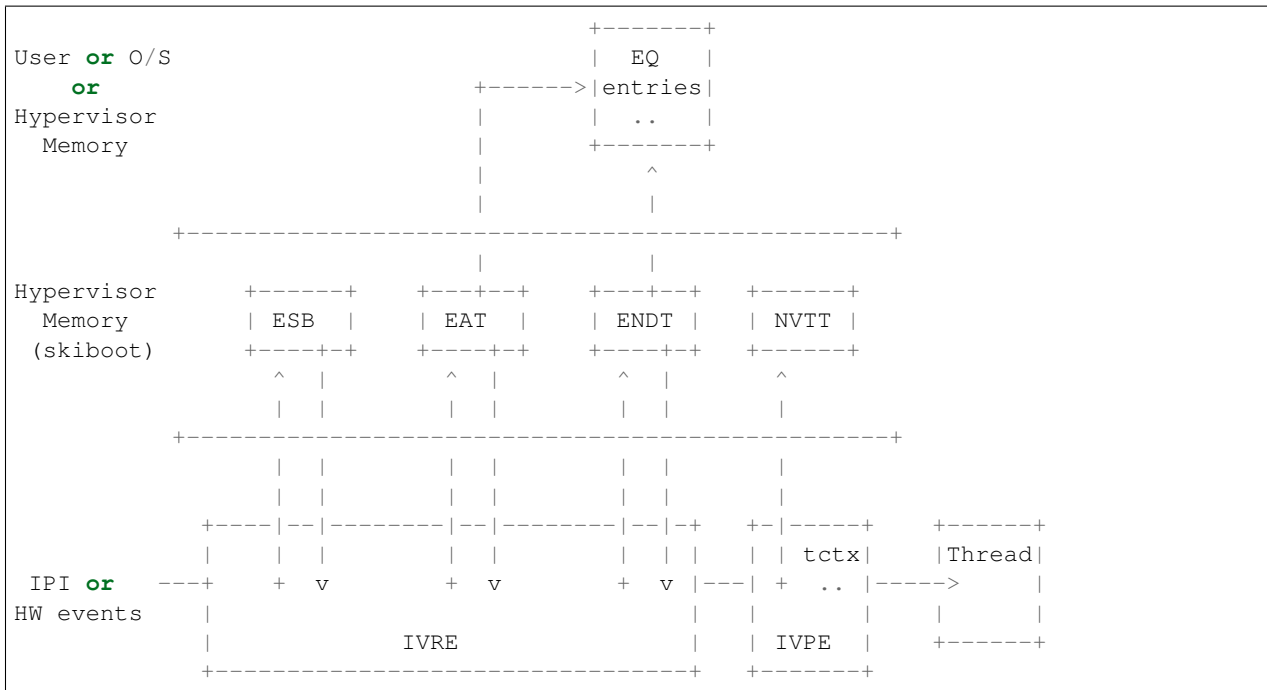
(continues on next page)

(continued from previous page)



XIVE internal tables

Each of the sub-engines uses a set of tables to redirect interrupts from event sources to CPU threads.



The IVSE have a 2-bits state machine, P for pending and Q for queued, for each source that allows events to be triggered. They are stored in an Event State Buffer (ESB) array and can be controlled by MMIOs.

If the event is let through, the IVRE looks up in the Event Assignment Structure (EAS) table for an Event Notification Descriptor (END) configured for the source. Each Event Notification Descriptor defines a notification path to a CPU and an in-memory Event Queue, in which will be enqueued an EQ data for the O/S to pull.

The IVPE determines if a Notification Virtual Target (NVT) can handle the event by scanning the thread contexts of the VCPUs dispatched on the processor HW threads. It maintains the interrupt context state of each thread in a NVT table.

XIVE thread interrupt context

The XIVE presenter can generate four different exceptions to its HW threads:

- hypervisor exception
- O/S exception
- Event-Based Branch (user level)
- msgsnd (doorbell)

Each exception has a state independent from the others called a Thread Interrupt Management context. This context is a set of registers which lets the thread handle priority management and interrupt acknowledgment among other things. The most important ones being :

- Interrupt Priority Register (PIPR)
- Interrupt Pending Buffer (IPB)
- Current Processor Priority (CPPR)
- Notification Source Register (NSR)

TIMA

The Thread Interrupt Management registers are accessible through a specific MMIO region, called the Thread Interrupt Management Area (TIMA), four aligned pages, each exposing a different view of the registers. First page (page address ending in 0b00) gives access to the entire context and is reserved for the ring 0 view for the physical thread context. The second (page address ending in 0b01) is for the hypervisor, ring 1 view. The third (page address ending in 0b10) is for the operating system, ring 2 view. The fourth (page address ending in 0b11) is for user level, ring 3 view.

Interrupt flow from an O/S perspective

After an event data has been enqueued in the O/S Event Queue, the IVPE raises the bit corresponding to the priority of the pending interrupt in the register IPB (Interrupt Pending Buffer) to indicate that an event is pending in one of the 8 priority queues. The Pending Interrupt Priority Register (PIPR) is also updated using the IPB. This register represent the priority of the most favored pending notification.

The PIPR is then compared to the Current Processor Priority Register (CPPR). If it is more favored (numerically less than), the CPU interrupt line is raised and the EO bit of the Notification Source Register (NSR) is updated to notify the presence of an exception for the O/S. The O/S acknowledges the interrupt with a special load in the Thread Interrupt Management Area.

The O/S handles the interrupt and when done, performs an EOI using a MMIO operation on the ESB management page of the associate source.

5.1.2 Overview of the QEMU models for XIVE

The XiveSource models the IVSE in general, internal and external. It handles the source ESBs and the MMIO interface to control them.

The XiveNotifier is a small helper interface interconnecting the XiveSource to the XiveRouter.

The XiveRouter is an abstract model acting as a combined IVRE and IVPE. It routes event notifications using the EAS and END tables to the IVPE sub-engine which does a CAM scan to find a CPU to deliver the exception. Storage should be provided by the inheriting classes.

XiveEndSource is a special source object. It exposes the END ESB MMIOs of the Event Queues which are used for coalescing event notifications and for escalation. Not used on the field, only to sync the EQ cache in OPAL.

Finally, the XiveTCTX contains the interrupt state context of a thread, four sets of registers, one for each exception that can be delivered to a CPU. These contexts are scanned by the IVPE to find a matching VP when a notification is triggered. It also models the Thread Interrupt Management Area (TIMA), which exposes the thread context registers to the CPU for interrupt management.

5.2 XIVE for sPAPR (pseries machines)

The POWER9 processor comes with a new interrupt controller architecture, called XIVE as “eXternal Interrupt Virtualization Engine”. It supports a larger number of interrupt sources and offers virtualization features which enables the HW to deliver interrupts directly to virtual processors without hypervisor assistance.

A QEMU `pseries` machine (which is PAPR compliant) using POWER9 processors can run under two interrupt modes:

- *Legacy Compatibility Mode*

the hypervisor provides identical interfaces and similar functionality to PAPR+ Version 2.7. This is the default mode

It is also referred as *XICS* in QEMU.

- *XIVE native exploitation mode*

the hypervisor provides new interfaces to manage the XIVE control structures, and provides direct control for interrupt management through MMIO pages.

Which interrupt modes can be used by the machine is negotiated with the guest O/S during the Client Architecture Support negotiation sequence. The two modes are mutually exclusive.

Both interrupt mode share the same IRQ number space. See below for the layout.

5.2.1 CAS Negotiation

QEMU advertises the supported interrupt modes in the device tree property `ibm,arch-vec-5-platform-support` in byte 23 and the OS Selection for XIVE is indicated in the `ibm,architecture-vec-5` property byte 23.

The interrupt modes supported by the machine depend on the CPU type (POWER9 is required for XIVE) but also on the machine property `ic-mode` which can be set on the command line. It can take the following values: `xics`, `xive`, and `dual` which is the default mode. `dual` means that both modes XICS **and** XIVE are supported and if the guest OS supports XIVE, this mode will be selected.

The chosen interrupt mode is activated after a reconfiguration done in a machine reset.

5.2.2 KVM negotiation

When the guest starts under KVM, the capabilities of the host kernel and QEMU are also negotiated. Depending on the version of the host kernel, KVM will advertise the XIVE capability to QEMU or not.

Nevertheless, the available interrupt modes in the machine should not depend on the XIVE KVM capability of the host. On older kernels without XIVE KVM support, QEMU will use the emulated XIVE device as a fallback and on newer kernels (≥ 5.2), the KVM XIVE device.

XIVE native exploitation mode is not supported for KVM nested guests, VMs running under a L1 hypervisor (KVM on pSeries). In that case, the hypervisor will not advertise the KVM capability and QEMU will use the emulated XIVE device, same as for older versions of KVM.

As a final refinement, the user can also switch the use of the KVM device with the machine option `kernel_irqchip`.

XIVE support in KVM

For guest OSes supporting XIVE, the resulting interrupt modes on host kernels with XIVE KVM support are the following:

ic-mode	kernel_irqchip		
/	allowed (default)	off	on
dual (default)	XIVE KVM	XIVE emul.	XIVE KVM
xive	XIVE KVM	XIVE emul.	XIVE KVM
xics	XICS KVM	XICS emul.	XICS KVM

For legacy guest OSes without XIVE support, the resulting interrupt modes are the following:

ic-mode	kernel_irqchip		
/	allowed (default)	off	on
dual (default)	XICS KVM	XICS emul.	XICS KVM
xive	QEMU error(3)	QEMU error(3)	QEMU error(3)
xics	XICS KVM	XICS emul.	XICS KVM

- (3) QEMU fails at CAS with Guest requested unavailable interrupt mode (XICS), either don't set the `ic-mode` machine property or try `ic-mode=xics` or `ic-mode=dual`

No XIVE support in KVM

For guest OSes supporting XIVE, the resulting interrupt modes on host kernels without XIVE KVM support are the following:

ic-mode	kernel_irqchip		
/	allowed (default)	off	on
dual (default)	XIVE emul.(1)	XIVE emul.	QEMU error (2)
xive	XIVE emul.(1)	XIVE emul.	QEMU error (2)
xics	XICS KVM	XICS emul.	XICS KVM

- (1) QEMU warns with warning: `kernel_irqchip requested but unavailable: IRQ_XIVE capability must be present for KVM` In some cases (old host kernels or KVM nested guests), one may hit a QEMU/KVM incompatibility due to device destruction in reset. QEMU fails with KVM is incompatible with `ic-mode=dual, kernel-irqchip=on`

- (2) QEMU fails with `kernel_irqchip` requested but unavailable: `IRQ_XIVE` capability must be present for KVM

For legacy guest OSes without XIVE support, the resulting interrupt modes are the following:

ic-mode	kernel_irqchip		
/	allowed (default)	off	on
dual (default)	QEMU error(4)	XICS emul.	QEMU error(4)
xive	QEMU error(3)	QEMU error(3)	QEMU error(3)
xics	XICS KVM	XICS emul.	XICS KVM

- (3) QEMU fails at CAS with Guest requested unavailable interrupt mode (XICS), either don't set the `ic-mode` machine property or try `ic-mode=xics` or `ic-mode=dual`
- (4) QEMU/KVM incompatibility due to device destruction in reset. QEMU fails with KVM is incompatible with `ic-mode=dual, kernel-irqchip=on`

5.2.3 XIVE Device tree properties

The properties for the PAPR interrupt controller node when the *XIVE native exploitation mode* is selected should contain:

- `device_type`
value should be “power-ivpe”.
- `compatible`
value should be “ibm,power-ivpe”.
- `reg`
contains the base address and size of the thread interrupt management areas (TIMA), for the User level and for the Guest OS level. Only the Guest OS level is taken into account today.
- `ibm,xive-eq-sizes`
the size of the event queues. One cell per size supported, contains log2 of size, in ascending order.
- `ibm,xive-lisn-ranges`
the IRQ interrupt number ranges assigned to the guest for the IPIs.

The root node also exports :

- `ibm,plat-res-int-priorities`
contains a list of priorities that the hypervisor has reserved for its own use.

5.2.4 IRQ number space

IRQ Number space of the `pseries` machine is 8K wide and is the same for both interrupt mode. The different ranges are defined as follow :

- `0x0000 .. 0x0FFF` 4K CPU IPIs (only used under XIVE)
- `0x1000 .. 0x1000` 1 EPOW
- `0x1001 .. 0x1001` 1 HOTPLUG
- `0x1002 .. 0x10FF` unused

- 0x1100 .. 0x11FF 256 VIO devices
- 0x1200 .. 0x127F 32x4 LSIs for PHB devices
- 0x1280 .. 0x12FF unused
- 0x1300 .. 0x1FFF PHB MSIs (dynamically allocated)

5.2.5 Monitoring XIVE

The state of the XIVE interrupt controller can be queried through the monitor commands `info pic`. The output comes in two parts.

First, the state of the thread interrupt context registers is dumped for each CPU :

```
(qemu) info pic
CPU[0000]:  QW   NSR CPPR IPB LSMFB ACK#  INC AGE PIPR  W2
CPU[0000]: USER    00  00  00    00  00  00  00  00 00000000
CPU[0000]: OS      00  ff  00    00  ff  00  ff  ff 80000400
CPU[0000]: POOL    00  00  00    00  00  00  00  00 00000000
CPU[0000]: PHYS    00  00  00    00  00  00  00  ff 00000000
...
```

In the case of a `pseries` machine, QEMU acts as the hypervisor and only the `O/S` and `USER` register rings make sense. `W2` contains the vCPU CAM line which is set to the VP identifier.

Then comes the routing information which aggregates the `EAS` and the `END` configuration:

```
...
LISN      PQ      EISN      CPU/PRIO EQ
00000000 MSI  --      00000010  0/6    380/16384 @1fe3e0000 ^1 [ 80000010 ... ]
00000001 MSI  --      00000010  1/6    305/16384 @1fc230000 ^1 [ 80000010 ... ]
00000002 MSI  --      00000010  2/6    220/16384 @1fc2f0000 ^1 [ 80000010 ... ]
00000003 MSI  --      00000010  3/6    201/16384 @1fc390000 ^1 [ 80000010 ... ]
00000004 MSI  -Q M 00000000
00000005 MSI  -Q M 00000000
00000006 MSI  -Q M 00000000
00000007 MSI  -Q M 00000000
00001000 MSI  --      00000012  0/6    380/16384 @1fe3e0000 ^1 [ 80000010 ... ]
00001001 MSI  --      00000013  0/6    380/16384 @1fe3e0000 ^1 [ 80000010 ... ]
00001100 MSI  --      00000100  1/6    305/16384 @1fc230000 ^1 [ 80000010 ... ]
00001101 MSI  -Q M 00000000
00001200 LSI  -Q M 00000000
00001201 LSI  -Q M 00000000
00001202 LSI  -Q M 00000000
00001203 LSI  -Q M 00000000
00001300 MSI  --      00000102  1/6    305/16384 @1fc230000 ^1 [ 80000010 ... ]
00001301 MSI  --      00000103  2/6    220/16384 @1fc2f0000 ^1 [ 80000010 ... ]
00001302 MSI  --      00000104  3/6    201/16384 @1fc390000 ^1 [ 80000010 ... ]
```

The source information and configuration:

- The `LISN` column outputs the interrupt number of the source in range [0x0 ... 0x1FFF] and its type : `MSI` or `LSI`
- The `PQ` column reflects the state of the `PQ` bits of the source :
 - -- source is ready to take events
 - P- an event was sent and an `EOI` is `PENDING`

- PQ an event was QUEUED
- -Q source is OFF

a M indicates that source is *MASKED* at the EAS level,

The targeting configuration :

- The EISN column is the event data that will be queued in the event queue of the O/S.
- The CPU/PRIO column is the tuple defining the CPU number and priority queue serving the source.
- The EQ column outputs :
 - the current index of the event queue/ the max number of entries
 - the O/S event queue address
 - the toggle bit
 - the last entries that were pushed in the event queue.

5.3 NUMA mechanics for sPAPR (pseries machines)

NUMA in sPAPR works different than the System Locality Distance Information Table (SLIT) in ACPI. The logic is explained in the LOPAPR 1.1 chapter 15, “Non Uniform Memory Access (NUMA) Option”. This document aims to complement this specification, providing details of the elements that impacts how QEMU views NUMA in pseries.

5.3.1 Associativity and ibm,associativity property

Associativity is defined as a group of platform resources that has similar mean performance (or in our context here, distance) relative to everyone else outside of the group.

The format of the ibm,associativity property varies with the value of bit 0 of byte 5 of the ibm,architecture-vec-5 property. The format with bit 0 equal to zero is deprecated. The current format, with the bit 0 with the value of one, makes ibm,associativity property represent the physical hierarchy of the platform, as one or more lists that starts with the highest level grouping up to the smallest. Considering the following topology:

Mem M1 ---- Proc P1		
-----		Socket S1 ---
chip C1		
		HW module 1 (MOD1)
Mem M2 ---- Proc P2		
-----		Socket S2 ---
chip C2		

The ibm,associativity property for the processors would be:

- P1: {MOD1, S1, C1, P1}
- P2: {MOD1, S2, C2, P2}

Each allocable resource has an ibm,associativity property. The LOPAPR specification allows multiple lists to be present in this property, considering that the same resource can have multiple connections to the platform.

5.3.2 Relative Performance Distance and ibm,associativity-reference-points

The `ibm,associativity-reference-points` property is an array that is used to define the relevant performance/distance related boundaries, defining the NUMA levels for the platform.

The definition of its elements also varies with the value of bit 0 of byte 5 of the `ibm,architecture-vec-5` property. The format with bit 0 equal to zero is also deprecated. With the current format, each integer of the `ibm,associativity-reference-points` represents an 1 based ordinal index (i.e. the first element is 1) of the `ibm,associativity` array. The first boundary is the most significant to application performance, followed by less significant boundaries. Allocated resources that belongs to the same performance boundaries are expected to have relative NUMA distance that matches the relevancy of the boundary itself. Resources that belongs to the same first boundary will have the shortest distance from each other. Subsequent boundaries represents greater distances and degraded performance.

Using the previous example, the following setting reference points defines three NUMA levels:

- `ibm,associativity-reference-points = {0x3, 0x2, 0x1}`

The first NUMA level (0x3) is interpreted as the third element of each `ibm,associativity` array, the second level is the second element and the third level is the first element. Let's also consider that elements belonging to the first NUMA level have distance equal to 10 from each other, and each NUMA level doubles the distance from the previous. This means that the second would be 20 and the third level 40. For the P1 and P2 processors, we would have the following NUMA levels:

```
* ibm,associativity-reference-points = {0x3, 0x2, 0x1}

* P1: associativity{MOD1, S1, C1, P1}

First NUMA level (0x3) => associativity[2] = C1
Second NUMA level (0x2) => associativity[1] = S1
Third NUMA level (0x1) => associativity[0] = MOD1

* P2: associativity{MOD1, S2, C2, P2}

First NUMA level (0x3) => associativity[2] = C2
Second NUMA level (0x2) => associativity[1] = S2
Third NUMA level (0x1) => associativity[0] = MOD1

P1 and P2 have the same third NUMA level, MOD1: Distance between them = 40
```

Changing the `ibm,associativity-reference-points` array changes the performance distance attributes for the same associativity arrays, as the following example illustrates:

```
* ibm,associativity-reference-points = {0x2}

* P1: associativity{MOD1, S1, C1, P1}

First NUMA level (0x2) => associativity[1] = S1

* P2: associativity{MOD1, S2, C2, P2}

First NUMA level (0x2) => associativity[1] = S2

P1 and P2 does not have a common performance boundary. Since this is a one level
NUMA configuration, distance between them is one boundary above the first
level, 20.
```

In a hypothetical platform where all resources inside the same hardware module is considered to be on the same performance boundary:

```
* ibm,associativity-reference-points = {0x1}

* P1: associativity{MOD1, S1, C1, P1}

First NUMA level (0x1) => associativity[0] = MOD0

* P2: associativity{MOD1, S2, C2, P2}

First NUMA level (0x1) => associativity[0] = MOD0

P1 and P2 belongs to the same first order boundary. The distance between them
is 10.
```

5.4 How the pseries Linux guest calculates NUMA distances

Another key difference between ACPI SLIT and the LOPAPR regarding NUMA is how the distances are expressed. The SLIT table provides the NUMA distance value between the relevant resources. LOPAPR does not provide a standard way to calculate it. We have the `ibm,associativity` for each resource, which provides a common-performance hierarchy, and the `ibm,associativity-reference-points` array that tells which level of associativity is considered to be relevant or not.

The result is that each OS is free to implement and to interpret the distance as it sees fit. For the pseries Linux guest, each level of NUMA duplicates the distance of the previous level, and the maximum amount of levels is limited to `MAX_DISTANCE_REF_POINTS = 4` (from `arch/powerpc/mm/numa.c` in the kernel tree). This results in the following distances:

- both resources in the first NUMA level: 10
- resources one NUMA level apart: 20
- resources two NUMA levels apart: 40
- resources three NUMA levels apart: 80
- resources four NUMA levels apart: 160

5.5 pseries NUMA mechanics

Starting in QEMU 5.2, the pseries machine considers user input when setting NUMA topology of the guest. The overall design is:

- `ibm,associativity-reference-points` is set to `{0x4, 0x3, 0x2, 0x1}`, allowing for 4 distinct NUMA distance values based on the NUMA levels
- `ibm,max-associativity-domains` supports multiple associativity domains in all NUMA levels, granting user flexibility
- `ibm,associativity` for all resources varies with user input

These changes are only effective for pseries-5.2 and newer machines that are created with more than one NUMA node (disconsidering NUMA nodes created by the machine itself, e.g. NVLink 2 GPUs). The now legacy support has been around for such a long time, with users seeing NUMA distances 10 and 40 (and 80 if using NVLink2 GPUs), and there is no need to disrupt the existing experience of those guests.

To bring the user experience x86 users have when tuning up NUMA, we had to operate under the current pseries Linux kernel logic described in *How the pseries Linux guest calculates NUMA distances*. The result is that we needed to translate NUMA distance user input to pseries Linux kernel input.

5.5.1 Translating user distance to kernel distance

User input for NUMA distance can vary from 10 to 254. We need to translate that to the values that the Linux kernel operates on (10, 20, 40, 80, 160). This is how it is being done:

- user distance 11 to 30 will be interpreted as 20
- user distance 31 to 60 will be interpreted as 40
- user distance 61 to 120 will be interpreted as 80
- user distance 121 and beyond will be interpreted as 160
- user distance 10 stays 10

The reasoning behind this approximation is to avoid any round up to the local distance (10), keeping it exclusive to the 4th NUMA level (which is still exclusive to the node_id). All other ranges were chosen under the developer discretion of what would be (somewhat) sensible considering the user input. Any other strategy can be used here, but in the end the reality is that we'll have to accept that a large array of values will be translated to the same NUMA topology in the guest, e.g. this user input:

```

    0   1   2
0  10  31 120
1  31  10  30
2 120  30  10
```

And this other user input:

```

    0   1   2
0  10  60  61
1  60  10  11
2  61  11  10
```

Will both be translated to the same values internally:

```

    0   1   2
0  10  40  80
1  40  10  20
2  80  20  10
```

Users are encouraged to use only the kernel values in the NUMA definition to avoid being taken by surprise with that the guest is actually seeing in the topology. There are enough potential surprises that are inherent to the associativity domain assignment process, discussed below.

5.5.2 How associativity domains are assigned

LOPAPR allows more than one associativity array (or 'string') per allocated resource. This would be used to represent that the resource has multiple connections with the board, and then the operational system, when deciding NUMA distancing, should consider the associativity information that provides the shortest distance.

The spapr implementation does not support multiple associativity arrays per resource, neither does the pseries Linux kernel. We'll have to represent the NUMA topology using one associativity per resource, which means that choices and compromises are going to be made.

Consider the following NUMA topology entered by user input:

	0	1	2	3
0	10	40	20	40
1	40	10	80	40
2	20	80	10	20
3	40	40	20	10

All the associativity arrays are initialized with NUMA id in all associativity domains:

- node 0: 0 0 0 0
- node 1: 1 1 1 1
- node 2: 2 2 2 2
- node 3: 3 3 3 3

Honoring just the relative distances of node 0 to every other node, we find the NUMA level matches (considering the reference points {0x4, 0x3, 0x2, 0x1}) for each distance:

- distance from 0 to 1 is 40 (no match at 0x4 and 0x3, will match at 0x2)
- distance from 0 to 2 is 20 (no match at 0x4, will match at 0x3)
- distance from 0 to 3 is 40 (no match at 0x4 and 0x3, will match at 0x2)

We'll copy the associativity domains of node 0 to all other nodes, based on the NUMA level matches. Between 0 and 1, a match in 0x2, we'll also copy the domains 0x2 and 0x1 from 0 to 1 as well. This will give us:

- node 0: 0 0 0 0
- node 1: 0 0 1 1

Doing the same to node 2 and node 3, these are the associativity arrays after considering all matches with node 0:

- node 0: 0 0 0 0
- node 1: 0 0 1 1
- node 2: 0 0 0 2
- node 3: 0 0 3 3

The distances related to node 0 are accounted for. For node 1, and keeping in mind that we don't need to revisit node 0 again, the distance from node 1 to 2 is 80, matching at 0x1, and distance from 1 to 3 is 40, match in 0x2. Repeating the same logic of copying all domains up to the NUMA level match:

- node 0: 0 0 0 0
- node 1: 1 0 1 1
- node 2: 1 0 0 2
- node 3: 1 0 3 3

In the last step we will analyze just nodes 2 and 3. The desired distance between 2 and 3 is 20, i.e. a match in 0x3:

- node 0: 0 0 0 0
- node 1: 1 0 1 1
- node 2: 1 0 0 2
- node 3: 1 0 0 3

The kernel will read these arrays and will calculate the following NUMA topology for the guest:

	0	1	2	3
0	10	40	20	20
1	40	10	40	40
2	20	40	10	20
3	20	40	20	10

Note that this is not what the user wanted - the desired distance between 0 and 3 is 40, we calculated it as 20. This is what the current logic and implementation constraints of the kernel and QEMU will provide inside the LOPAPR specification.

Users are welcome to use this knowledge and experiment with the input to get the NUMA topology they want, or as closer as they want. The important thing is to keep expectations up to par with what we are capable of provide at this moment: an approximation.

5.5.3 Limitations of the implementation

As mentioned above, the pSeries NUMA distance logic is, in fact, a way to approximate user choice. The Linux kernel, and PAPR itself, does not provide QEMU with the ways to fully map user input to actual NUMA distance the guest will use. These limitations creates two notable limitations in our support:

- Asymmetrical topologies aren't supported. We only support NUMA topologies where the distance from node A to B is always the same as B to A. We do not support any A-B pair where the distance back and forth is asymmetric. For example, the following topology isn't supported and the pSeries guest will not boot with this user input:

	0	1
0	10	40
1	20	10

- 'non-transitive' topologies will be poorly translated to the guest. This is the kind of topology where the distance from a node A to B is X, B to C is X, but the distance A to C is not X. E.g.:

	0	1	2	3
0	10	20	20	40
1	20	10	80	40
2	20	80	10	20
3	40	40	20	10

In the example above, distance 0 to 2 **is** 20, 2 to 3 **is** 20, but 0 to 3 **is** 40. The kernel will always match **with** the shortest associativity domain possible, **and** we're attempting to retain the previous established relations between the nodes. This means that a distance equal to 20 between nodes 0 **and** 2 **and** the same distance 20 between nodes 2 **and** 3 will cause the distance between 0 **and** 3 to also be 20.

5.6 Legacy (5.1 and older) pseries NUMA mechanics

In short, we can summarize the NUMA distances seem in pseries Linux guests, using QEMU up to 5.1, as follows:

- local distance, i.e. the distance of the resource to its own NUMA node: 10
- if it's a NVLink GPU device, distance: 80
- every other resource, distance: 40

The way the pseries Linux guest calculates NUMA distances has a direct effect on what QEMU users can expect when doing NUMA tuning. As of QEMU 5.1, this is the default `ibm,associativity-reference-points` being used in the pseries machine:

```
ibm,associativity-reference-points = {0x4, 0x4, 0x2}
```

The first and second level are equal, 0x4, and a third one was added in commit a6030d7e0b35 exclusively for NVLink GPUs support. This means that regardless of how the `ibm,associativity` properties are being created in the device tree, the pseries Linux guest will only recognize three scenarios as far as NUMA distance goes:

- if the resources belongs to the same first NUMA level = 10
- second level is skipped since it's equal to the first
- all resources that aren't a NVLink GPU, it is guaranteed that they will belong to the same third NUMA level, having distance = 40
- for NVLink GPUs, distance = 80 from everything else

This also means that user input in QEMU command line does not change the NUMA distancing inside the guest for the pseries machine.

5.7 QEMU and ACPI BIOS Generic Event Device interface

The *ACPI Generic Event Device* (GED) is a HW reduced platform specific device introduced in ACPI v6.1 that handles all platform events, including the hotplug ones. GED is modelled as a device in the namespace with a `_HID` defined to be `ACPI0013`. This document describes the interface between QEMU and the ACPI BIOS.

GED allows HW reduced platforms to handle interrupts in ACPI ASL statements. It follows a very similar approach to the `_EVT` method from GPIO events. All interrupts are listed in `_CRS` and the handler is written in `_EVT` method. However, the QEMU implementation uses a single interrupt for the GED device, relying on an IO memory region to communicate the type of device affected by the interrupt. This way, we can support up to 32 events with a unique interrupt.

Here is an example,

```
Device (\_SB.GED)
{
    Name (_HID, "ACPI0013")
    Name (_UID, Zero)
    Name (_CRS, ResourceTemplate ()
    {
        Interrupt (ResourceConsumer, Edge, ActiveHigh, Exclusive, ,, )
        {
            0x00000029,
        }
    })
    OperationRegion (EREG, SystemMemory, 0x09080000, 0x04)
    Field (EREG, DWordAcc, NoLock, WriteAsZeros)
    {
        ESEL, 32
    }
    Method (_EVT, 1, Serialized)
    {
        Local0 = ESEL // ESEL = IO memory region which specifies the
                        // device type.
        If ((Local0 & One) == One)
        {
```

(continues on next page)

(continued from previous page)

```

        MethodEvent1()
    }
    If ((Local0 & 0x2) == 0x2)
    {
        MethodEvent2()
    }
    ...
}

```

5.7.1 GED IO interface (4 byte access)

read access:

```

[0x0-0x3] Event selector bit field (32 bit) set by QEMU.

bits:
    0: Memory hotplug event
    1: System power down event
    2: NVDIMM hotplug event
    3-31: Reserved

```

write access:

Nothing is expected to be written into GED IO memory

5.8 QEMU TPM Device

5.8.1 Guest-side hardware interface

TIS interface

The QEMU TPM emulation implements a TPM TIS hardware interface following the Trusted Computing Group's specification "TCG PC Client Specific TPM Interface Specification (TIS)", Specification Version 1.3, 21 March 2013. (see the [TIS specification](#), or a later version of it).

The TIS interface makes a memory mapped IO region in the area 0xfed40000-0xfed44fff available to the guest operating system.

QEMU files related to TPM TIS interface:

- hw/tpm/tpm_tis_common.c
- hw/tpm/tpm_tis_isa.c
- hw/tpm/tpm_tis_sysbus.c
- hw/tpm/tpm_tis.h

Both an ISA device and a sysbus device are available. The former is used with pc/q35 machine while the latter can be instantiated in the Arm virt machine.

CRB interface

QEMU also implements a TPM CRB interface following the Trusted Computing Group’s specification “TCG PC Client Platform TPM Profile (PTP) Specification”, Family “2.0”, Level 00 Revision 01.03 v22, May 22, 2017. (see the [CRB specification](#), or a later version of it)

The CRB interface makes a memory mapped IO region in the area 0xfed40000-0xfed40fff (1 locality) available to the guest operating system.

QEMU files related to TPM CRB interface:

- hw/tpm/tpm_crb.c

SPAPR interface

pSeries (ppc64) machines offer a tpm-spapr device model.

QEMU files related to the SPAPR interface:

- hw/tpm/tpm_spapr.c

5.8.2 fw_cfg interface

The bios/firmware may read the "etc/tpm/config" fw_cfg entry for configuring the guest appropriately.

The entry of 6 bytes has the following content, in little-endian:

```
#define TPM_VERSION_UNSPEC      0
#define TPM_VERSION_1_2        1
#define TPM_VERSION_2_0        2

#define TPM_PPI_VERSION_NONE    0
#define TPM_PPI_VERSION_1_30   1

struct FwCfgTPMConfig {
    uint32_t tpmppi_address;      /* PPI memory location */
    uint8_t  tpm_version;        /* TPM version */
    uint8_t  tpmppi_version;     /* PPI version */
};
```

5.8.3 ACPI interface

The TPM device is defined with ACPI ID “PNP0C31”. QEMU builds a SSDT and passes it into the guest through the fw_cfg device. The device description contains the base address of the TIS interface 0xfed40000 and the size of the MMIO area (0x5000). In case a TPM2 is used by QEMU, a TPM2 ACPI table is also provided. The device is described to be used in polling mode rather than interrupt mode primarily because no unused IRQ could be found.

To support measurement logs to be written by the firmware, e.g. SeaBIOS, a TCPA table is implemented. This table provides a 64kb buffer where the firmware can write its log into. For TPM 2 only a more recent version of the TPM2 table provides support for measurements logs and a TCPA table does not need to be created.

The TCPA and TPM2 ACPI tables follow the Trusted Computing Group specification “TCG ACPI Specification” Family “1.2” and “2.0”, Level 00 Revision 00.37. (see the [ACPI specification](#), or a later version of it)

ACPI PPI Interface

QEMU supports the Physical Presence Interface (PPI) for TPM 1.2 and TPM 2. This interface requires ACPI and firmware support. (see the [PPI specification](#))

PPI enables a system administrator (root) to request a modification to the TPM upon reboot. The PPI specification defines the operation requests and the actions the firmware has to take. The system administrator passes the operation request number to the firmware through an ACPI interface which writes this number to a memory location that the firmware knows. Upon reboot, the firmware finds the number and sends commands to the TPM. The firmware writes the TPM result code and the operation request number to a memory location that ACPI can read from and pass the result on to the administrator.

The PPI specification defines a set of mandatory and optional operations for the firmware to implement. The ACPI interface also allows an administrator to list the supported operations. In QEMU the ACPI code is generated by QEMU, yet the firmware needs to implement support on a per-operations basis, and different firmwares may support a different subset. Therefore, QEMU introduces the virtual memory device for PPI where the firmware can indicate which operations it supports and ACPI can enable the ones that are supported and disable all others. This interface lies in main memory and has the following layout:

Field	Length	Off-set	Description
func	0x100	0x000	Firmware sets values for each supported operation. See defined values below.
ppin	0x1	0x100	SMI interrupt to use. Set by firmware. Not supported.
ppip	0x4	0x101	ACPI function index to pass to SMM code. Set by ACPI. Not supported.
pprp	0x4	0x105	Result of last executed operation. Set by firmware. See function index 5 for values.
pprq	0x4	0x109	Operation request number to execute. See ‘Physical Presence Interface Operation Summary’ tables in specs. Set by ACPI.
pprm	0x4	0x10d	Operation request optional parameter. Values depend on operation. Set by ACPI.
lppr	0x4	0x111	Last executed operation request number. Copied from pprq field by firmware.
fret	0x4	0x115	Result code from SMM function. Not supported.
res1	0x40	0x119	Reserved for future use
next_step	0x1	0x159	Operation to execute after reboot by firmware. Used by firmware.
movv	0x1	0x15a	Memory overwrite variable

The following values are supported for the `func` field. They correspond to the values used by ACPI function index 8.

Value	Description
0	Operation is not implemented.
1	Operation is only accessible through firmware.
2	Operation is blocked for OS by firmware configuration.
3	Operation is allowed and physically present user required.
4	Operation is allowed and physically present user is not required.

The location of the table is given by the `fw_cfg tpmppi_address` field. The PPI memory region size is 0x400 (`TPM_PPI_ADDR_SIZE`) to leave enough room for future updates.

QEMU files related to TPM ACPI tables:

- `hw/i386/acpi-build.c`

- `include/hw/acpi/tpm.h`

5.8.4 TPM backend devices

The TPM implementation is split into two parts, frontend and backend. The frontend part is the hardware interface, such as the TPM TIS interface described earlier, and the other part is the TPM backend interface. The backend interfaces implement the interaction with a TPM device, which may be a physical or an emulated device. The split between the front- and backend devices allows a frontend to be connected with any available backend. This enables the TIS interface to be used with the passthrough backend or the swtpm backend.

QEMU files related to TPM backends:

- `backends/tpm.c`
- `include/sysemu/tpm.h`
- `include/sysemu/tpm_backend.h`

The QEMU TPM passthrough device

In case QEMU is run on Linux as the host operating system it is possible to make the hardware TPM device available to a single QEMU guest. In this case the user must make sure that no other program is using the device, e.g., `/dev/tpm0`, before trying to start QEMU with it.

The passthrough driver uses the host's TPM device for sending TPM commands and receiving responses from. Besides that it accesses the TPM device's sysfs entry for support of command cancellation. Since none of the state of a hardware TPM can be migrated between hosts, virtual machine migration is disabled when the TPM passthrough driver is used.

Since the host's TPM device will already be initialized by the host's firmware, certain commands, e.g. `TPM_Startup()`, sent by the virtual firmware for device initialization, will fail. In this case the firmware should not use the TPM.

Sharing the device with the host is generally not a recommended usage scenario for a TPM device. The primary reason for this is that two operating systems can then access the device's single set of resources, such as platform configuration registers (PCRs). Applications or kernel security subsystems, such as the Linux Integrity Measurement Architecture (IMA), are not expecting to share PCRs.

QEMU files related to the TPM passthrough device:

- `backends/tpm/tpm_passthrough.c`
- `backends/tpm/tpm_util.c`
- `include/sysemu/tpm_util.h`

Command line to start QEMU with the TPM passthrough device using the host's hardware TPM `/dev/tpm0`:

```
qemu-system-x86_64 -display sdl -accel kvm \
-m 1024 -boot d -bios bios-256k.bin -boot menu=on \
-tpmdev passthrough,id=tpm0,path=/dev/tpm0 \
-device tpm-tis,tpmdev=tpm0 test.img
```

The following commands should result in similar output inside the VM with a Linux kernel that either has the TPM TIS driver built-in or available as a module:

```
# dmesg | grep -i tpm
[    0.711310] tpm_tis 00:06: 1.2 TPM (device=id 0x1, rev-id 1)

# dmesg | grep TCPA
[    0.000000] ACPI: TCPA 0x00000000003FFD191C 000032 (v02 BOCHS \
    BXPCTCPA 0000001 BXPC 00000001)

# ls -l /dev/tpm*
crw-----. 1 root root 10, 224 Jul 11 10:11 /dev/tpm0

# find /sys/devices/ | grep pcrcs$ | xargs cat
PCR-00: 35 4E 3B CE 23 9F 38 59 ...
...
PCR-23: 00 00 00 00 00 00 00 00 ...
```

The QEMU TPM emulator device

The TPM emulator device uses an external TPM emulator called ‘swtpm’ for sending TPM commands to and receiving responses from. The swtpm program must have been started before trying to access it through the TPM emulator with QEMU.

The TPM emulator implements a command channel for transferring TPM commands and responses as well as a control channel over which control commands can be sent. (see the [SWTPM protocol](#) specification)

The control channel serves the purpose of resetting, initializing, and migrating the TPM state, among other things.

The swtpm program behaves like a hardware TPM and therefore needs to be initialized by the firmware running inside the QEMU virtual machine. One necessary step for initializing the device is to send the TPM_Startup command to it. SeaBIOS, for example, has been instrumented to initialize a TPM 1.2 or TPM 2 device using this command.

QEMU files related to the TPM emulator device:

- backends/tpm/tpm_emulator.c
- backends/tpm/tpm_util.c
- include/sysemu/tpm_util.h

The following commands start the swtpm with a UnixIO control channel over a socket interface. They do not need to be run as root.

```
mkdir /tmp/mytpm1
swtpm socket --tpmstate dir=/tmp/mytpm1 \
  --ctrl type=unixio,path=/tmp/mytpm1/swtpm-sock \
  --log level=20
```

Command line to start QEMU with the TPM emulator device communicating with the swtpm (x86):

```
qemu-system-x86_64 -display sdl -accel kvm \
  -m 1024 -boot d -bios bios-256k.bin -boot menu=on \
  -chardev socket,id=chrtpm,path=/tmp/mytpm1/swtpm-sock \
  -tpmdev emulator,id=tpm0,chardev=chrtpm \
  -device tpm-tis,tpmdev=tpm0 test.img
```

In case a pSeries machine is emulated, use the following command line:

```
qemu-system-ppc64 -display sdl -machine pseries,accel=kvm \
  -m 1024 -bios slof.bin -boot menu=on \
```

(continues on next page)

(continued from previous page)

```
-nodefaults -device VGA -device pci-ohci -device usb-kbd \  
-chardev socket,id=chrtpm,path=/tmp/mytpm1/swtpm-sock \  
-tpmdev emulator,id=tpm0,chardev=chrtpm \  
-device tpm-spapr,tpmdev=tpm0 \  
-device spapr-vscsi,id=scsi0,reg=0x00002000 \  
-device virtio-blk-pci,scsi=off,bus=pci.0,addr=0x3,drive=drive-virtio-disk0,  
→id=virtio-disk0 \  
-drive file=test.img,format=raw,if=none,id=drive-virtio-disk0
```

In case an Arm virt machine is emulated, use the following command line:

```
qemu-system-aarch64 -machine virt,gic-version=3,accel=kvm \  
-cpu host -m 4G \  
-nographic -no-acpi \  
-chardev socket,id=chrtpm,path=/tmp/mytpm1/swtpm-sock \  
-tpmdev emulator,id=tpm0,chardev=chrtpm \  
-device tpm-tis-device,tpmdev=tpm0 \  
-device virtio-blk-pci,drive=drv0 \  
-drive format=qcow2,file=hda.qcow2,if=none,id=drv0 \  
-drive if=pflash,format=raw,file=flash0.img,readonly=on \  
-drive if=pflash,format=raw,file=flash1.img
```

In case SeaBIOS is used as firmware, it should show the TPM menu item after entering the menu with ‘ESC’.

```
Select boot device:  
1. DVD/CD [ata1-0: QEMU DVD-ROM ATAPI-4 DVD/CD]  
[...]  
5. Legacy option rom  
  
t. TPM Configuration
```

The following commands should result in similar output inside the VM with a Linux kernel that either has the TPM TIS driver built-in or available as a module:

```
# dmesg | grep -i tpm  
[    0.711310] tpm_tis 00:06: 1.2 TPM (device=id 0x1, rev-id 1)  
  
# dmesg | grep TCPA  
[    0.000000] ACPI: TCPA 0x0000000003FFD191C 000032 (v02 BOCHS \  
    BXPCTCPA 0000001 BXPC 00000001)  
  
# ls -l /dev/tpm*  
crw-----. 1 root root 10, 224 Jul 11 10:11 /dev/tpm0  
  
# find /sys/devices/ | grep pcrs$ | xargs cat  
PCR-00: 35 4E 3B CE 23 9F 38 59 ...  
...  
PCR-23: 00 00 00 00 00 00 00 00 ...
```

5.8.5 Migration with the TPM emulator

The TPM emulator supports the following types of virtual machine migration:

- VM save / restore (migration into a file)
- Network migration

- Snapshotting (migration into storage like QoW2 or QED)

The following command sequences can be used to test VM save / restore.

In a 1st terminal start an instance of a swtpm using the following command:

```
mkdir /tmp/mytpm1
swtpm socket --tpmstate dir=/tmp/mytpm1 \
  --ctrl type=unixio,path=/tmp/mytpm1/swtpm-sock \
  --log level=20 --tpm2
```

In a 2nd terminal start the VM:

```
qemu-system-x86_64 -display sdl -accel kvm \
  -m 1024 -boot d -bios bios-256k.bin -boot menu=on \
  -chardev socket,id=chrtpm,path=/tmp/mytpm1/swtpm-sock \
  -tpmdev emulator,id=tpm0,chardev=chrtpm \
  -device tpm-tis,tpmdev=tpm0 \
  -monitor stdio \
  test.img
```

Verify that the attached TPM is working as expected using applications inside the VM.

To store the state of the VM use the following command in the QEMU monitor in the 2nd terminal:

```
(qemu) migrate "exec:cat > testvm.bin"
(qemu) quit
```

At this point a file called `testvm.bin` should exist and the swtpm and QEMU processes should have ended.

To test ‘VM restore’ you have to start the swtpm with the same parameters as before. If previously a TPM 2 [`-tpm2`] was saved, `-tpm2` must now be passed again on the command line.

In the 1st terminal restart the swtpm with the same command line as before:

```
swtpm socket --tpmstate dir=/tmp/mytpm1 \
  --ctrl type=unixio,path=/tmp/mytpm1/swtpm-sock \
  --log level=20 --tpm2
```

In the 2nd terminal restore the state of the VM using the additional ‘-incoming’ option.

```
qemu-system-x86_64 -display sdl -accel kvm \
  -m 1024 -boot d -bios bios-256k.bin -boot menu=on \
  -chardev socket,id=chrtpm,path=/tmp/mytpm1/swtpm-sock \
  -tpmdev emulator,id=tpm0,chardev=chrtpm \
  -device tpm-tis,tpmdev=tpm0 \
  -incoming "exec:cat < testvm.bin" \
  test.img
```

Troubleshooting migration

There are several reasons why migration may fail. In case of problems, please ensure that the command lines adhere to the following rules and, if possible, that identical versions of QEMU and swtpm are used at all times.

VM save and restore:

- QEMU command line parameters should be identical apart from the ‘-incoming’ option on VM restore
- swtpm command line parameters should be identical

VM migration to ‘localhost’:

- QEMU command line parameters should be identical apart from the ‘-incoming’ option on the destination side
- swtpm command line parameters should point to two different directories on the source and destination swtpm (–tpmstate dir=...) (especially if different versions of libtpms were to be used on the same machine).

VM migration across the network:

- QEMU command line parameters should be identical apart from the ‘-incoming’ option on the destination side
- swtpm command line parameters should be identical

VM Snapshotting:

- QEMU command line parameters should be identical
- swtpm command line parameters should be identical

Besides that, migration failure reasons on the swtpm level may include the following:

- the versions of the swtpm on the source and destination sides are incompatible
 - downgrading of TPM state may not be supported
 - the source and destination libtpms were compiled with different compile-time options and the destination side refuses to accept the state
- different migration keys are used on the source and destination side and the destination side cannot decrypt the migrated state (swtpm ... –migration-key ...)

5.9 APEI tables generating and CPER record

5.9.1 Design Details

etc/acpi/tables				etc/hardware_errors			
=====				=====			
+ +-----+ +				+ +-----+ +			
HEST		+----->		error_block_address1		-----+	
+-----+ +				+-----+ +			
GHES1		+----->		error_block_address2		-----+ +	
+-----+ +				+-----+ +			
.....						
error_status_address-----+ +				+-----+ +			
.....		+-->		error_block_addressN		-----+ + + +	
read_ack_register-----+ +				+-----+ +			
read_ack_preserve		+ +----->		read_ack_register1			
read_ack_write				+-----+ +			
+ +-----+ +		+ +----->		read_ack_register2			
GHES2				+-----+ +			
+ +-----+ +						
.....				+-----+ +			
error_status_address-----+ +		+-->		read_ack_registerN			
.....				+-----+ +			
read_ack_register-----+ +				Generic Error Status Block 1 <-----+			
read_ack_preserve				+-----+ +			
read_ack_write							
+ +-----+ +				CPER			
.....				CPER			
.....						

(continues on next page)

(continued from previous page)

+ +-----+ CPER		
GHESN -+-----+		
+ +-----+ Generic Error Status Block 2 <-----+		
..... -+-----+ +		
error_status_address-----+-----+ CPER		
..... CPER		
read_ack_register-----+-----+ 		
read_ack_preserve CPER		
read_ack_write +-----+ +		
+ +-----+ 		
	-----+	
	Generic Error Status Block N <-----+	
	-+-----+ +	
	CPER	
	CPER	
	
	CPER	
	+-----+ +	

- (1) QEMU generates the ACPI HEST table. This table goes in the current “etc/acpi/tables” fw_cfg blob. Each error source has different notification types.
- (2) A new fw_cfg blob called “etc/hardware_errors” is introduced. QEMU also needs to populate this blob. The “etc/hardware_errors” fw_cfg blob contains an address registers table and an Error Status Data Block table.
- (3) The address registers table contains N Error Block Address entries and N Read Ack Register entries. The size for each entry is 8-byte. The Error Status Data Block table contains N Error Status Data Block entries. The size for each entry is 4096(0x1000) bytes. The total size for the “etc/hardware_errors” fw_cfg blob is $(N * 8 * 2 + N * 4096)$ bytes. N is the number of the kinds of hardware error sources.
- (4) QEMU generates the ACPI linker/loader script for the firmware. The firmware pre-allocates memory for “etc/acpi/tables”, “etc/hardware_errors” and copies blob contents there.
- (5) QEMU generates N ADD_POINTER commands, which patch addresses in the “error_status_address” fields of the HEST table with a pointer to the corresponding “address registers” in the “etc/hardware_errors” blob.
- (6) QEMU generates N ADD_POINTER commands, which patch addresses in the “read_ack_register” fields of the HEST table with a pointer to the corresponding “read_ack_register” within the “etc/hardware_errors” blob.
- (7) QEMU generates N ADD_POINTER commands for the firmware, which patch addresses in the “error_block_address” fields with a pointer to the respective “Error Status Data Block” in the “etc/hardware_errors” blob.
- (8) QEMU defines a third and write-only fw_cfg blob which is called “etc/hardware_errors_addr”. Through that blob, the firmware can send back the guest-side allocation addresses to QEMU. The “etc/hardware_errors_addr” blob contains a 8-byte entry. QEMU generates a single WRITE_POINTER command for the firmware. The firmware will write back the start address of “etc/hardware_errors” blob to the fw_cfg file “etc/hardware_errors_addr”.
- (9) When QEMU gets a SIGBUS from the kernel, QEMU writes CPER into corresponding “Error Status Data Block”, guest memory, and then injects platform specific interrupt (in case of arm/virt machine it’s Synchronous External Abort) as a notification which is necessary for notifying the guest.
- (10) This notification (in virtual hardware) will be handled by the guest kernel, on receiving notification, guest APEI driver could read the CPER error and take appropriate action.
- (11) kvm_arch_on_sigbus_vcpu() uses source_id as index in “etc/hardware_errors” to find out “Error Status Data Block” entry corresponding to error source. So supported source_id values should be assigned here and not

be changed afterwards to make sure that guest will write error into expected “Error Status Data Block” even if guest was migrated to a newer QEMU.

This manual documents various parts of the internals of QEMU. You only need to read it if you are interested in reading or modifying QEMU's source code.

Contents:

6.1 The QEMU build system architecture

This document aims to help developers understand the architecture of the QEMU build system. As with projects using GNU autotools, the QEMU build system has two stages, first the developer runs the “configure” script to determine the local build environment characteristics, then they run “make” to build the project. There is about where the similarities with GNU autotools end, so try to forget what you know about them.

6.1.1 Stage 1: configure

The QEMU configure script is written directly in shell, and should be compatible with any POSIX shell, hence it uses `#!/bin/sh`. An important implication of this is that it is important to avoid using bash-isms on development platforms where bash is the primary host.

In contrast to autoconf scripts, QEMU's configure is expected to be silent while it is checking for features. It will only display output when an error occurs, or to show the final feature enablement summary on completion.

Because QEMU uses the Meson build system under the hood, only VPATH builds are supported. There are two general ways to invoke configure & perform a build:

- VPATH, build artifacts outside of QEMU source tree entirely:

```
cd ../
mkdir build
cd build
../qemu/configure
make
```

- VPATH, build artifacts in a subdir of QEMU source tree:

```
mkdir build
cd build
../configure
make
```

For now, checks on the compilation environment are found in `configure` rather than `meson.build`, though this is expected to change. The command line is parsed in the `configure` script and, whenever needed, converted into the appropriate options to Meson.

New checks should be added to Meson, which usually comprises the following tasks:

- Add a Meson build option to `meson_options.txt`.
- Add support to the command line arg parser to handle any new `-enable-XXX/-disable-XXX` flags required by the feature.
- Add information to the help output message to report on the new feature flag.
- Add code to perform the actual feature check.
- Add code to include the feature status in `config-host.h`
- Add code to print out the feature status in the configure summary upon completion.

Taking the probe for `SDL2_Image` as an example, we have the following pieces in `configure`:

```
# Initial variable state
sdl_image=auto

..snip..

# Configure flag processing
--disable-sdl-image) sdl_image=disabled
;;
--enable-sdl-image) sdl_image=enabled
;;

..snip..

# Help output feature message
sdl-image          SDL Image support for icons

..snip..

# Meson invocation
-Dsdl_image=$sdl_image
```

In `meson_options.txt`:

```
option('sdl', type : 'feature', value : 'auto',
       description: 'SDL Image support for icons')
```

In `meson.build`:

```
# Detect dependency
sdl_image = dependency('SDL2_image', required: get_option('sdl_image'),
                      method: 'pkg-config',
                      kwargs: static_kwargs)
```

(continues on next page)

(continued from previous page)

```
# Create config-host.h (if applicable)
config_host_data.set('CONFIG_SDL_IMAGE', sdl_image.found())

# Summary
summary_info += {'SDL image support': sdl_image.found() }
```

Helper functions

The configure script provides a variety of helper functions to assist developers in checking for system features:

do_cc \$ARGS... Attempt to run the system C compiler passing it \$ARGS...

do_cxx \$ARGS... Attempt to run the system C++ compiler passing it \$ARGS...

compile_object \$CFLAGS Attempt to compile a test program with the system C compiler using \$CFLAGS. The test program must have been previously written to a file called \$TMPC. The replacement in Meson is the compiler object *cc*, which has methods such as *cc.compiles()*, *cc.check_header()*, *cc.has_function()*.

compile_prog \$CFLAGS \$LDFLAGS Attempt to compile a test program with the system C compiler using \$CFLAGS and link it with the system linker using \$LDFLAGS. The test program must have been previously written to a file called \$TMPC. The replacement in Meson is *cc.find_library()* and *cc.links()*.

has \$COMMAND Determine if \$COMMAND exists in the current environment, either as a shell builtin, or executable binary, returning 0 on success. The replacement in Meson is *find_program()*.

check_define \$NAME Determine if the macro \$NAME is defined by the system C compiler

check_include \$NAME Determine if the include \$NAME file is available to the system C compiler. The replacement in Meson is *cc.has_header()*.

write_c_skeleton Write a minimal C program main() function to the temporary file indicated by \$TMPC

feature_not_found \$NAME \$REMEDY Print a message to stderr that the feature \$NAME was not available on the system, suggesting the user try \$REMEDY to address the problem.

error_exit \$MESSAGE \$MORE... Print \$MESSAGE to stderr, followed by \$MORE... and then exit from the configure script with non-zero status

query_pkg_config \$ARGS... Run pkg-config passing it \$ARGS. If QEMU is doing a static build, then `--static` will be automatically added to \$ARGS

6.1.2 Stage 2: Meson

The Meson build system is currently used to describe the build process for:

- 1) executables, which include:
 - Tools - qemu-img, qemu-nbd, qga (guest agent), etc
 - System emulators - qemu-system-\$ARCH
 - Userspace emulators - qemu-\$ARCH
 - Unit tests
- 2) documentation
- 3) ROMs, which can be either installed as binary blobs or compiled
- 4) other data files, such as icons or desktop files

All executables are built by default, except for some *contrib/* binaries that are known to fail to build on some platforms (for example 32-bit or big-endian platforms). Tests are also built by default, though that might change in the future.

The source code is highly modularized, split across many files to facilitate building of all of these components with as little duplicated compilation as possible. Using the Meson “sourceset” functionality, *meson.build* files group the source files in rules that are enabled according to the available system libraries and to various configuration symbols. Sourcesets belong to one of four groups:

Subsystem sourcesets: Various subsystems that are common to both tools and emulators have their own sourceset, for example *block_ss* for the block device subsystem, *chardev_ss* for the character device subsystem, etc. These sourcesets are then turned into static libraries as follows:

```
libchardev = static_library('chardev', chardev_ss.sources(),
                           name_suffix: 'fa',
                           build_by_default: false)

chardev = declare_dependency(link_whole: libchardev)
```

As of Meson 0.55.1, the special *.fa* suffix should be used for everything that is used with *link_whole*, to ensure that the link flags are placed correctly in the command line.

Target-independent emulator sourcesets: Various general purpose helper code is compiled only once and the *.o* files are linked into all output binaries that need it. This includes error handling infrastructure, standard data structures, platform portability wrapper functions, etc.

Target-independent code lives in the *common_ss*, *softmmu_ss* and *user_ss* sourcesets. *common_ss* is linked into all emulators, *softmmu_ss* only in system emulators, *user_ss* only in user-mode emulators.

Target-independent sourcesets must exercise particular care when using *if_false* rules. The *if_false* rule will be used correctly when linking emulator binaries; however, when *compiling* target-independent files into *.o* files, Meson may need to pick *both* the *if_true* and *if_false* sides to cater for targets that want either side. To achieve that, you can add a special rule using the *CONFIG_ALL* symbol:

```
# Some targets have CONFIG_ACPI, some don't, so this is not enough
softmmu_ss.add(when: 'CONFIG_ACPI', if_true: files('acpi.c'),
               if_false: files('acpi-stub.c'))

# This is required as well:
softmmu_ss.add(when: 'CONFIG_ALL', if_true: files('acpi-stub.c'))
```

Target-dependent emulator sourcesets: In the target-dependent set lives CPU emulation, some device emulation and much glue code. This sometimes also has to be compiled multiple times, once for each target being built. Target-dependent files are included in the *specific_ss* sourceset.

Each emulator also includes sources for files in the *hw/* and *target/* subdirectories. The subdirectory used for each emulator comes from the target’s definition of *TARGET_BASE_ARCH* or (if missing) *TARGET_ARCH*, as found in *default-configs/targets/*.mak*.

Each subdirectory in *hw/* adds one sourceset to the *hw_arch* dictionary, for example:

```
arm_ss = ss.source_set()
arm_ss.add(files('boot.c'), fdt)
...
hw_arch += {'arm': arm_ss}
```

The sourceset is only used for system emulators.

Each subdirectory in *target/* instead should add one sourceset to each of the *target_arch* and *target_softmmu_arch*, which are used respectively for all emulators and for system emulators only. For example:

```
arm_ss = ss.source_set()
arm_softmmu_ss = ss.source_set()
...
target_arch += {'arm': arm_ss}
target_softmmu_arch += {'arm': arm_softmmu_ss}
```

Utility source sets: All binaries link with a static library *libqemuutil.a*. This library is built from several source sets; most of them however host generated code, and the only two of general interest are *util_ss* and *stub_ss*.

The separation between these two is purely for documentation purposes. *util_ss* contains generic utility files. Even though this code is only linked in some binaries, sometimes it requires hooks only in some of these and depend on other functions that are not fully implemented by all QEMU binaries. *stub_ss* links dummy stubs that will only be linked into the binary if the real implementation is not present. In a way, the stubs can be thought of as a portable implementation of the weak symbols concept.

The following files concur in the definition of which files are linked into each emulator:

default-configs/devices/*.mak The files under *default-configs/devices/* control the boards and devices that are built into each QEMU system emulation targets. They merely contain a list of config variable definitions such as:

```
include arm-softmmu.mak
CONFIG_XLNX_ZYNQMP_ARM=y
CONFIG_XLNX_VERSAL=y
```

****/Kconfig*** These files are processed together with *default-configs/devices/*.mak* and describe the dependencies between various features, subsystems and device models. They are described in [QEMU and Kconfig](#)

default-configs/targets/*.mak These files mostly define symbols that appear in the **-config-target.h* file for each emulator¹. However, the `TARGET_ARCH` and `TARGET_BASE_ARCH` will also be used to select the *hw/* and *target/* subdirectories that are compiled into each target.

These files rarely need changing unless you are adding a completely new target, or enabling new devices or hardware for a particular system/userspace emulation target

Support scripts

Meson has a special convention for invoking Python scripts: if their first line is `#!/usr/bin/env python3` and the file is *not* executable, `find_program()` arranges to invoke the script under the same Python interpreter that was used to invoke Meson. This is the most common and preferred way to invoke support scripts from Meson build files, because it automatically uses the value of `configure's --python=` option.

In case the script is not written in Python, use a `#!/usr/bin/env ...` line and make the script executable.

Scripts written in Python, where it is desirable to make the script executable (for example for test scripts that developers may want to invoke from the command line, such as `tests/qapi-schema/test-qapi.py`), should be invoked through the *python* variable in `meson.build`. For example:

```
test('QAPI schema regression tests', python,
     args: files('test-qapi.py'),
     env: test_env, suite: ['qapi-schema', 'qapi-frontend'])
```

This is needed to obey the `--python=` option passed to the `configure` script, which may point to something other than the first `python3` binary on the path.

¹ This header is included by *qemu/osdep.h* when compiling files from the target-specific source sets.

6.1.3 Stage 3: makefiles

The use of GNU make is required with the QEMU build system.

The output of Meson is a build.ninja file, which is used with the Ninja build system. QEMU uses a different approach, where Makefile rules are synthesized from the build.ninja file. The main Makefile includes these rules and wraps them so that e.g. submodules are built before QEMU. The resulting build system is largely non-recursive in nature, in contrast to common practices seen with automake.

Tests are also ran by the Makefile with the traditional *make check* phony target, while benchmarks are run with *make bench*. Meson test suites such as *unit* can be ran with *make check-unit* too. It is also possible to run tests defined in meson.build with *meson test*.

6.1.4 Important files for the build system

Statically defined files

The following key files are statically defined in the source tree, with the rules needed to build QEMU. Their behaviour is influenced by a number of dynamically created files listed later.

Makefile The main entry point used when invoking make to build all the components of QEMU. The default ‘all’ target will naturally result in the build of every component. Makefile takes care of recursively building submodules directly via a non-recursive set of rules.

***/meson.build** The meson.build file in the root directory is the main entry point for the Meson build system, and it coordinates the configuration and build of all executables. Build rules for various subdirectories are included in other meson.build files spread throughout the QEMU source tree.

tests/Makefile.include Rules for external test harnesses. These include the TCG tests, *qemu-iotests* and the Avocado-based acceptance tests.

tests/docker/Makefile.include Rules for Docker tests. Like tests/Makefile, this file is included directly by the top level Makefile, anything defined in this file will influence the entire build system.

tests/vm/Makefile.include Rules for VM-based tests. Like tests/Makefile, this file is included directly by the top level Makefile, anything defined in this file will influence the entire build system.

Dynamically created files

The following files are generated dynamically by configure in order to control the behaviour of the statically defined makefiles. This avoids the need for QEMU makefiles to go through any pre-processing as seen with autotools, where Makefile.am generates Makefile.in which generates Makefile.

Built by configure:

config-host.mak When configure has determined the characteristics of the build host it will write a long list of variables to config-host.mak file. This provides the various install directories, compiler / linker flags and a variety of *CONFIG_** variables related to optionally enabled features. This is imported by the top level Makefile and meson.build in order to tailor the build output.

config-host.mak is also used as a dependency checking mechanism. If make sees that the modification timestamp on configure is newer than that on config-host.mak, then configure will be re-run.

The variables defined here are those which are applicable to all QEMU build outputs. Variables which are potentially different for each emulator target are defined by the next file...

\$TARGET-NAME/config-target.mak TARGET-NAME is the name of a system or userspace emulator, for example, x86_64-softhmmu denotes the system emulator for the x86_64 architecture. This file contains the variables which

need to vary on a per-target basis. For example, it will indicate whether KVM or Xen are enabled for the target and any other potential custom libraries needed for linking the target.

Built by Meson:

\${TARGET-NAME}-config-devices.mak *TARGET-NAME* is again the name of a system or userspace emulator. The *config-devices.mak* file is automatically generated by make using the *scripts/make_device_config.sh* program, feeding it the *default-configs/\${TARGET-NAME}* file as input.

config-host.h, \${TARGET-NAME}/config-target.h, \${TARGET-NAME}/config-devices.h These files are used by source code to determine what features are enabled. They are generated from the contents of the corresponding **.h* files using the *scripts/create_config* program. This extracts relevant variables and formats them as C preprocessor macros.

build.ninja The build rules.

Built by Makefile:

Makefile.ninja A Makefile include that bridges to *ninja* for the actual build. The Makefile is mostly a list of targets that Meson included in *build.ninja*.

Makefile.mtest The Makefile definitions that let “make check” run tests defined in *meson.build*. The rules are produced from Meson’s JSON description of tests (obtained with “meson introspect –tests”) through the script *scripts/mtest2make.py*.

Useful make targets

help Print a help message for the most common build targets.

print-VAR Print the value of the variable *VAR*. Useful for debugging the build system.

6.2 QEMU and Kconfig

QEMU is a very versatile emulator; it can be built for a variety of targets, where each target can emulate various boards and at the same time different targets can share large amounts of code. For example, a POWER and an x86 board can run the same code to emulate a PCI network card, even though the boards use different PCI host bridges, and they can run the same code to emulate a SCSI disk while using different SCSI adapters. Arm, s390 and x86 boards can all present a virtio-blk disk to their guests, but with three different virtio guest interfaces.

Each QEMU target enables a subset of the boards, devices and buses that are included in QEMU’s source code. As a result, each QEMU executable only links a small subset of the files that form QEMU’s source code; anything that is not needed to support a particular target is culled.

QEMU uses a simple domain-specific language to describe the dependencies between components. This is useful for two reasons:

- new targets and boards can be added without knowing in detail the architecture of the hardware emulation subsystems. Boards only have to list the components they need, and the compiled executable will include all the required dependencies and all the devices that the user can add to that board;
- users can easily build reduced versions of QEMU that support only a subset of boards or devices. For example, by default most targets will include all emulated PCI devices that QEMU supports, but the build process is configurable and it is easy to drop unnecessary (or otherwise unwanted) code to make a leaner binary.

This domain-specific language is based on the Kconfig language that originated in the Linux kernel, though it was heavily simplified and the handling of dependencies is stricter in QEMU.

Unlike Linux, there is no user interface to edit the configuration, which is instead specified in per-target files under the *default-configs/* directory of the QEMU source tree. This is because, unlike Linux, configuration and

dependencies can be treated as a black box when building QEMU; the default configuration that QEMU ships with should be okay in almost all cases.

6.2.1 The Kconfig language

Kconfig defines configurable components in files named `hw/*/Kconfig`. Note that configurable components are `_not_` visible in C code as preprocessor symbols; they are only visible in the Makefile. Each configurable component defines a Makefile variable whose name starts with `CONFIG_`.

All elements have boolean (true/false) type; truth is written as `y`, while falsehood is written `n`. They are defined in a Kconfig stanza like the following:

```
config ARM_VIRT
    bool
    imply PCI_DEVICES
    imply VFIO_AMD_XGBE
    imply VFIO_XGMAC
    select A15MPCORE
    select ACPI
    select ARM_SMMUV3
```

The `config` keyword introduces a new configuration element. In the example above, Makefiles will have access to a variable named `CONFIG_ARM_VIRT`, with value `y` or `n` (respectively for boolean true and false).

Boolean expressions can be used within the language, whenever `<expr>` is written in the remainder of this section. The `&&`, `||` and `!` operators respectively denote conjunction (AND), disjunction (OR) and negation (NOT).

The `bool` data type declaration is optional, but it is suggested to include it for clarity and future-proofing. After `bool` the following directives can be included:

dependencies: `depends on <expr>`

This defines a dependency for this configurable element. Dependencies evaluate an expression and force the value of the variable to false if the expression is false.

reverse dependencies: `select <symbol> [if <expr>]`

While `depends on` can force a symbol to false, reverse dependencies can be used to force another symbol to true. In the following example, `CONFIG_BAZ` will be true whenever `CONFIG_FOO` is true:

```
config FOO
    select BAZ
```

The optional expression will prevent `select` from having any effect unless it is true.

Note that unlike Linux's Kconfig implementation, QEMU will detect contradictions between `depends on` and `select` statements and prevent you from building such a configuration.

default value: `default <value> [if <expr>]`

Default values are assigned to the config symbol if no other value was set by the user via `default-configs/*.mak` files, and only if `select` or `depends on` directives do not force the value to true or false respectively. `<value>` can be `y` or `n`; it cannot be an arbitrary Boolean expression. However, a condition for applying the default value can be added with `if`.

A configuration element can have any number of default values (usually, if more than one default is present, they will have different conditions). If multiple default values satisfy their condition, only the first defined one is active.

reverse default (weak reverse dependency): `imply <symbol> [if <expr>]`

This is similar to `select` as it applies a lower limit of `y` to another symbol. However, the lower limit is only a default and the “implied” symbol’s value may still be set to `n` from a `default-configs/*.mak` files. The following two examples are equivalent:

```
config FOO
    bool
    imply BAZ

config BAZ
    bool
    default y if FOO
```

The next section explains where to use `imply` or `default y`.

6.2.2 Guidelines for writing Kconfig files

Configurable elements in QEMU fall under five broad groups. Each group declares its dependencies in different ways:

subsystems, of which **buses** are a special case

Example:

```
config SCSI
    bool
```

Subsystems always default to false (they have no `default` directive) and are never visible in `default-configs/*.mak` files. It’s up to other symbols to `select` whatever subsystems they require.

They sometimes have `select` directives to bring in other required subsystems or buses. For example, `AUX` (the DisplayPort auxiliary channel “bus”) `selects` `I2C` because it can act as an `I2C` master too.

devices

Example:

```
config MEGASAS_SCSI_PCI
    bool
    default y if PCI_DEVICES
    depends on PCI
    select SCSI
```

Devices are the most complex of the five. They can have a variety of directives that cooperate so that a default configuration includes all the devices that can be accessed from QEMU.

Devices *depend on* the bus that they lie on, for example a PCI device would specify `depends on PCI`. An MMIO device will likely have no `depends on` directive. Devices also *select* the buses that the device provides, for example a SCSI adapter would specify `select SCSI`. Finally, devices are usually `default y` if and only if they have at least one `depends on`; the default could be conditional on a device group.

Devices also select any optional subsystem that they use; for example a video card might specify `select EDID` if it needs to build EDID information and publish it to the guest.

device groups

Example:

```
config PCI_DEVICES
    bool
```

Device groups provide a convenient mechanism to enable/disable many devices in one go. This is useful when a set of devices is likely to be enabled/disabled by several targets. Device groups usually need no directive and are not used in the Makefile either; they only appear as conditions for `default y` directives.

QEMU currently has two device groups, `PCI_DEVICES` and `TEST_DEVICES`. PCI devices usually have a `default y` if `PCI_DEVICES` directive rather than just `default y`. This lets some boards (notably s390) easily support a subset of PCI devices, for example only VFIO (passthrough) and virtio-pci devices. `TEST_DEVICES` instead is used for devices that are rarely used on production virtual machines, but provide useful hooks to test QEMU or KVM.

boards

Example:

```
config SUN4M
    bool
    imply TCX
    imply CG3
    select CS4231
    select ECCMEMCTL
    select EMPTY_SLOT
    select ESCC
    select ESP
    select FDC
    select SLAVIO
    select LANCE
    select M48T59
    select STP2000
```

Boards specify their constituent devices using `imply` and `select` directives. A device should be listed under `select` if the board cannot be started at all without it. It should be listed under `imply` if (depending on the QEMU command line) the board may or may not be started without it. Boards also default to false; they are enabled by the `default-configs/*.mak` for the target they apply to.

internal elements

Example:

```
config ECCMEMCTL
    bool
    select ECC
```

Internal elements group code that is useful in several boards or devices. They are usually enabled with `select` and in turn select other elements; they are never visible in `default-configs/*.mak` files, and often not even in the Makefile.

6.2.3 Writing and modifying default configurations

In addition to the Kconfig files under `hw/`, each target also includes a file called `default-configs/TARGETNAME-softmmu.mak`. These files initialize some Kconfig variables to non-default values and provide the starting point to turn on devices and subsystems.

A file in `default-configs/` looks like the following example:

```
# Default configuration for alpha-softmmu

# Uncomment the following lines to disable these optional devices:
#
#CONFIG_PCI_DEVICES=n
#CONFIG_TEST_DEVICES=n

# Boards:
#
CONFIG_DP264=y
```

The first part, consisting of commented-out `=n` assignments, tells the user which devices or device groups are implied by the boards. The second part, consisting of `=y` assignments, tells the user which boards are supported by the target. The user will typically modify the default configuration by uncommenting lines in the first group, or commenting out lines in the second group.

It is also possible to run QEMU's configure script with the `--without-default-devices` option. When this is done, everything defaults to `n` unless it is select`ed or explicitly switched on in the ``.mak files. In other words, default and imply directives are disabled. When QEMU is built with this option, the user will probably want to change some lines in the first group, for example like this:

```
CONFIG_PCI_DEVICES=y
#CONFIG_TEST_DEVICES=n
```

and/or pick a subset of the devices in those device groups. Right now there is no single place that lists all the optional devices for `CONFIG_PCI_DEVICES` and `CONFIG_TEST_DEVICES`. In the future, we expect that .mak files will be automatically generated, so that they will include all these symbols and some help text on what they do.

6.2.4 Kconfig.host

In some special cases, a configurable element depends on host features that are detected by QEMU's configure or meson.build scripts; for example some devices depend on the availability of KVM or on the presence of a library on the host.

These symbols should be listed in `Kconfig.host` like this:

```
config TPM
    bool
```

and also listed as follows in the top-level meson.build's `host_kconfig` variable:

```
host_kconfig = \
    ('CONFIG_TPM' in config_host ? ['CONFIG_TPM=y'] : []) + \
    ('CONFIG_SPICE' in config_host ? ['CONFIG_SPICE=y'] : []) + \
    ('CONFIG_IVSHMEM' in config_host ? ['CONFIG_IVSHMEM=y'] : []) + \
    ...
```

6.3 Testing in QEMU

This document describes the testing infrastructure in QEMU.

6.3.1 Testing with “make check”

The “make check” testing family includes most of the C based tests in QEMU. For a quick help, run `make check-help` from the source tree.

The usual way to run these tests is:

```
make check
```

which includes QAPI schema tests, unit tests, QTestS and some iotests. Different sub-types of “make check” tests will be explained below.

Before running tests, it is best to build QEMU programs first. Some tests expect the executables to exist and will fail with obscure messages if they cannot find them.

Unit tests

Unit tests, which can be invoked with `make check-unit`, are simple C tests that typically link to individual QEMU object files and exercise them by calling exported functions.

If you are writing new code in QEMU, consider adding a unit test, especially for utility modules that are relatively stateless or have few dependencies. To add a new unit test:

1. Create a new source file. For example, `tests/foo-test.c`.
2. Write the test. Normally you would include the header file which exports the module API, then verify the interface behaves as expected from your test. The test code should be organized with the glib testing framework. Copying and modifying an existing test is usually a good idea.
3. Add the test to `tests/meson.build`. The unit tests are listed in a dictionary called `tests`. The values are any additional sources and dependencies to be linked with the test. For a simple test whose source is in `tests/foo-test.c`, it is enough to add an entry like:

```
{
    ...
    'foo-test': [],
    ...
}
```

Since unit tests don’t require environment variables, the simplest way to debug a unit test failure is often directly invoking it or even running it under `gdb`. However there can still be differences in behavior between `make` invocations and your manual run, due to `$MALLOCPERTURB` environment variable (which affects memory reclamation and catches invalid pointers better) and `gtester` options. If necessary, you can run

```
make check-unit V=1
```

and copy the actual command line which executes the unit test, then run it from the command line.

QTest

QTest is a device emulation testing framework. It can be very useful to test device models; it could also control certain aspects of QEMU (such as virtual clock stepping), with a special purpose “qtest” protocol. Refer to [QTest Device Emulation Testing Framework](#) for more details.

QTest cases can be executed with

```
make check-qtest
```

QAPI schema tests

The QAPI schema tests validate the QAPI parser used by QMP, by feeding predefined input to the parser and comparing the result with the reference output.

The input/output data is managed under the `tests/qapi-schema` directory. Each test case includes four files that have a common base name:

- `${casename}.json` - the file contains the JSON input for feeding the parser
- `${casename}.out` - the file contains the expected stdout from the parser
- `${casename}.err` - the file contains the expected stderr from the parser
- `${casename}.exit` - the expected error code

Consider adding a new QAPI schema test when you are making a change on the QAPI parser (either fixing a bug or extending/modifying the syntax). To do this:

1. Add four files for the new case as explained above. For example:

```
$EDITOR tests/qapi-schema/foo.{json,out,err,exit}.
```

2. Add the new test in `tests/Makefile.include`. For example:

```
qapi-schema += foo.json
```

check-block

`make check-block` runs a subset of the block layer iotests (the tests that are in the “auto” group). See the “QEMU iotests” section below for more information.

GCC gcov support

`gcov` is a GCC tool to analyze the testing coverage by instrumenting the tested code. To use it, configure QEMU with `--enable-gcov` option and build. Then run `make check` as usual.

If you want to gather coverage information on a single test the `make clean-gcda` target can be used to delete any existing coverage information before running a single test.

You can generate a HTML coverage report by executing `make coverage-html` which will create `meson-logs/coveragereport/index.html`.

Further analysis can be conducted by running the `gcov` command directly on the various `.gcda` output files. Please read the `gcov` documentation for more information.

6.3.2 QEMU iotests

QEMU iotests, under the directory `tests/qemu-iotests`, is the testing framework widely used to test block layer related features. It is higher level than “make check” tests and 99% of the code is written in bash or Python scripts. The testing success criteria is golden output comparison, and the test files are named with numbers.

To run iotests, make sure QEMU is built successfully, then switch to the `tests/qemu-iotests` directory under the build directory, and run `./check` with desired arguments from there.

By default, “raw” format and “file” protocol is used; all tests will be executed, except the unsupported ones. You can override the format and protocol with arguments:

```
# test with qcow2 format
./check -qcow2
# or test a different protocol
./check -nbd
```

It’s also possible to list test numbers explicitly:

```
# run selected cases with qcow2 format
./check -qcow2 001 030 153
```

Cache mode can be selected with the “-c” option, which may help reveal bugs that are specific to certain cache mode.

More options are supported by the `./check` script, run `./check -h` for help.

Writing a new test case

Consider writing a tests case when you are making any changes to the block layer. An iotest case is usually the choice for that. There are already many test cases, so it is possible that extending one of them may achieve the goal and save the boilerplate to create one. (Unfortunately, there isn’t a 100% reliable way to find a related one out of hundreds of tests. One approach is using `git grep`.)

Usually an iotest case consists of two files. One is an executable that produces output to stdout and stderr, the other is the expected reference output. They are given the same number in file names. E.g. Test script `055` and reference output `055.out`.

In rare cases, when outputs differ between cache mode `none` and others, a `.out.nocache` file is added. In other cases, when outputs differ between image formats, more than one `.out` files are created ending with the respective format names, e.g. `178.out.qcow2` and `178.out.raw`.

There isn’t a hard rule about how to write a test script, but a new test is usually a (copy and) modification of an existing case. There are a few commonly used ways to create a test:

- A Bash script. It will make use of several environmental variables related to the testing procedure, and could source a group of `common.*` libraries for some common helper routines.
- A Python unittest script. Import `iotests` and create a subclass of `iotests.QMPTestCase`, then call `iotests.main` method. The downside of this approach is that the output is too scarce, and the script is considered harder to debug.
- A simple Python script without using unittest module. This could also import `iotests` for launching QEMU and utilities etc, but it doesn’t inherit from `iotests.QMPTestCase` therefore doesn’t use the Python unittest execution. This is a combination of 1 and 2.

Pick the language per your preference since both Bash and Python have comparable library support for invoking and interacting with QEMU programs. If you opt for Python, it is strongly recommended to write Python 3 compatible code.

Both Python and Bash frameworks in `iotests` provide helpers to manage test images. They can be used to create and clean up images under the test directory. If no I/O or any protocol specific feature is needed, it is often more convenient to use the pseudo block driver, `null-co://`, as the test image, which doesn’t require image creation or cleaning up. Avoid system-wide devices or files whenever possible, such as `/dev/null` or `/dev/zero`. Otherwise, image locking implications have to be considered. For example, another application on the host may have locked the file, possibly leading to a test failure. If using such devices are explicitly desired, consider adding `locking=off` option to disable image locking.

Test case groups

“Tests may belong to one or more test groups, which are defined in the form of a comment in the test source file. By convention, test groups are listed in the second line of the test file, after the “#!/...” line, like this:

```
#!/usr/bin/env python3
# group: auto quick
#
...
```

Another way of defining groups is creating the tests/qemu-iotests/group.local file. This should be used only for downstream (this file should never appear in upstream). This file may be used for defining some downstream test groups or for temporarily disabling tests, like this:

```
# groups for some company downstream process
#
# ci - tests to run on build
# down - our downstream tests, not for upstream
#
# Format of each line is:
# TEST_NAME TEST_GROUP [TEST_GROUP ]...

013 ci
210 disabled
215 disabled
our-ugly-workaround-test down ci
```

Note that the following group names have a special meaning:

- quick: Tests in this group should finish within a few seconds.
- auto: Tests in this group are used during “make check” and should be runnable in any case. That means they should run with every QEMU binary (also non-x86), with every QEMU configuration (i.e. must not fail if an optional feature is not compiled in - but reporting a “skip” is ok), work at least with the qcow2 file format, work with all kind of host filesystems and users (e.g. “nobody” or “root”) and must not take too much memory and disk space (since CI pipelines tend to fail otherwise).
- disabled: Tests in this group are disabled and ignored by check.

6.3.3 Docker based tests

Introduction

The Docker testing framework in QEMU utilizes public Docker images to build and test QEMU in predefined and widely accessible Linux environments. This makes it possible to expand the test coverage across distros, toolchain flavors and library versions.

Prerequisites

Install “docker” with the system package manager and start the Docker service on your development machine, then make sure you have the privilege to run Docker commands. Typically it means setting up passwordless `sudo docker` command or login as root. For example:

```
$ sudo yum install docker
$ # or `apt-get install docker` for Ubuntu, etc.
$ sudo systemctl start docker
$ sudo docker ps
```

The last command should print an empty table, to verify the system is ready.

An alternative method to set up permissions is by adding the current user to “docker” group and making the docker daemon socket file (by default `/var/run/docker.sock`) accessible to the group:

```
$ sudo groupadd docker
$ sudo usermod $USER -a -G docker
$ sudo chown :docker /var/run/docker.sock
```

Note that any one of above configurations makes it possible for the user to exploit the whole host with Docker bind mounting or other privileged operations. So only do it on development machines.

Quickstart

From source tree, type `make docker` to see the help. Testing can be started without configuring or building QEMU (configure and make are done in the container, with parameters defined by the make target):

```
make docker-test-build@min-glib
```

This will create a container instance using the `min-glib` image (the image is downloaded and initialized automatically), in which the `test-build` job is executed.

Images

Along with many other images, the `min-glib` image is defined in a Dockerfile in `tests/docker/dockerfiles/`, called `min-glib.docker`. `make docker` command will list all the available images.

To add a new image, simply create a new `.docker` file under the `tests/docker/dockerfiles/` directory.

A `.pre` script can be added beside the `.docker` file, which will be executed before building the image under the build context directory. This is mainly used to do necessary host side setup. One such setup is `binfmt_misc`, for example, to make `qemu-user` powered cross build containers work.

Tests

Different tests are added to cover various configurations to build and test QEMU. Docker tests are the executables under `tests/docker` named `test-*`. They are typically shell scripts and are built on top of a shell library, `tests/docker/common.rc`, which provides helpers to find the QEMU source and build it.

The full list of tests is printed in the `make docker help`.

Debugging a Docker test failure

When CI tasks, maintainers or yourself report a Docker test failure, follow the below steps to debug it:

1. Locally reproduce the failure with the reported command line. E.g. `run make docker-test-mingw@fedora J=8`.
2. Add “V=1” to the command line, try again, to see the verbose output.

3. Further add “DEBUG=1” to the command line. This will pause in a shell prompt in the container right before testing starts. You could either manually build QEMU and run tests from there, or press Ctrl-D to let the Docker testing continue.
4. If you press Ctrl-D, the same building and testing procedure will begin, and will hopefully run into the error again. After that, you will be dropped to the prompt for debug.

Options

Various options can be used to affect how Docker tests are done. The full list is in the `make docker help` text. The frequently used ones are:

- `V=1`: the same as in top level `make`. It will be propagated to the container and enable verbose output.
- `J=$N`: the number of parallel tasks in `make` commands in the container, similar to the `-j $N` option in top level `make`. (The `-j` option in top level `make` will not be propagated into the container.)
- `DEBUG=1`: enables debug. See the previous “Debugging a Docker test failure” section.

6.3.4 Thread Sanitizer

Thread Sanitizer (TSan) is a tool which can detect data races. QEMU supports building and testing with this tool.

For more information on TSan:

<https://github.com/google/sanitizers/wiki/ThreadSanitizerCppManual>

Thread Sanitizer in Docker

TSan is currently supported in the `ubuntu2004` docker.

The `test-tsan` test will build using TSan and then run `make check`.

```
make docker-test-tsan@ubuntu2004
```

TSan warnings under docker are placed in files located at `build/tsan/`.

We recommend using `DEBUG=1` to allow launching the test from inside the docker, and to allow review of the warnings generated by TSan.

Building and Testing with TSan

It is possible to build and test with TSan, with a few additional steps. These steps are normally done automatically in the docker.

There is a one time patch needed in clang-9 or clang-10 at this time:

```
sed -i 's/^const/static const/g' \
    /usr/lib/llvm-10/lib/clang/10.0.0/include/sanitizer/tsan_interface.h
```

To configure the build for TSan:

```
../configure --enable-tsan --cc=clang-10 --cxx=clang++-10 \
    --disable-werror --extra-cflags="-O0"
```

The runtime behavior of TSAN is controlled by the `TSAN_OPTIONS` environment variable.

More information on the `TSAN_OPTIONS` can be found here:

<https://github.com/google/sanitizers/wiki/ThreadSanitizerFlags>

For example:

```
export TSAN_OPTIONS=suppressions=<path to qemu>/tests/tsan/suppressions.tsan \
        detect_deadlocks=false history_size=7 exitcode=0 \
        log_path=<build path>/tsan/tsan_warning
```

The above `exitcode=0` has TSan continue without error if any warnings are found. This allows for running the test and then checking the warnings afterwards. If you want TSan to stop and exit with error on warnings, use `exitcode=66`.

TSan Suppressions

Keep in mind that for any data race warning, although there might be a data race detected by TSan, there might be no actual bug here. TSan provides several different mechanisms for suppressing warnings. In general it is recommended to fix the code if possible to eliminate the data race rather than suppress the warning.

A few important files for suppressing warnings are:

`tests/tsan/suppressions.tsan` - Has TSan warnings we wish to suppress at runtime. The comment on each suppression will typically indicate why we are suppressing it. More information on the file format can be found here:

<https://github.com/google/sanitizers/wiki/ThreadSanitizerSuppressions>

`tests/tsan/blacklist.tsan` - Has TSan warnings we wish to disable at compile time for test or debug. Add flags to configure to enable:

“`-extra-cflags=-fsanitize-blacklist=<src path>/tests/tsan/blacklist.tsan`”

More information on the file format can be found here under “Blacklist Format”:

<https://github.com/google/sanitizers/wiki/ThreadSanitizerFlags>

TSan Annotations

`include/qemu/tsan.h` defines annotations. See this file for more descriptions of the annotations themselves. Annotations can be used to suppress TSan warnings or give TSan more information so that it can detect proper relationships between accesses of data.

Annotation examples can be found here:

<https://github.com/llvm/llvm-project/tree/master/compiler-rt/test/tsan/>

Good files to start with are: `annotate_happens_before.cpp` and `ignore_race.cpp`

The full set of annotations can be found here:

https://github.com/llvm/llvm-project/blob/master/compiler-rt/lib/tsan/rtl/tsan_interface_ann.cpp

6.3.5 VM testing

This test suite contains scripts that bootstrap various guest images that have necessary packages to build QEMU. The basic usage is documented in `Makefile` help which is displayed with `make vm-help`.

Quickstart

Run `make vm-help` to list available make targets. Invoke a specific make command to run build test in an image. For example, `make vm-build-freebsd` will build the source tree in the FreeBSD image. The command can be executed from either the source tree or the build dir; if the former, `./configure` is not needed. The command will then generate the test image in `./tests/vm/` under the working directory.

Note: images created by the scripts accept a well-known RSA key pair for SSH access, so they SHOULD NOT be exposed to external interfaces if you are concerned about attackers taking control of the guest and potentially exploiting a QEMU security bug to compromise the host.

QEMU binaries

By default, `qemu-system-x86_64` is searched in `$PATH` to run the guest. If there isn't one, or if it is older than 2.10, the test won't work. In this case, provide the QEMU binary in env var: `QEMU=/path/to/qemu-2.10+`.

Likewise the path to `qemu-img` can be set in `QEMU_IMG` environment variable.

Make jobs

The `-j$X` option in the make command line is not propagated into the VM, specify `J=$X` to control the make jobs in the guest.

Debugging

Add `DEBUG=1` and/or `V=1` to the make command to allow interactive debugging and verbose output. If this is not enough, see the next section. `V=1` will be propagated down into the make jobs in the guest.

Manual invocation

Each guest script is an executable script with the same command line options. For example to work with the `netbsd` guest, use `$QEMU_SRC/tests/vm/netbsd`:

```
$ cd $QEMU_SRC/tests/vm

# To bootstrap the image
$ ./netbsd --build-image --image /var/tmp/netbsd.img
<...>

# To run an arbitrary command in guest (the output will not be echoed unless
# --debug is added)
$ ./netbsd --debug --image /var/tmp/netbsd.img uname -a

# To build QEMU in guest
$ ./netbsd --debug --image /var/tmp/netbsd.img --build-qemu $QEMU_SRC

# To get to an interactive shell
$ ./netbsd --interactive --image /var/tmp/netbsd.img sh
```

Adding new guests

Please look at existing guest scripts for how to add new guests.

Most importantly, create a subclass of `BaseVM` and implement `build_image()` method and define `BUILD_SCRIPT`, then finally call `basevm.main()` from the script's `main()`.

- Usually in `build_image()`, a template image is downloaded from a predefined URL. `BaseVM._download_with_cache()` takes care of the cache and the checksum, so consider using it.
- Once the image is downloaded, users, SSH server and QEMU build deps should be set up:
 - Root password set to `BaseVM.ROOT_PASS`
 - User `BaseVM.GUEST_USER` is created, and password set to `BaseVM.GUEST_PASS`
 - SSH service is enabled and started on boot, `$QEMU_SRC/tests/keys/id_rsa.pub` is added to `ssh's authorized_keys` file of both root and the normal user
 - DHCP client service is enabled and started on boot, so that it can automatically configure the virtio-net-pci NIC and communicate with QEMU user net (10.0.2.2)
 - Necessary packages are installed to untar the source tarball and build QEMU
- Write a proper `BUILD_SCRIPT` template, which should be a shell script that untars a raw virtio-blk block device, which is the tarball data blob of the QEMU source tree, then configure/build it. Running “make check” is also recommended.

6.3.6 Image fuzzer testing

An image fuzzer was added to exercise format drivers. Currently only `qcow2` is supported. To start the fuzzer, run

```
tests/image-fuzzer/runner.py -c '["qemu-img", "info", "$test_img"]' /tmp/test qcow2
```

Alternatively, some command different from “`qemu-img info`” can be tested, by changing the `-c` option.

6.3.7 Acceptance tests using the Avocado Framework

The `tests/acceptance` directory hosts functional tests, also known as acceptance level tests. They're usually higher level tests, and may interact with external resources and with various guest operating systems.

These tests are written using the Avocado Testing Framework (which must be installed separately) in conjunction with a the `avocado_qemu.Test` class, implemented at `tests/acceptance/avocado_qemu`.

Tests based on `avocado_qemu.Test` can easily:

- Customize the command line arguments given to the convenience `self.vm` attribute (a `QEMUMachine` instance)
- Interact with the QEMU monitor, send QMP commands and check their results
- Interact with the guest OS, using the convenience console device (which may be useful to assert the effectiveness and correctness of command line arguments or QMP commands)
- Interact with external data files that accompany the test itself (see `self.get_data()`)
- Download (and cache) remote data files, such as firmware and kernel images
- Have access to a library of guest OS images (by means of the `avocado.utils.vmimage` library)
- Make use of various other test related utilities available at the test class itself and at the utility library:
 - <http://avocado-framework.readthedocs.io/en/latest/api/test/avocado.html#avocado.Test>
 - <http://avocado-framework.readthedocs.io/en/latest/api/utis/avocado.utils.html>

Running tests

You can run the acceptance tests simply by executing:

```
make check-acceptance
```

This involves the automatic creation of Python virtual environment within the build tree (at `tests/venv`) which will have all the right dependencies, and will save tests results also within the build tree (at `tests/results`).

Note: the build environment must be using a Python 3 stack, and have the `venv` and `pip` packages installed. If necessary, make sure `configure` is called with `--python=` and that those modules are available. On Debian and Ubuntu based systems, depending on the specific version, they may be on packages named `python3-venv` and `python3-pip`.

The scripts installed inside the virtual environment may be used without an “activation”. For instance, the Avocado test runner may be invoked by running:

```
tests/venv/bin/avocado run $OPTION1 $OPTION2 tests/acceptance/
```

Manual Installation

To manually install Avocado and its dependencies, run:

```
pip install --user avocado-framework
```

Alternatively, follow the instructions on this link:

<https://avocado-framework.readthedocs.io/en/latest/guides/user/chapters/installing.html>

Overview

The `tests/acceptance/avocado_qemu` directory provides the `avocado_qemu` Python module, containing the `avocado_qemu.Test` class. Here’s a simple usage example:

```
from avocado_qemu import Test

class Version(Test):
    """
    :avocado: tags=quick
    """
    def test_qmp_human_info_version(self):
        self.vm.launch()
        res = self.vm.command('human-monitor-command',
                              command_line='info version')
        self.assertRegexMatches(res, r'^(\d+\.\d+\.\d+)')
```

To execute your test, run:

```
avocado run version.py
```

Tests may be classified according to a convention by using docstring directives such as `:avocado: tags=TAG1, TAG2`. To run all tests in the current directory, tagged as “quick”, run:

```
avocado run -t quick .
```

The `avocado_qemu.Test` base test class

The `avocado_qemu.Test` class has a number of characteristics that are worth being mentioned right away.

First of all, it attempts to give each test a ready to use `QEMUMachine` instance, available at `self.vm`. Because many tests will tweak the QEMU command line, launching the `QEMUMachine` (by using `self.vm.launch()`) is left to the test writer.

The base test class has also support for tests with more than one `QEMUMachine`. The way to get machines is through the `self.get_vm()` method which will return a `QEMUMachine` instance. The `self.get_vm()` method accepts arguments that will be passed to the `QEMUMachine` creation and also an optional *name* attribute so you can identify a specific machine and get it more than once through the tests methods. A simple and hypothetical example follows:

```
from avocado_qemu import Test

class MultipleMachines(Test):
    def test_multiple_machines(self):
        first_machine = self.get_vm()
        second_machine = self.get_vm()
        self.get_vm(name='third_machine').launch()

        first_machine.launch()
        second_machine.launch()

        first_res = first_machine.command(
            'human-monitor-command',
            command_line='info version')

        second_res = second_machine.command(
            'human-monitor-command',
            command_line='info version')

        third_res = self.get_vm(name='third_machine').command(
            'human-monitor-command',
            command_line='info version')

        self.assertEqual(first_res, second_res, third_res)
```

At test “tear down”, `avocado_qemu.Test` handles all the `QEMUMachines` shutdown.

QEMUMachine

The `QEMUMachine` API is already widely used in the Python iotests, device-crash-test and other Python scripts. It’s a wrapper around the execution of a QEMU binary, giving its users:

- the ability to set command line arguments to be given to the QEMU binary
- a ready to use QMP connection and interface, which can be used to send commands and inspect its results, as well as asynchronous events
- convenience methods to set commonly used command line arguments in a more succinct and intuitive way

QEMU binary selection

The QEMU binary used for the `self.vm` `QEMUMachine` instance will primarily depend on the value of the `qemu_bin` parameter. If it’s not explicitly set, its default value will be the result of a dynamic probe in the same

source tree. A suitable binary will be one that targets the architecture matching host machine.

Based on this description, test writers will usually rely on one of the following approaches:

- 1) Set `qemu_bin`, and use the given binary
- 2) Do not set `qemu_bin`, and use a QEMU binary named like “`qemu-system-${arch}`”, either in the current working directory, or in the current source tree.

The resulting `qemu_bin` value will be preserved in the `avocado_qemu.Test` as an attribute with the same name.

Attribute reference

Besides the attributes and methods that are part of the base `avocado.Test` class, the following attributes are available on any `avocado_qemu.Test` instance.

vm

A `QEMUMachine` instance, initially configured according to the given `qemu_bin` parameter.

arch

The architecture can be used on different levels of the stack, e.g. by the framework or by the test itself. At the framework level, it will currently influence the selection of a QEMU binary (when one is not explicitly given).

Tests are also free to use this attribute value, for their own needs. A test may, for instance, use the same value when selecting the architecture of a kernel or disk image to boot a VM with.

The `arch` attribute will be set to the test parameter of the same name. If one is not given explicitly, it will either be set to `None`, or, if the test is tagged with one (and only one) `:avocado: tags=arch:VALUE` tag, it will be set to `VALUE`.

machine

The machine type that will be set to all `QEMUMachine` instances created by the test.

The `machine` attribute will be set to the test parameter of the same name. If one is not given explicitly, it will either be set to `None`, or, if the test is tagged with one (and only one) `:avocado: tags=machine:VALUE` tag, it will be set to `VALUE`.

qemu_bin

The preserved value of the `qemu_bin` parameter or the result of the dynamic probe for a QEMU binary in the current working directory or source tree.

Parameter reference

To understand how Avocado parameters are accessed by tests, and how they can be passed to tests, please refer to:

[https://avocado-framework.readthedocs.io/en/latest/guides/writer/chapters/writing.html](https://avocado-framework.readthedocs.io/en/latest/guides/writer/chapters/writing.html#accessing-test-parameters)
[↪ #accessing-test-parameters](#)

Parameter values can be easily seen in the log files, and will look like the following:

```
PARAMS (key=qemu_bin, path=*, default=./qemu-system-x86_64) => './qemu-system-x86_64
```

arch

The architecture that will influence the selection of a QEMU binary (when one is not explicitly given).

Tests are also free to use this parameter value, for their own needs. A test may, for instance, use the same value when selecting the architecture of a kernel or disk image to boot a VM with.

This parameter has a direct relation with the `arch` attribute. If not given, it will default to `None`.

machine

The machine type that will be set to all `QEMUMachine` instances created by the test.

qemu_bin

The exact QEMU binary to be used on `QEMUMachine`.

Skipping tests

The Avocado framework provides Python decorators which allow for easily skip tests running under certain conditions. For example, on the lack of a binary on the test system or when the running environment is a CI system. For further information about those decorators, please refer to:

```
https://avocado-framework.readthedocs.io/en/latest/guides/writer/chapters/writing.html  
↪ #skipping-tests
```

While the conditions for skipping tests are often specifics of each one, there are recurring scenarios identified by the QEMU developers and the use of environment variables became a kind of standard way to enable/disable tests.

Here is a list of the most used variables:

AVOCADO_ALLOW_LARGE_STORAGE

Tests which are going to fetch or produce assets considered *large* are not going to run unless that `AVOCADO_ALLOW_LARGE_STORAGE=1` is exported on the environment.

The definition of *large* is a bit arbitrary here, but it usually means an asset which occupies at least 1GB of size on disk when uncompressed.

AVOCADO_ALLOW_UNTRUSTED_CODE

There are tests which will boot a kernel image or firmware that can be considered not safe to run on the developer's workstation, thus they are skipped by default. The definition of *not safe* is also arbitrary but usually it means a blob which either its source or build process aren't public available.

You should export `AVOCADO_ALLOW_UNTRUSTED_CODE=1` on the environment in order to allow tests which make use of those kind of assets.

AVOCADO_TIMEOUT_EXPECTED

The Avocado framework has a timeout mechanism which interrupts tests to avoid the test suite of getting stuck. The timeout value can be set via test parameter or property defined in the test class, for further details:

```
https://avocado-framework.readthedocs.io/en/latest/guides/writer/chapters/writing.html
↪ #setting-a-test-timeout
```

Even though the timeout can be set by the test developer, there are some tests that may not have a well-defined limit of time to finish under certain conditions. For example, tests that take longer to execute when QEMU is compiled with debug flags. Therefore, the *AVOCADO_TIMEOUT_EXPECTED* variable has been used to determine whether those tests should run or not.

GITLAB_CI

A number of tests are flagged to not run on the GitLab CI. Usually because they proved to be flaky or there are constraints on the CI environment which would make them fail. If you encounter a similar situation then use that variable as shown on the code snippet below to skip the test:

```
@skipIf(os.getenv('GITLAB_CI'), 'Running on GitLab')
def test(self):
    do_something()
```

Uninstalling Avocado

If you’ve followed the manual installation instructions above, you can easily uninstall Avocado. Start by listing the packages you have installed:

```
pip list --user
```

And remove any package you want with:

```
pip uninstall <package_name>
```

If you’ve used `make check-acceptance`, the Python virtual environment where Avocado is installed will be cleaned up as part of `make check-clean`.

6.3.8 Testing with “make check-tcg”

The check-tcg tests are intended for simple smoke tests of both linux-user and softmmu TCG functionality. However to build test programs for guest targets you need to have cross compilers available. If your distribution supports cross compilers you can do something as simple as:

```
apt install gcc-aarch64-linux-gnu
```

The configure script will automatically pick up their presence. Sometimes compilers have slightly odd names so the availability of them can be prompted by passing in the appropriate configure option for the architecture in question, for example:

```
$(configure) --cross-cc-aarch64=aarch64-cc
```

There is also a `--cross-cc-flags-ARCH` flag in case additional compiler flags are needed to build for a given target.

If you have the ability to run containers as the user you can also take advantage of the build systems “Docker” support. It will then use containers to build any test case for an enabled guest where there is no system compiler available. See *Docker based tests* for details.

Running subset of tests

You can build the tests for one architecture:

```
make build-tcg-tests-$TARGET
```

And run with:

```
make run-tcg-tests-$TARGET
```

Adding `V=1` to the invocation will show the details of how to invoke QEMU for the test which is useful for debugging tests.

TCG test dependencies

The TCG tests are deliberately very light on dependencies and are either totally bare with minimal gcc lib support (for softmmu tests) or just glibc (for linux-user tests). This is because getting a cross compiler to work with additional libraries can be challenging.

Other TCG Tests

There are a number of out-of-tree test suites that are used for more extensive testing of processor features.

KVM Unit Tests

The KVM unit tests are designed to run as a Guest OS under KVM but there is no reason why they can’t exercise the TCG as well. It provides a minimal OS kernel with hooks for enabling the MMU as well as reporting test results via a special device:

```
https://git.kernel.org/pub/scm/virt/kvm/kvm-unit-tests.git
```

Linux Test Project

The LTP is focused on exercising the syscall interface of a Linux kernel. It checks that syscalls behave as documented and strives to exercise as many corner cases as possible. It is a useful test suite to run to exercise QEMU’s linux-user code:

```
https://linux-test-project.github.io/
```

6.4 Fuzzing

This document describes the virtual-device fuzzing infrastructure in QEMU and how to use it to implement additional fuzzers.

6.4.1 Basics

Fuzzing operates by passing inputs to an entry point/target function. The fuzzer tracks the code coverage triggered by the input. Based on these findings, the fuzzer mutates the input and repeats the fuzzing.

To fuzz QEMU, we rely on libfuzzer. Unlike other fuzzers such as AFL, libfuzzer is an *in-process* fuzzer. For the developer, this means that it is their responsibility to ensure that state is reset between fuzzing-runs.

6.4.2 Building the fuzzers

NOTE: If possible, build a 32-bit binary. When forking, the 32-bit fuzzer is much faster, since the page-map has a smaller size. This is due to the fact that AddressSanitizer maps ~20TB of memory, as part of its detection. This results in a large page-map, and a much slower `fork()`.

To build the fuzzers, install a recent version of clang: Configure with (substitute the clang binaries with the version you installed). Here, `enable-sanitizers`, is optional but it allows us to reliably detect bugs such as out-of-bounds accesses, use-after-frees, double-frees etc.:

```
CC=clang-8 CXX=clang++-8 /path/to/configure --enable-fuzzing \
                                         --enable-sanitizers
```

Fuzz targets are built similarly to system targets:

```
make qemu-fuzz-i386
```

This builds `./qemu-fuzz-i386`

The first option to this command is: `--fuzz-target=FUZZ_NAME` To list all of the available fuzzers run `qemu-fuzz-i386` with no arguments.

For example:

```
./qemu-fuzz-i386 --fuzz-target=virtio-scsi-fuzz
```

Internally, libfuzzer parses all arguments that do not begin with `--`. Information about these is available by passing `-help=1`

Now the only thing left to do is wait for the fuzzer to trigger potential crashes.

6.4.3 Useful libFuzzer flags

As mentioned above, libFuzzer accepts some arguments. Passing `-help=1` will list the available arguments. In particular, these arguments might be helpful:

- `CORPUS_DIR/` : Specify a directory as the last argument to libFuzzer. libFuzzer stores each “interesting” input in this corpus directory. The next time you run libFuzzer, it will read all of the inputs from the corpus, and continue fuzzing from there. You can also specify multiple directories. libFuzzer loads existing inputs from all specified directories, but will only write new ones to the first one specified.
- `-max_len=4096` : specify the maximum byte-length of the inputs libFuzzer will generate.

- `-close_fd_mask={1, 2, 3}` : close, stderr, or both. Useful for targets that trigger many debug/error messages, or create output on the serial console.
- `-jobs=4 -workers=4` : These arguments configure libFuzzer to run 4 fuzzers in parallel (4 fuzzing jobs in 4 worker processes). Alternatively, with only `-jobs=N`, libFuzzer automatically spawns a number of workers less than or equal to half the available CPU cores. Replace 4 with a number appropriate for your machine. Make sure to specify a `CORPUS_DIR`, which will allow the parallel fuzzers to share information about the interesting inputs they find.
- `-use_value_profile=1` : For each comparison operation, libFuzzer computes `(caller_pc & 4095) | (popcnt(Arg1 ^ Arg2) << 12)` and places this in the coverage table. Useful for targets with “magic” constants. If Arg1 came from the fuzzer’s input and Arg2 is a magic constant, then each time the Hamming distance between Arg1 and Arg2 decreases, libFuzzer adds the input to the corpus.
- `-shrink=1` : Tries to make elements of the corpus “smaller”. Might lead to better coverage performance, depending on the target.

Note that libFuzzer’s exact behavior will depend on the version of clang and libFuzzer used to build the device fuzzers.

6.4.4 Generating Coverage Reports

Code coverage is a crucial metric for evaluating a fuzzer’s performance. libFuzzer’s output provides a “cov: ” column that provides a total number of unique blocks/edges covered. To examine coverage on a line-by-line basis we can use Clang coverage:

1. Configure libFuzzer to store a corpus of all interesting inputs (see `CORPUS_DIR` above)
2. `./configure` the QEMU build with

```
--enable-fuzzing \  
--extra-cflags="-fprofile-instr-generate -fcoverage-mapping"
```

3. Re-run the fuzzer. Specify `$CORPUS_DIR/*` as an argument, telling libfuzzer to execute all of the inputs in `$CORPUS_DIR` and exit. Once the process exits, you should find a file, “default.profrw” in the working directory.
4. Execute these commands to generate a detailed HTML coverage-report:

```
llvm-profdata merge -output=default.profdata default.profrw  
llvm-cov show ./path/to/qemu-fuzz-i386 -instr-profile=default.profdata \  
--format html -output-dir=/path/to/output/report
```

6.4.5 Adding a new fuzzer

Coverage over virtual devices can be improved by adding additional fuzzers. Fuzzers are kept in `tests/qtest/fuzz/` and should be added to `tests/qtest/fuzz/meson.build`

Fuzzers can rely on both `qtest` and `libqos` to communicate with virtual devices.

1. Create a new source file. For example `tests/qtest/fuzz/foo-device-fuzz.c`.
2. Write the fuzzing code using the `libqtest/libqos` API. See existing fuzzers for reference.
3. Add the fuzzer to `tests/qtest/fuzz/meson.build`.

Fuzzers can be more-or-less thought of as special `qtest` programs which can modify the `qtest` commands and/or `qtest` command arguments based on inputs provided by libfuzzer. Libfuzzer passes a byte array and length. Commonly the fuzzer loops over the byte-array interpreting it as a list of `qtest` commands, addresses, or values.

6.4.6 The Generic Fuzzer

Writing a fuzz target can be a lot of effort (especially if a device driver has not been built-out within libqos). Many devices can be fuzzed to some degree, without any device-specific code, using the generic-fuzz target.

The generic-fuzz target is capable of fuzzing devices over their PIO, MMIO, and DMA input-spaces. To apply the generic-fuzz to a device, we need to define two env-variables, at minimum:

- `QEMU_FUZZ_ARGS=` is the set of QEMU arguments used to configure a machine, with the device attached. For example, if we want to fuzz the virtio-net device attached to a pc-i440fx machine, we can specify:

```
QEMU_FUZZ_ARGS="-M pc -nodefaults -netdev user,id=user0 \
-device virtio-net,netdev=user0"
```

- `QEMU_FUZZ_OBJECTS=` is a set of space-delimited strings used to identify the MemoryRegions that will be fuzzed. These strings are compared against MemoryRegion names and MemoryRegion owner names, to decide whether each MemoryRegion should be fuzzed. These strings support globbing. For the virtio-net example, we could use one of

```
QEMU_FUZZ_OBJECTS='virtio-net'
QEMU_FUZZ_OBJECTS='virtio*'
QEMU_FUZZ_OBJECTS='virtio* pcspk' # Fuzz the virtio devices and the speaker
QEMU_FUZZ_OBJECTS='*' # Fuzz the whole machine``
```

The "info mtree" and "info qom-tree" monitor commands can be especially useful for identifying the MemoryRegion and Object names used for matching.

As a generic rule-of-thumb, the more MemoryRegions/Devices we match, the greater the input-space, and the smaller the probability of finding crashing inputs for individual devices. As such, it is usually a good idea to limit the fuzzer to only a few MemoryRegions.

To ensure that these env variables have been configured correctly, we can use:

```
./qemu-fuzz-i386 --fuzz-target=generic-fuzz -runs=0
```

The output should contain a complete list of matched MemoryRegions.

6.4.7 OSS-Fuzz

QEMU is continuously fuzzed on *OSS-Fuzz* (<https://github.com/google/oss-fuzz>). By default, the OSS-Fuzz build will try to fuzz every fuzz-target. Since the generic-fuzz target requires additional information provided in environment variables, we pre-define some generic-fuzz configs in `tests/qtest/fuzz/generic_fuzz_configs.h`. Each config must specify:

- `.name`: To identify the fuzzer config
- `.args` OR `.argfunc`: A string or pointer to a function returning a string. These strings are used to specify the `QEMU_FUZZ_ARGS` environment variable. `argfunc` is useful when the config relies on e.g. a dynamically created temp directory, or a free tcp/udp port.
- `.objects`: A string that specifies the `QEMU_FUZZ_OBJECTS` environment variable.

To fuzz additional devices/device configuration on OSS-Fuzz, send patches for either a new device-specific fuzzer or a new generic-fuzz config.

Build details:

- The Dockerfile that sets up the environment for building QEMU's fuzzers on OSS-Fuzz can be found in the OSS-Fuzz repository (<https://github.com/google/oss-fuzz/blob/master/projects/qemu/Dockerfile>)

- The script responsible for building the fuzzers can be found in the QEMU source tree at `scripts/oss-fuzz/build.sh`

6.4.8 Implementation Details / Fuzzer Lifecycle

The fuzzer has two entrypoints that libfuzzer calls. libfuzzer provides its own `main()`, which performs some setup, and calls the entrypoints:

`LLVMFuzzerInitialize`: called prior to fuzzing. Used to initialize all of the necessary state

`LLVMFuzzerTestOneInput`: called for each fuzzing run. Processes the input and resets the state at the end of each run.

In more detail:

`LLVMFuzzerInitialize` parses the arguments to the fuzzer (must start with two dashes, so they are ignored by libfuzzer `main()`). Currently, the arguments select the fuzz target. Then, the QTest client is initialized. If the target requires QOS, QGraph is set up and the QOM/LIBQOS modules are initialized. Then the QGraph is walked and the QEMU `cmd_line` is determined and saved.

After this, the `vl.c:qemu_main` is called to set up the guest. There are target-specific hooks that can be called before and after `qemu_main`, for additional setup (e.g. PCI setup, or VM snapshotting).

`LLVMFuzzerTestOneInput`: Uses QTest/QOS functions to act based on the fuzz input. It is also responsible for manually calling `main_loop_wait` to ensure that bottom halves are executed and any cleanup required before the next input.

Since the same process is reused for many fuzzing runs, QEMU state needs to be reset at the end of each run. There are currently two implemented options for resetting state:

- Reboot the guest between runs. - *Pros*: Straightforward and fast for simple fuzz targets.
 - *Cons*: Depending on the device, does not reset all device state. If the device requires some initialization prior to being ready for fuzzing (common for QOS-based targets), this initialization needs to be done after each reboot.
 - *Example target*: `i440fx-qtest-reboot-fuzz`
- **Run each test case in a separate forked process and copy the coverage** information back to the parent. This is fairly similar to AFL's "deferred" fork-server mode [3]
 - *Pros*: Relatively fast. Devices only need to be initialized once. No need to do slow reboots or vmloads.
 - *Cons*: **Not officially supported by libfuzzer. Does not work well for** devices that rely on dedicated threads.
 - *Example target*: `virtio-net-fork-fuzz`

6.5 Control-Flow Integrity (CFI)

This document describes the current control-flow integrity (CFI) mechanism in QEMU. How it can be enabled, its benefits and deficiencies, and how it affects new and existing code in QEMU

6.5.1 Basics

CFI is a hardening technique that focusing on guaranteeing that indirect function calls have not been altered by an attacker. The type used in QEMU is a forward-edge control-flow integrity that ensures function calls performed

through function pointers, always call a “compatible” function. A compatible function is a function with the same signature of the function pointer declared in the source code.

This type of CFI is entirely compiler-based and relies on the compiler knowing the signature of every function and every function pointer used in the code. As of now, the only compiler that provides support for CFI is Clang.

CFI is best used on production binaries, to protect against unknown attack vectors.

In case of a CFI violation (i.e. call to a non-compatible function) QEMU will terminate abruptly, to stop the possible attack.

6.5.2 Building with CFI

NOTE: CFI requires the use of link-time optimization. Therefore, when CFI is selected, LTO will be automatically enabled.

To build with CFI, the minimum requirement is Clang 6+. If you are planning to also enable fuzzing, then Clang 11+ is needed (more on this later).

Given the use of LTO, a version of AR that supports LLVM IR is required. The easiest way of doing this is by selecting the AR provided by LLVM:

```
AR=llvm-ar-9 CC=clang-9 CXX=clang++-9 /path/to/configure --enable-cfi
```

CFI is enabled on every binary produced.

If desired, an additional flag to increase the verbosity of the output in case of a CFI violation is offered (`--enable-debug-cfi`).

6.5.3 Using QEMU built with CFI

A binary with CFI will work exactly like a standard binary. In case of a CFI violation, the binary will terminate with an illegal instruction signal.

6.5.4 Incompatible code with CFI

As mentioned above, CFI is entirely compiler-based and therefore relies on compile-time knowledge of the code. This means that, while generally supported for most code, some specific use pattern can break CFI compatibility, and create false-positives. The two main patterns that can cause issues are:

- Just-in-time compiled code: since such code is created at runtime, the jump to the buffer containing JIT code will fail.
- Libraries loaded dynamically, e.g. with `dlopen/dlsym`, since the library was not known at compile time.

Current areas of QEMU that are not entirely compatible with CFI are:

1. TCG, since the idea of TCG is to pre-compile groups of instructions at runtime to speed-up interpretation, quite similarly to a JIT compiler
2. TCI, where the interpreter has to interpret the generic *call* operation
3. Plugins, since a plugin is implemented as an external library
4. Modules, since they are implemented as an external library
5. Directly calling signal handlers from the QEMU source code, since the signal handler may have been provided by an external library or even plugged at runtime.

6.5.5 Disabling CFI for a specific function

If you are working on function that is performing a call using an incompatible way, as described before, you can selectively disable CFI checks for such function by using the decorator `QEMU_DISABLE_CFI` at function definition, and add an explanation on why the function is not compatible with CFI. An example of the use of `QEMU_DISABLE_CFI` is provided here:

```
/*
 * Disable CFI checks.
 * TCG creates binary blobs at runtime, with the transformed code.
 * A TB is a blob of binary code, created at runtime and called with an
 * indirect function call. Since such function did not exist at compile time,
 * the CFI runtime has no way to verify its signature and would fail.
 * TCG is not considered a security-sensitive part of QEMU so this does not
 * affect the impact of CFI in environment with high security requirements
 */
QEMU_DISABLE_CFI
static inline tcg_target_ulong cpu_tb_exec(CPUState *cpu, TranslationBlock *itb)
```

NOTE: CFI needs to be disabled at the **caller** function, (i.e. a compatible cfi function that calls a non-compatible one), since the check is performed when the function call is performed.

6.5.6 CFI and fuzzing

There is generally no advantage of using CFI and fuzzing together, because they target different environments (production for CFI, debug for fuzzing).

CFI could be used in conjunction with fuzzing to identify a broader set of bugs that may not end immediately in a segmentation fault or triggering an assertion. However, other sanitizers such as address and ub sanitizers can identify such bugs in a more precise way than CFI.

There is, however, an interesting use case in using CFI in conjunction with fuzzing, that is to make sure that CFI is not triggering any false positive in remote-but-possible parts of the code.

CFI can be enabled with fuzzing, but with some caveats: 1. Fuzzing relies on the linker performing function wrapping at link-time. The standard BFD linker does not support function wrapping when LTO is also enabled. The workaround is to use LLVM's lld linker. 2. Fuzzing also relies on a custom linker script, which is only supported by lld with version 11+.

In other words, to compile with fuzzing and CFI, clang 11+ is required, and lld needs to be used as a linker:

```
AR=llvm-ar-11 CC=clang-11 CXX=lang++-11 /path/to/configure --enable-cfi \
    -enable-fuzzing --extra-ldflags="-fuse-ld=lld"
```

and then, compile the fuzzers as usual.

6.6 Load and Store APIs

QEMU internally has multiple families of functions for performing loads and stores. This document attempts to enumerate them all and indicate when to use them. It does not provide detailed documentation of each API – for that you should look at the documentation comments in the relevant header files.

6.6.1 `ld*_p` and `st*_p`

These functions operate on a host pointer, and should be used when you already have a pointer into host memory (corresponding to guest ram or a local buffer). They deal with doing accesses with the desired endianness and with correctly handling potentially unaligned pointer values.

Function names follow the pattern:

load: `ld{sign}{size}_{endian}_p(ptr)`

store: `st{size}_{endian}_p(ptr, val)`

sign

- (empty) : for 32 or 64 bit sizes
- `u` : unsigned
- `s` : signed

size

- `b` : 8 bits
- `w` : 16 bits
- `l` : 32 bits
- `q` : 64 bits

endian

- `he` : host endian
- `be` : big endian
- `le` : little endian

The `_{endian}` infix is omitted for target-endian accesses.

The target endian accessors are only available to source files which are built per-target.

There are also functions which take the size as an argument:

load: `ldn{endian}_p(ptr, sz)`

which performs an unsigned load of `sz` bytes from `ptr` as an `{endian}` order value and returns it in a `uint64_t`.

store: `stn{endian}_p(ptr, sz, val)`

which stores `val` to `ptr` as an `{endian}` order value of size `sz` bytes.

Regexes for `git grep`

- `<ld[us]?[bwlq]\(_[hbl]e\)?_p\>`
- `<st[bwlq]\(_[hbl]e\)?_p\>`
- `<ldn_\([hbl]e\)?_p\>`
- `<stn_\([hbl]e\)?_p\>`

6.6.2 `cpu_{ld,st}*_mmuidx_ra`

These functions operate on a guest virtual address plus a context, known as a “mmu index” or `mmuidx`, which controls how that virtual address is translated. The meaning of the indexes are target specific, but specifying a particular index

might be necessary if, for instance, the helper requires an “always as non-privileged” access rather than the default access for the current state of the guest CPU.

These functions may cause a guest CPU exception to be taken (e.g. for an alignment fault or MMU fault) which will result in guest CPU state being updated and control longjmp’ing out of the function call. They should therefore only be used in code that is implementing emulation of the guest CPU.

The `retaddr` parameter is used to control unwinding of the guest CPU state in case of a guest CPU exception. This is passed to `cpu_restore_state()`. Therefore the value should either be 0, to indicate that the guest CPU state is already synchronized, or the result of `GETPC()` from the top level `HELPER(foo)` function, which is a return address into the generated code¹.

Function names follow the pattern:

load: `cpu_ld{sign}{size}{end}_mmuidx_ra(env, ptr, mmuidx, retaddr)`

store: `cpu_st{size}{end}_mmuidx_ra(env, ptr, val, mmuidx, retaddr)`

sign

- (empty) : for 32 or 64 bit sizes
- `u` : unsigned
- `s` : signed

size

- `b` : 8 bits
- `w` : 16 bits
- `l` : 32 bits
- `q` : 64 bits

end

- (empty) : for target endian, or 8 bit sizes
- `_be` : big endian
- `_le` : little endian

Regexes for git grep:

- `<cpu_ld[us]\?[bwlq](_[ble])\?_mmuidx_ra\>`
- `<cpu_st[bwlq](_[ble])\?_mmuidx_ra\>`

6.6.3 `cpu_{ld,st}*_data_ra`

These functions work like the `cpu_{ld,st}_mmuidx_ra` functions except that the `mmuidx` parameter is taken from the current mode of the guest CPU, as determined by `cpu_mmu_index(env, false)`.

These are generally the preferred way to do accesses by guest virtual address from helper functions, unless the access should be performed with a context other than the default.

Function names follow the pattern:

load: `cpu_ld{sign}{size}{end}_data_ra(env, ptr, ra)`

store: `cpu_st{size}{end}_data_ra(env, ptr, val, ra)`

¹ Note that `GETPC()` should be used with great care: calling it in other functions that are *not* the top level `HELPER(foo)` will cause unexpected behavior. Instead, the value of `GETPC()` should be read from the helper and passed if needed to the functions that the helper calls.

sign

- (empty) : for 32 or 64 bit sizes
- u : unsigned
- s : signed

size

- b : 8 bits
- w : 16 bits
- l : 32 bits
- q : 64 bits

end

- (empty) : for target endian, or 8 bit sizes
- _be : big endian
- _le : little endian

Regexes for git grep:

- `<cpu_ld[us]\?[bwlq](_[ble])\?_data_ra\>`
- `<cpu_st[bwlq](_[ble])\?_data_ra\>`

6.6.4 `cpu_{ld,st}*_data`

These functions work like the `cpu_{ld,st}_data_ra` functions except that the `retaddr` parameter is 0, and thus does not unwind guest CPU state.

This means they must only be used from helper functions where the translator has saved all necessary CPU state. These functions are the right choice for calls made from hooks like the CPU `do_interrupt` hook or when you know for certain that the translator had to save all the CPU state anyway.

Function names follow the pattern:

load: `cpu_ld{sign}{size}{end}_data(env, ptr)`

store: `cpu_st{size}{end}_data(env, ptr, val)`

sign

- (empty) : for 32 or 64 bit sizes
- u : unsigned
- s : signed

size

- b : 8 bits
- w : 16 bits
- l : 32 bits
- q : 64 bits

end

- (empty) : for target endian, or 8 bit sizes

- `_be` : big endian
- `_le` : little endian

Regexes for git grep

- `\<cpu_ld[us]\?[bwlq](_[bl]e)\?_data\>`
- `\<cpu_st[bwlq](_[bl]e)\?_data\+\>`

6.6.5 `cpu_ld*_code`

These functions perform a read for instruction execution. The `mmuidx` parameter is taken from the current mode of the guest CPU, as determined by `cpu_mmu_index(env, true)`. The `retaddr` parameter is 0, and thus does not unwind guest CPU state, because CPU state is always synchronized while translating instructions. Any guest CPU exception that is raised will indicate an instruction execution fault rather than a data read fault.

In general these functions should not be used directly during translation. There are wrapper functions that are to be used which also take care of plugins for tracing.

Function names follow the pattern:

```
load: cpu_ld{sign}{size}_code(env, ptr)
```

sign

- (empty) : for 32 or 64 bit sizes
- `u` : unsigned
- `s` : signed

size

- `b` : 8 bits
- `w` : 16 bits
- `l` : 32 bits
- `q` : 64 bits

Regexes for git grep:

- `\<cpu_ld[us]\?[bwlq]_code\>`

6.6.6 `translator_ld*`

These functions are a wrapper for `cpu_ld*_code` which also perform any actions required by any tracing plugins. They are only to be called during the translator callback `translate_insn`.

There is a set of functions ending in `_swap` which, if the parameter is true, returns the value in the endianness that is the reverse of the guest native endianness, as determined by `TARGET_WORDS_BIGENDIAN`.

Function names follow the pattern:

```
load: translator_ld{sign}{size}(env, ptr)
```

```
swap: translator_ld{sign}{size}_swap(env, ptr, swap)
```

sign

- (empty) : for 32 or 64 bit sizes
- `u` : unsigned

- `s` : signed

size

- `b` : 8 bits
- `w` : 16 bits
- `l` : 32 bits
- `q` : 64 bits

Regexes for git grep

- `<translator_ld[us]\?[bwlq]\(_swap\)\?\>`

6.6.7 helper_*_{ld,st}*_mmu

These functions are intended primarily to be called by the code generated by the TCG backend. They may also be called by target CPU helper function code. Like the `cpu_{ld,st}_mmuidx_ra` functions they perform accesses by guest virtual address, with a given `mmuidx`.

These functions specify an `opindex` parameter which encodes (among other things) the mmu index to use for the access. This parameter should be created by calling `make_memop_idx()`.

The `retaddr` parameter should be the result of `GETPC()` called directly from the top level `HELPER(foo)` function (or 0 if no guest CPU state unwinding is required).

TODO The names of these functions are a bit odd for historical reasons because they were originally expected to be called only from within generated code. We should rename them to bring them more in line with the other memory access functions. The explicit endianness is the only feature they have beyond `*_mmuidx_ra`.

```
load: helper_{endian}_ld{sign}{size}_mmu(env, addr, opindex, retaddr)
```

```
store: helper_{endian}_st{size}_mmu(env, addr, val, opindex, retaddr)
```

sign

- (empty) : for 32 or 64 bit sizes
- `u` : unsigned
- `s` : signed

size

- `b` : 8 bits
- `w` : 16 bits
- `l` : 32 bits
- `q` : 64 bits

endian

- `le` : little endian
- `be` : big endian
- `ret` : target endianness

Regexes for git grep

- `<helper_\(le\|be\|ret\) _ld[us]\?[bwlq]_mmu\>`
- `<helper_\(le\|be\|ret\) _st[bwlq]_mmu\>`

6.6.8 address_space_*

These functions are the primary ones to use when emulating CPU or device memory accesses. They take an AddressSpace, which is the way QEMU defines the view of memory that a device or CPU has. (They generally correspond to being the “master” end of a hardware bus or bus fabric.)

Each CPU has an AddressSpace. Some kinds of CPU have more than one AddressSpace (for instance Arm guest CPUs have an AddressSpace for the Secure world and one for NonSecure if they implement TrustZone). Devices which can do DMA-type operations should generally have an AddressSpace. There is also a “system address space” which typically has all the devices and memory that all CPUs can see. (Some older device models use the “system address space” rather than properly modelling that they have an AddressSpace of their own.)

Functions are provided for doing byte-buffer reads and writes, and also for doing one-data-item loads and stores.

In all cases the caller provides a MemTxAttrs to specify bus transaction attributes, and can check whether the memory transaction succeeded using a MemTxResult return code.

```
address_space_read(address_space, addr, attrs, buf, len)
address_space_write(address_space, addr, attrs, buf, len)
address_space_rw(address_space, addr, attrs, buf, len, is_write)
address_space_ld{sign}{size}_{endian}(address_space, addr, attrs, txresult)
address_space_st{size}_{endian}(address_space, addr, val, attrs, txresult)
```

sign

- (empty) : for 32 or 64 bit sizes
- u : unsigned

(No signed load operations are provided.)

size

- b : 8 bits
- w : 16 bits
- l : 32 bits
- q : 64 bits

endian

- le : little endian
- be : big endian

The `_{endian}` suffix is omitted for byte accesses.

Regexes for git grep

- `\<address_space_\(read\|write\|rw\) \>`
- `\<address_space_ldu\?[bwql]\(_[lb]e\) \?\>`
- `\<address_space_st[bwql]\(_[lb]e\) \?\>`

6.6.9 address_space_write_rom

This function performs a write by physical address like `address_space_write`, except that if the write is to a ROM then the ROM contents will be modified, even though a write by the guest CPU to the ROM would be ignored. This is used for non-guest writes like writes from the gdb debug stub or initial loading of ROM contents.

Note that portions of the write which attempt to write data to a device will be silently ignored – only real RAM and ROM will be written to.

Regexes for git grep

- `address_space_write_rom`

6.6.10 `{ld, st}*_phys`

These are functions which are identical to `address_space_{ld, st}*`, except that they always pass `MEMTXATTRS_UNSPECIFIED` for the transaction attributes, and ignore whether the transaction succeeded or failed.

The fact that they ignore whether the transaction succeeded means they should not be used in new code, unless you know for certain that your code will only be used in a context where the CPU or device doing the access has no way to report such an error.

load: `ld{sign}{size}_{endian}_phys`

store: `st{size}_{endian}_phys`

sign

- (empty) : for 32 or 64 bit sizes
- `u` : unsigned

(No signed load operations are provided.)

size

- `b` : 8 bits
- `w` : 16 bits
- `l` : 32 bits
- `q` : 64 bits

endian

- `le` : little endian
- `be` : big endian

The `_{endian}_` infix is omitted for byte accesses.

Regexes for git grep

- `\<ldu\?[bwlq]\(_[bl]e\)\?_phys\>`
- `\<st[bwlq]\(_[bl]e\)\?_phys\>`

6.6.11 `cpu_physical_memory_*`

These are convenience functions which are identical to `address_space_*` but operate specifically on the system address space, always pass a `MEMTXATTRS_UNSPECIFIED` set of memory attributes and ignore whether the memory transaction succeeded or failed. For new code they are better avoided:

- there is likely to be behaviour you need to model correctly for a failed read or write operation
- a device should usually perform operations on its own `AddressSpace` rather than using the system address space

```
cpu_physical_memory_read
cpu_physical_memory_write
cpu_physical_memory_rw
```

Regexes for git grep

- `\<cpu_physical_memory_\(read\|write\|rw\)\>`

6.6.12 `cpu_memory_rw_debug`

Access CPU memory by virtual address for debug purposes.

This function is intended for use by the GDB stub and similar code. It takes a virtual address, converts it to a physical address via an MMU lookup using the current settings of the specified CPU, and then performs the access (using `address_space_rw` for reads or `cpu_physical_memory_write_rom` for writes). This means that if the access is a write to a ROM then this function will modify the contents (whereas a normal guest CPU access would ignore the write attempt).

```
cpu_memory_rw_debug
```

6.6.13 `dma_memory_*`

These behave like `address_space_*`, except that they perform a DMA barrier operation first.

TODO: We should provide guidance on when you need the DMA barrier operation and when it's OK to use `address_space_*`, and make sure our existing code is doing things correctly.

```
dma_memory_read
dma_memory_write
dma_memory_rw
```

Regexes for git grep

- `\<dma_memory_\(read\|write\|rw\)\>`
- `\<ldu\?[bwlq]\(_[bl]e\)\?_dma\>`
- `\<st[bwlq]\(_[bl]e\)\?_dma\>`

6.6.14 `pci_dma_*` and `{ld,st}_*_pci_dma`

These functions are specifically for PCI device models which need to perform accesses where the PCI device is a bus master. You pass them a `PCIDevice *` and they will do `dma_memory_*` operations on the correct address space for that device.

```
pci_dma_read
pci_dma_write
pci_dma_rw
load:  ld{sign}{size}_{endian}_pci_dma
store: st{size}_{endian}_pci_dma
```

sign

- (empty) : for 32 or 64 bit sizes

- `u` : unsigned

(No signed load operations are provided.)

size

- `b` : 8 bits
- `w` : 16 bits
- `l` : 32 bits
- `q` : 64 bits

endian

- `le` : little endian
- `be` : big endian

The `_{endian}_` infix is omitted for byte accesses.

Regexes for `git grep`

- `\<pci_dma_\(read\|write\|rw\)\>`
- `\<ldu\?[bwlq]\(_[ble]\)\?_pci_dma\>`
- `\<st[bwlq]\(_[ble]\)\?_pci_dma\>`

6.7 The memory API

The memory API models the memory and I/O buses and controllers of a QEMU machine. It attempts to allow modelling of:

- ordinary RAM
- memory-mapped I/O (MMIO)
- memory controllers that can dynamically reroute physical memory regions to different destinations

The memory model provides support for

- tracking RAM changes by the guest
- setting up coalesced memory for kvm
- setting up ioeventfd regions for kvm

Memory is modelled as an acyclic graph of `MemoryRegion` objects. Sinks (leaves) are RAM and MMIO regions, while other nodes represent buses, memory controllers, and memory regions that have been rerouted.

In addition to `MemoryRegion` objects, the memory API provides `AddressSpace` objects for every root and possibly for intermediate `MemoryRegions` too. These represent memory as seen from the CPU or a device's viewpoint.

6.7.1 Types of regions

There are multiple types of memory regions (all represented by a single C type `MemoryRegion`):

- **RAM:** a RAM region is simply a range of host memory that can be made available to the guest. You typically initialize these with `memory_region_init_ram()`. Some special purposes require the variants `memory_region_init_resizeable_ram()`, `memory_region_init_ram_from_file()`, or `memory_region_init_ram_ptr()`.

- **MMIO**: a range of guest memory that is implemented by host callbacks; each read or write causes a callback to be called on the host. You initialize these with `memory_region_init_io()`, passing it a `MemoryRegionOps` structure describing the callbacks.
- **ROM**: a ROM memory region works like RAM for reads (directly accessing a region of host memory), and forbids writes. You initialize these with `memory_region_init_rom()`.
- **ROM device**: a ROM device memory region works like RAM for reads (directly accessing a region of host memory), but like MMIO for writes (invoking a callback). You initialize these with `memory_region_init_rom_device()`.
- **IOMMU region**: an IOMMU region translates addresses of accesses made to it and forwards them to some other target memory region. As the name suggests, these are only needed for modelling an IOMMU, not for simple devices. You initialize these with `memory_region_init_iommu()`.
- **container**: a container simply includes other memory regions, each at a different offset. Containers are useful for grouping several regions into one unit. For example, a PCI BAR may be composed of a RAM region and an MMIO region.

A container's subregions are usually non-overlapping. In some cases it is useful to have overlapping regions; for example a memory controller that can overlay a subregion of RAM with MMIO or ROM, or a PCI controller that does not prevent card from claiming overlapping BARs.

You initialize a pure container with `memory_region_init()`.

- **alias**: a subsection of another region. Aliases allow a region to be split apart into discontinuous regions. Examples of uses are memory banks used when the guest address space is smaller than the amount of RAM addressed, or a memory controller that splits main memory to expose a "PCI hole". Aliases may point to any type of region, including other aliases, but an alias may not point back to itself, directly or indirectly. You initialize these with `memory_region_init_alias()`.
- **reservation region**: a reservation region is primarily for debugging. It claims I/O space that is not supposed to be handled by QEMU itself. The typical use is to track parts of the address space which will be handled by the host kernel when KVM is enabled. You initialize these by passing a NULL callback parameter to `memory_region_init_io()`.

It is valid to add subregions to a region which is not a pure container (that is, to an MMIO, RAM or ROM region). This means that the region will act like a container, except that any addresses within the container's region which are not claimed by any subregion are handled by the container itself (ie by its MMIO callbacks or RAM backing). However it is generally possible to achieve the same effect with a pure container one of whose subregions is a low priority "background" region covering the whole address range; this is often clearer and is preferred. Subregions cannot be added to an alias region.

6.7.2 Migration

Where the memory region is backed by host memory (RAM, ROM and ROM device memory region types), this host memory needs to be copied to the destination on migration. These APIs which allocate the host memory for you will also register the memory so it is migrated:

- `memory_region_init_ram()`
- `memory_region_init_rom()`
- `memory_region_init_rom_device()`

For most devices and boards this is the correct thing. If you have a special case where you need to manage the migration of the backing memory yourself, you can call the functions:

- `memory_region_init_ram_nomigrate()`
- `memory_region_init_rom_nomigrate()`

- `memory_region_init_rom_device_nomigrate()`

which only initialize the `MemoryRegion` and leave handling migration to the caller.

The functions:

- `memory_region_init_resizeable_ram()`
- `memory_region_init_ram_from_file()`
- `memory_region_init_ram_from_fd()`
- `memory_region_init_ram_ptr()`
- `memory_region_init_ram_device_ptr()`

are for special cases only, and so they do not automatically register the backing memory for migration; the caller must manage migration if necessary.

6.7.3 Region names

Regions are assigned names by the constructor. For most regions these are only used for debugging purposes, but RAM regions also use the name to identify live migration sections. This means that RAM region names need to have ABI stability.

6.7.4 Region lifecycle

A region is created by one of the `memory_region_init*()` functions and attached to an object, which acts as its owner or parent. QEMU ensures that the owner object remains alive as long as the region is visible to the guest, or as long as the region is in use by a virtual CPU or another device. For example, the owner object will not die between an `address_space_map` operation and the corresponding `address_space_unmap`.

After creation, a region can be added to an address space or a container with `memory_region_add_subregion()`, and removed using `memory_region_del_subregion()`.

Various region attributes (read-only, dirty logging, coalesced mmio, `ioeventfd`) can be changed during the region lifecycle. They take effect as soon as the region is made visible. This can be immediately, later, or never.

Destruction of a memory region happens automatically when the owner object dies.

If however the memory region is part of a dynamically allocated data structure, you should call `object_unparent()` to destroy the memory region before the data structure is freed. For an example see `VFIOMSIInfo` and `VFIOQuirk` in `hw/vfio/pci.c`.

You must not destroy a memory region as long as it may be in use by a device or CPU. In order to do this, as a general rule do not create or destroy memory regions dynamically during a device's lifetime, and only call `object_unparent()` in the memory region owner's `instance_finalize` callback. The dynamically allocated data structure that contains the memory region then should obviously be freed in the `instance_finalize` callback as well.

If you break this rule, the following situation can happen:

- the memory region's owner had a reference taken via `memory_region_ref` (for example by `address_space_map`)
- the region is unparented, and has no owner anymore
- when `address_space_unmap` is called, the reference to the memory region's owner is leaked.

There is an exception to the above rule: it is okay to call `object_unparent` at any time for an alias or a container region. It is therefore also okay to create or destroy alias and container regions dynamically during a device's lifetime.

This exceptional usage is valid because aliases and containers only help QEMU building the guest’s memory map; they are never accessed directly. `memory_region_ref` and `memory_region_unref` are never called on aliases or containers, and the above situation then cannot happen. Exploiting this exception is rarely necessary, and therefore it is discouraged, but nevertheless it is used in a few places.

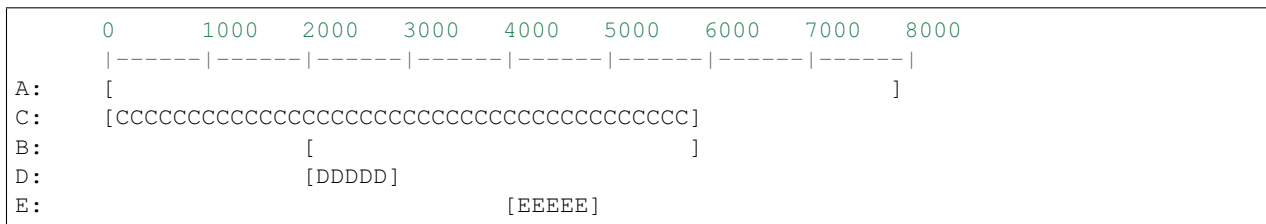
For regions that “have no owner” (NULL is passed at creation time), the machine object is actually used as the owner. Since `instance_finalize` is never called for the machine object, you must never call `object_unparent` on regions that have no owner, unless they are aliases or containers.

6.7.5 Overlapping regions and priority

Usually, regions may not overlap each other; a memory address decodes into exactly one target. In some cases it is useful to allow regions to overlap, and sometimes to control which of an overlapping regions is visible to the guest. This is done with `memory_region_add_subregion_overlap()`, which allows the region to overlap any other region in the same container, and specifies a priority that allows the core to decide which of two regions at the same address are visible (highest wins). Priority values are signed, and the default value is zero. This means that you can use `memory_region_add_subregion_overlap()` both to specify a region that must sit ‘above’ any others (with a positive priority) and also a background region that sits ‘below’ others (with a negative priority).

If the higher priority region in an overlap is a container or alias, then the lower priority region will appear in any “holes” that the higher priority region has left by not mapping subregions to that area of its address range. (This applies recursively – if the subregions are themselves containers or aliases that leave holes then the lower priority region will appear in these holes too.)

For example, suppose we have a container A of size 0x8000 with two subregions B and C. B is a container mapped at 0x2000, size 0x4000, priority 2; C is an MMIO region mapped at 0x0, size 0x6000, priority 1. B currently has two of its own subregions: D of size 0x1000 at offset 0 and E of size 0x1000 at offset 0x2000. As a diagram:



The regions that will be seen within this address range then are:

```
[CCCCCCCCCCCC] [DDDDD] [CCCCC] [EEEEEE] [CCCCC]
```

Since B has higher priority than C, its subregions appear in the flat map even where they overlap with C. In ranges where B has not mapped anything C’s region appears.

If B had provided its own MMIO operations (ie it was not a pure container) then these would be used for any addresses in its range not handled by D or E, and the result would be:

```
[CCCCCCCCCCCC] [DDDDD] [BBBBB] [EEEEEE] [BBBBB]
```

Priority values are local to a container, because the priorities of two regions are only compared when they are both children of the same container. This means that the device in charge of the container (typically modelling a bus or a memory controller) can use them to manage the interaction of its child regions without any side effects on other parts of the system. In the example above, the priorities of D and E are unimportant because they do not overlap each other. It is the relative priority of B and C that causes D and E to appear on top of C: D and E’s priorities are never compared against the priority of C.

6.7.6 Visibility

The memory core uses the following rules to select a memory region when the guest accesses an address:

- all direct subregions of the root region are matched against the address, in descending priority order
 - if the address lies outside the region offset/size, the subregion is discarded
 - if the subregion is a leaf (RAM or MMIO), the search terminates, returning this leaf region
 - if the subregion is a container, the same algorithm is used within the subregion (after the address is adjusted by the subregion offset)
 - if the subregion is an alias, the search is continued at the alias target (after the address is adjusted by the subregion offset and alias offset)
 - if a recursive search within a container or alias subregion does not find a match (because of a “hole” in the container’s coverage of its address range), then if this is a container with its own MMIO or RAM backing the search terminates, returning the container itself. Otherwise we continue with the next subregion in priority order
- if none of the subregions match the address then the search terminates with no match found

6.7.7 Example memory map

```
system_memory: container@0-2^48-1
|
+---- lomem: alias@0-0xdfffffff ---> #ram (0-0xdfffffff)
|
+---- himem: alias@0x10000000-0x1ffffffff ---> #ram (0xe0000000-0xffffffff)
|
+---- vga-window: alias@0xa0000-0xbffff ---> #pci (0xa0000-0xbffff)
|      (prio 1)
|
+---- pci-hole: alias@0xe0000000-0xffffffff ---> #pci (0xe0000000-0xffffffff)

pci (0-2^32-1)
|
+---- vga-area: container@0xa0000-0xbffff
|
|      +---- alias@0x00000-0x7fff ---> #vram (0x010000-0x017fff)
|      |
|      +---- alias@0x08000-0xffff ---> #vram (0x020000-0x027fff)
|
+---- vram: ram@0xe1000000-0xe1fffffff
|
+---- vga-mmio: mmio@0xe2000000-0xe200ffff

ram: ram@0x00000000-0xffffffff
```

This is a (simplified) PC memory map. The 4GB RAM block is mapped into the system address space via two aliases: “lomem” is a 1:1 mapping of the first 3.5GB; “himem” maps the last 0.5GB at address 4GB. This leaves 0.5GB for the so-called PCI hole, that allows a 32-bit PCI bus to exist in a system with 4GB of memory.

The memory controller diverts addresses in the range 640K-768K to the PCI address space. This is modelled using the “vga-window” alias, mapped at a higher priority so it obscures the RAM at the same addresses. The vga window can be removed by programming the memory controller; this is modelled by removing the alias and exposing the RAM underneath.

The pci address space is not a direct child of the system address space, since we only want parts of it to be visible (we accomplish this using aliases). It has two subregions: vga-area models the legacy vga window and is occupied by two 32K memory banks pointing at two sections of the framebuffer. In addition the vram is mapped as a BAR at address e1000000, and an additional BAR containing MMIO registers is mapped after it.

Note that if the guest maps a BAR outside the PCI hole, it would not be visible as the pci-hole alias clips it to a 0.5GB range.

6.7.8 MMIO Operations

MMIO regions are provided with `->read()` and `->write()` callbacks, which are sufficient for most devices. Some devices change behaviour based on the attributes used for the memory transaction, or need to be able to respond that the access should provoke a bus error rather than completing successfully; those devices can use the `->read_with_attrs()` and `->write_with_attrs()` callbacks instead.

In addition various constraints can be supplied to control how these callbacks are called:

- `.valid.min_access_size`, `.valid.max_access_size` define the access sizes (in bytes) which the device accepts; accesses outside this range will have device and bus specific behaviour (ignored, or machine check)
- `.valid.unaligned` specifies that the *device being modelled* supports unaligned accesses; if false, unaligned accesses will invoke the appropriate bus or CPU specific behaviour.
- `.impl.min_access_size`, `.impl.max_access_size` define the access sizes (in bytes) supported by the *implementation*; other access sizes will be emulated using the ones available. For example a 4-byte write will be emulated using four 1-byte writes, if `.impl.max_access_size = 1`.
- `.impl.unaligned` specifies that the *implementation* supports unaligned accesses; if false, unaligned accesses will be emulated by two aligned accesses.

6.7.9 API Reference

struct **MemoryListener**

callbacks structure for updates to the physical memory map

Definition

```
struct MemoryListener {
    void (*begin)(MemoryListener *listener);
    void (*commit)(MemoryListener *listener);
    void (*region_add)(MemoryListener *listener, MemoryRegionSection *section);
    void (*region_del)(MemoryListener *listener, MemoryRegionSection *section);
    void (*region_nop)(MemoryListener *listener, MemoryRegionSection *section);
    void (*log_start)(MemoryListener *listener, MemoryRegionSection *section, int old,
↪int new);
    void (*log_stop)(MemoryListener *listener, MemoryRegionSection *section, int old,
↪int new);
    void (*log_sync)(MemoryListener *listener, MemoryRegionSection *section);
    void (*log_clear)(MemoryListener *listener, MemoryRegionSection *section);
    void (*log_global_start)(MemoryListener *listener);
    void (*log_global_stop)(MemoryListener *listener);
    void (*log_global_after_sync)(MemoryListener *listener);
    void (*eventfd_add)(MemoryListener *listener, MemoryRegionSection *section, bool
↪match_data, uint64_t data, EventNotifier *e);
    void (*eventfd_del)(MemoryListener *listener, MemoryRegionSection *section, bool
↪match_data, uint64_t data, EventNotifier *e);
    void (*coalesced_io_add)(MemoryListener *listener, MemoryRegionSection *section,
↪hwaddr addr, hwaddr len);
```

(continues on next page)

(continued from previous page)

```

void (*coalesced_io_del)(MemoryListener *listener, MemoryRegionSection *section,
↪ hwaddr addr, hwaddr len);
    unsigned priority;
};

```

Members

begin Called at the beginning of an address space update transaction. Followed by calls to *MemoryListener.region_add()*, *MemoryListener.region_del()*, *MemoryListener.region_nop()*, *MemoryListener.log_start()* and *MemoryListener.log_stop()* in increasing address order.

listener: The *MemoryListener*.

commit Called at the end of an address space update transaction, after the last call to *MemoryListener.region_add()*, *MemoryListener.region_del()* or *MemoryListener.region_nop()*, *MemoryListener.log_start()* and *MemoryListener.log_stop()*.

listener: The *MemoryListener*.

region_add Called during an address space update transaction, for a section of the address space that is new in this address space since the last transaction.

listener: The *MemoryListener*. **section:** The new *MemoryRegionSection*.

region_del Called during an address space update transaction, for a section of the address space that has disappeared in the address space since the last transaction.

listener: The *MemoryListener*. **section:** The old *MemoryRegionSection*.

region_nop Called during an address space update transaction, for a section of the address space that is in the same place in the address space as in the last transaction.

listener: The *MemoryListener*. **section:** The *MemoryRegionSection*.

log_start Called during an address space update transaction, after one of *MemoryListener.region_add()*, *MemoryListener.region_del()* or *MemoryListener.region_nop()*, if dirty memory logging clients have become active since the last transaction.

listener: The *MemoryListener*. **section:** The *MemoryRegionSection*. **old:** A bitmap of dirty memory logging clients that were active in the previous transaction. **new:** A bitmap of dirty memory logging clients that are active in the current transaction.

log_stop Called during an address space update transaction, after one of *MemoryListener.region_add()*, *MemoryListener.region_del()* or *MemoryListener.region_nop()* and possibly after *MemoryListener.log_start()*, if dirty memory logging clients have become inactive since the last transaction.

listener: The *MemoryListener*. **section:** The *MemoryRegionSection*. **old:** A bitmap of dirty memory logging clients that were active in the previous transaction. **new:** A bitmap of dirty memory logging clients that are active in the current transaction.

log_sync Called by *memory_region_snapshot_and_clear_dirty()* and *memory_global_dirty_log_sync()*, before accessing QEMU's "official" copy of the dirty memory bitmap for a *MemoryRegionSection*.

listener: The *MemoryListener*. **section:** The *MemoryRegionSection*.

log_clear Called before reading the dirty memory bitmap for a *MemoryRegionSection*.

listener: The *MemoryListener*. **section:** The *MemoryRegionSection*.

log_global_start Called by `memory_global_dirty_log_start()`, which enables the `DIRTY_LOG_MIGRATION` client on all memory regions in the address space. `MemoryListener.log_global_start()` is also called when a `MemoryListener` is added, if global dirty logging is active at that time.

listener: The `MemoryListener`.

log_global_stop Called by `memory_global_dirty_log_stop()`, which disables the `DIRTY_LOG_MIGRATION` client on all memory regions in the address space.

listener: The `MemoryListener`.

log_global_after_sync Called after reading the dirty memory bitmap for any `MemoryRegionSection`.

listener: The `MemoryListener`.

eventfd_add Called during an address space update transaction, for a section of the address space that has had a new ioeventfd registration since the last transaction.

listener: The `MemoryListener`. **section:** The new `MemoryRegionSection`. **match_data:** The **match_data** parameter for the new ioeventfd. **data:** The **data** parameter for the new ioeventfd. **e:** The `EventNotifier` parameter for the new ioeventfd.

eventfd_del Called during an address space update transaction, for a section of the address space that has dropped an ioeventfd registration since the last transaction.

listener: The `MemoryListener`. **section:** The new `MemoryRegionSection`. **match_data:** The **match_data** parameter for the dropped ioeventfd. **data:** The **data** parameter for the dropped ioeventfd. **e:** The `EventNotifier` parameter for the dropped ioeventfd.

coalesced_io_add Called during an address space update transaction, for a section of the address space that has had a new coalesced MMIO range registration since the last transaction.

listener: The `MemoryListener`. **section:** The new `MemoryRegionSection`. **addr:** The starting address for the coalesced MMIO range. **len:** The length of the coalesced MMIO range.

coalesced_io_del Called during an address space update transaction, for a section of the address space that has dropped a coalesced MMIO range since the last transaction.

listener: The `MemoryListener`. **section:** The new `MemoryRegionSection`. **addr:** The starting address for the coalesced MMIO range. **len:** The length of the coalesced MMIO range.

priority Govern the order in which memory listeners are invoked. Lower priorities are invoked earlier for “add” or “start” callbacks, and later for “delete” or “stop” callbacks.

Description

Allows a component to adjust to changes in the guest-visible memory map. Use with `memory_listener_register()` and `memory_listener_unregister()`.

struct **AddressSpace**

describes a mapping of addresses to `MemoryRegion` objects

Definition

```
struct AddressSpace {  
};
```

Members

struct **MemoryRegionSection**

describes a fragment of a `MemoryRegion`

Definition

```

struct MemoryRegionSection {
    Int128 size;
    MemoryRegion *mr;
    FlatView *fv;
    hwaddr offset_within_region;
    hwaddr offset_within_address_space;
    bool readonly;
    bool nonvolatile;
};

```

Members

size the size of the section; will not exceed **mr**'s boundaries

mr the region, or NULL if empty

fv the flat view of the address space the region is mapped in

offset_within_region the beginning of the section, relative to **mr**'s start

offset_within_address_space the address of the first byte of the section relative to the region's address space

readonly writes to this section are ignored

nonvolatile this section is non-volatile

void **memory_region_init** (MemoryRegion **mr*, struct *Object* **owner*, const char **name*, uint64_t *size*)
Initialize a memory region

Parameters

MemoryRegion *mr the MemoryRegion to be initialized

struct Object *owner the object that tracks the region's reference count

const char *name used for debugging; not visible to the user or ABI

uint64_t size size of the region; any subregions beyond this size will be clipped

Description

The region typically acts as a container for other memory regions. Use `memory_region_add_subregion()` to add subregions.

void **memory_region_ref** (MemoryRegion **mr*)
Add 1 to a memory region's reference count

Parameters

MemoryRegion *mr the MemoryRegion

Description

Whenever memory regions are accessed outside the BQL, they need to be preserved against hot-unplug. MemoryRegions actually do not have their own reference count; they piggyback on a QOM object, their "owner". This function adds a reference to the owner.

All MemoryRegions must have an owner if they can disappear, even if the device they belong to operates exclusively under the BQL. This is because the region could be returned at any time by `memory_region_find`, and this is usually under guest control.

void **memory_region_unref** (MemoryRegion **mr*)
Remove 1 to a memory region's reference count

Parameters

MemoryRegion *mr the MemoryRegion

Description

Whenever memory regions are accessed outside the BQL, they need to be preserved against hot-unplug. MemoryRegions actually do not have their own reference count; they piggyback on a QOM object, their “owner”. This function removes a reference to the owner and possibly destroys it.

```
void memory_region_init_io (MemoryRegion *mr, struct Object *owner, const MemoryRegionOps *ops, void *opaque, const char *name, uint64_t size)
```

Initialize an I/O memory region.

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region’s reference count

const MemoryRegionOps *ops a structure containing read and write callbacks to be used when I/O is performed on the region.

void *opaque passed to the read and write callbacks of the **ops** structure.

const char *name used for debugging; not visible to the user or ABI

uint64_t size size of the region.

Description

Accesses into the region will cause the callbacks in **ops** to be called. if **size** is nonzero, subregions will be clipped to **size**.

```
void memory_region_init_ram_nomigrate (MemoryRegion *mr, struct Object *owner, const char *name, uint64_t size, Error **errp)
```

Initialize RAM memory region. Accesses into the region will modify memory directly.

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region’s reference count

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

Error **errp pointer to Error*, to store an error if it happens.

Description

Note that this function does not do anything to cause the data in the RAM memory region to be migrated; that is the responsibility of the caller.

```
void memory_region_init_ram_shared_nomigrate (MemoryRegion *mr, struct Object *owner, const char *name, uint64_t size, bool share, Error **errp)
```

Initialize RAM memory region. Accesses into the region will modify memory directly.

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region’s reference count

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

bool share allow remapping RAM to different addresses

Error **errp pointer to Error*, to store an error if it happens.

Description

Note that this function is similar to `memory_region_init_ram_nomigrate`. The only difference is part of the RAM region can be remapped.

```
void memory_region_init_resizeable_ram(MemoryRegion *mr, struct Object *owner, const
                                     char *name, uint64_t size, uint64_t max_size, void (*re-
                                     sized)(const char*, uint64_t length, void *host), Er-
                                     ror **errp)
```

Initialize memory region with resizeable RAM. Accesses into the region will modify memory directly. Only an initial portion of this RAM is actually used. The used size can change across reboots.

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size used size of the region.

uint64_t max_size max size of the region.

void (*resized)(const char*, uint64_t length, void *host) callback to notify owner about used size change.

Error **errp pointer to Error*, to store an error if it happens.

Description

Note that this function does not do anything to cause the data in the RAM memory region to be migrated; that is the responsibility of the caller.

```
void memory_region_init_ram_from_file(MemoryRegion *mr, struct Object *owner,
                                     const char *name, uint64_t size, uint64_t align,
                                     uint32_t ram_flags, const char *path, bool readonly,
                                     Error **errp)
```

Initialize RAM memory region with a mmap-ed backend.

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

uint64_t align alignment of the region base address; if 0, the default alignment (`getpagesize()`) will be used.

uint32_t ram_flags Memory region features: - `RAM_SHARED`: memory must be mmaped with the `MAP_SHARED` flag - `RAM_PMEM`: the memory is persistent memory Other bits are ignored now.

const char *path the path in which to allocate the RAM.

bool readonly true to open `path` for reading, false for read/write.

Error **errp pointer to Error*, to store an error if it happens.

Description

Note that this function does not do anything to cause the data in the RAM memory region to be migrated; that is the responsibility of the caller.

```
void memory_region_init_ram_from_fd(MemoryRegion *mr, struct Object *owner, const
                                   char *name, uint64_t size, bool share, int fd, ram_addr_t off-
                                   set, Error **errp)
    Initialize RAM memory region with a mmap-ed backend.
```

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name the name of the region.

uint64_t size size of the region.

bool share true if memory must be mmaped with the MAP_SHARED flag

int fd the fd to mmap.

ram_addr_t offset offset within the file referenced by fd

Error **errp pointer to Error*, to store an error if it happens.

Description

Note that this function does not do anything to cause the data in the RAM memory region to be migrated; that is the responsibility of the caller.

```
void memory_region_init_ram_ptr(MemoryRegion *mr, struct Object *owner, const char *name,
                               uint64_t size, void *ptr)
    Initialize RAM memory region from a user-provided pointer. Accesses into the region will modify memory
    directly.
```

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

void *ptr memory to be mapped; must contain at least **size** bytes.

Description

Note that this function does not do anything to cause the data in the RAM memory region to be migrated; that is the responsibility of the caller.

```
void memory_region_init_ram_device_ptr(MemoryRegion *mr, struct Object *owner, const
                                       char *name, uint64_t size, void *ptr)
    Initialize RAM device memory region from a user-provided pointer.
```

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name the name of the region.

uint64_t size size of the region.

void *ptr memory to be mapped; must contain at least **size** bytes.

Description

A RAM device represents a mapping to a physical device, such as to a PCI MMIO BAR of an vfio-pci assigned device. The memory region may be mapped into the VM address space and access to the region will modify memory directly. However, the memory region should not be included in a memory dump (device may not be enabled/mapped at the time of the dump), and operations incompatible with manipulating MMIO should be avoided. Replaces skip_dump flag.

Note that this function does not do anything to cause the data in the RAM memory region to be migrated; that is the responsibility of the caller. (For RAM device memory regions, migrating the contents rarely makes sense.)

void **memory_region_init_alias** (MemoryRegion *mr, struct *Object* *owner, const char *name, MemoryRegion *orig, hwaddr offset, uint64_t size)

Initialize a memory region that aliases all or a part of another memory region.

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name used for debugging; not visible to the user or ABI

MemoryRegion *orig the region to be referenced; **mr** will be equivalent to **orig** between **offset** and **offset + size - 1**.

hwaddr offset start of the section in **orig** to be referenced.

uint64_t size size of the region.

void **memory_region_init_rom_nomigrate** (MemoryRegion *mr, struct *Object* *owner, const char *name, uint64_t size, Error **errp)

Initialize a ROM memory region.

Parameters

MemoryRegion *mr the MemoryRegion to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

Error **errp pointer to Error*, to store an error if it happens.

Description

This has the same effect as calling memory_region_init_ram_nomigrate() and then marking the resulting region read-only with memory_region_set_readonly().

Note that this function does not do anything to cause the data in the RAM side of the memory region to be migrated; that is the responsibility of the caller.

void **memory_region_init_rom_device_nomigrate** (MemoryRegion *mr, struct *Object* *owner, const MemoryRegionOps *ops, void *opaque, const char *name, uint64_t size, Error **errp)

Initialize a ROM memory region. Writes are handled via callbacks.

Parameters

MemoryRegion *mr the `MemoryRegion` to be initialized.

struct Object *owner the object that tracks the region's reference count

const MemoryRegionOps *ops callbacks for write access handling (must not be NULL).

void *opaque passed to the read and write callbacks of the `ops` structure.

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

Error **errp pointer to `Error*`, to store an error if it happens.

Description

Note that this function does not do anything to cause the data in the RAM side of the memory region to be migrated; that is the responsibility of the caller.

void **memory_region_init_iommu** (void *_iommu_mr, size_t instance_size, const char *mrtypename, *Object* *owner, const char *name, uint64_t size)
Initialize a memory region of a custom type that translates addresses

Parameters

void *_iommu_mr the `IOMMUMemoryRegion` to be initialized

size_t instance_size the `IOMMUMemoryRegion` subclass instance size

const char *mrtypename the type name of the `IOMMUMemoryRegion`

Object *owner the object that tracks the region's reference count

const char *name used for debugging; not visible to the user or ABI

uint64_t size size of the region.

Description

An IOMMU region translates addresses and forwards accesses to a target memory region.

The IOMMU implementation must define a subclass of `TYPE_IOMMU_MEMORY_REGION`. `_iommu_mr` should be a pointer to enough memory for an instance of that subclass, `instance_size` is the size of that subclass, and `mrtypename` is its name. This function will initialize `_iommu_mr` as an instance of the subclass, and its methods will then be called to handle accesses to the memory region. See the documentation of `IOMMUMemoryRegionClass` for further details.

void **memory_region_init_ram** (`MemoryRegion` *mr, struct *Object* *owner, const char *name, uint64_t size, `Error **errp`)
Initialize RAM memory region. Accesses into the region will modify memory directly.

Parameters

MemoryRegion *mr the `MemoryRegion` to be initialized

struct Object *owner the object that tracks the region's reference count (must be `TYPE_DEVICE` or a subclass of `TYPE_DEVICE`, or NULL)

const char *name name of the memory region

uint64_t size size of the region in bytes

Error **errp pointer to `Error*`, to store an error if it happens.

Description

This function allocates RAM for a board model or device, and arranges for it to be migrated (by calling `vmstate_register_ram()` if **owner** is a `DeviceState`, or `vmstate_register_ram_global()` if **owner** is `NULL`).

TODO: Currently we restrict **owner** to being either `NULL` (for global RAM regions with no owner) or devices, so that we can give the RAM block a unique name for migration purposes. We should lift this restriction and allow arbitrary Objects. If you pass a non-`NULL` non-device **owner** then we will assert.

```
void memory_region_init_rom (MemoryRegion *mr, struct Object *owner, const char *name,
                             uint64_t size, Error **errp)
```

Initialize a ROM memory region.

Parameters

MemoryRegion *mr the `MemoryRegion` to be initialized.

struct Object *owner the object that tracks the region's reference count

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

Error **errp pointer to `Error*`, to store an error if it happens.

Description

This has the same effect as calling `memory_region_init_ram()` and then marking the resulting region read-only with `memory_region_set_readonly()`. This includes arranging for the contents to be migrated.

TODO: Currently we restrict **owner** to being either `NULL` (for global RAM regions with no owner) or devices, so that we can give the RAM block a unique name for migration purposes. We should lift this restriction and allow arbitrary Objects. If you pass a non-`NULL` non-device **owner** then we will assert.

```
void memory_region_init_rom_device (MemoryRegion *mr, struct Object *owner, const MemoryRegionOps *ops, void *opaque, const char *name,
                                    uint64_t size, Error **errp)
```

Initialize a ROM memory region. Writes are handled via callbacks.

Parameters

MemoryRegion *mr the `MemoryRegion` to be initialized.

struct Object *owner the object that tracks the region's reference count

const MemoryRegionOps *ops callbacks for write access handling (must not be `NULL`).

void *opaque passed to the read and write callbacks of the **ops** structure.

const char *name Region name, becomes part of RAMBlock name used in migration stream must be unique within any device

uint64_t size size of the region.

Error **errp pointer to `Error*`, to store an error if it happens.

Description

This function initializes a memory region backed by RAM for reads and callbacks for writes, and arranges for the RAM backing to be migrated (by calling `vmstate_register_ram()` if **owner** is a `DeviceState`, or `vmstate_register_ram_global()` if **owner** is `NULL`).

TODO: Currently we restrict **owner** to being either `NULL` (for global RAM regions with no owner) or devices, so that we can give the RAM block a unique name for migration purposes. We should lift this restriction and allow arbitrary Objects. If you pass a non-`NULL` non-device **owner** then we will assert.

struct *Object* * **memory_region_owner** (MemoryRegion **mr*)
get a memory region's owner.

Parameters

MemoryRegion *mr the memory region being queried.

uint64_t **memory_region_size** (MemoryRegion **mr*)
get a memory region's size.

Parameters

MemoryRegion *mr the memory region being queried.

bool **memory_region_is_ram** (MemoryRegion **mr*)
check whether a memory region is random access

Parameters

MemoryRegion *mr the memory region being queried

Description

Returns `true` if a memory region is random access.

bool **memory_region_is_ram_device** (MemoryRegion **mr*)
check whether a memory region is a ram device

Parameters

MemoryRegion *mr the memory region being queried

Description

Returns `true` if a memory region is a device backed ram region

bool **memory_region_is_romd** (MemoryRegion **mr*)
check whether a memory region is in ROMD mode

Parameters

MemoryRegion *mr the memory region being queried

Description

Returns `true` if a memory region is a ROM device and currently set to allow direct reads.

IOMMUMemoryRegion * **memory_region_get_iommu** (MemoryRegion **mr*)
check whether a memory region is an iommu

Parameters

MemoryRegion *mr the memory region being queried

Description

Returns pointer to IOMMUMemoryRegion if a memory region is an iommu, otherwise NULL.

IOMMUMemoryRegionClass * **memory_region_get_iommu_class_nocheck** (IOMMUMemoryRegion **iommu_mr*)
returns iommu memory region class if an iommu or NULL if not

Parameters

IOMMUMemoryRegion *iommu_mr the memory region being queried

Description

Returns pointer to IOMMUMemoryRegionClass if a memory region is an iommu, otherwise NULL. This is fast path avoiding QOM checking, use with caution.

uint64_t **memory_region_iommu_get_min_page_size** (IOMMUMemoryRegion **iommu_mr*)
 get minimum supported page size for an iommu

Parameters

IOMMUMemoryRegion *iommu_mr the memory region being queried

Description

Returns minimum supported page size for an iommu.

void **memory_region_notify_iommu** (IOMMUMemoryRegion **iommu_mr*, int *iommu_idx*, IOMMUTLBEvent *event*)
 notify a change in an IOMMU translation entry.

Parameters

IOMMUMemoryRegion *iommu_mr the memory region that was changed

int iommu_idx the IOMMU index for the translation table which has changed

IOMMUTLBEvent event TLB event with the new entry in the IOMMU translation table. The entry replaces all old entries for the same virtual I/O address range.

Note

for any IOMMU implementation, an in-place mapping change should be notified with an UNMAP followed by a MAP.

void **memory_region_notify_iommu_one** (IOMMUNotifier **notifier*, IOMMUTLBEvent **event*)
 notify a change in an IOMMU translation entry to a single notifier

Parameters

IOMMUNotifier *notifier the notifier to be notified

IOMMUTLBEvent *event TLB event with the new entry in the IOMMU translation table. The entry replaces all old entries for the same virtual I/O address range.

Description

This works just like `memory_region_notify_iommu()`, but it only notifies a specific notifier, not all of them.

int **memory_region_register_iommu_notifier** (MemoryRegion **mr*, IOMMUNotifier **n*, Error ***errp*)
 register a notifier for changes to IOMMU translation entries.

Parameters

MemoryRegion *mr the memory region to observe

IOMMUNotifier *n the IOMMUNotifier to be added; the notify callback receives a pointer to an IOMMUTLBEntry as the opaque value; the pointer ceases to be valid on exit from the notifier.

Error **errp pointer to Error*, to store an error if it happens.

Description

Returns 0 on success, or a negative errno otherwise. In particular, -EINVAL indicates that at least one of the attributes of the notifier is not supported (flag/range) by the IOMMU memory region. In case of error the error object must be created.

void **memory_region_iommu_replay** (IOMMUMemoryRegion **iommu_mr*, IOMMUNotifier **n*)
 replay existing IOMMU translations to a notifier with the minimum page granularity returned by `mr->iommu_ops->get_page_size()`.

Parameters

IOMMUMemoryRegion *iommu_mr the memory region to observe

IOMMUNotifier *n the notifier to which to replay iommu mappings

Note

this is not related to record-and-replay functionality.

void **memory_region_unregister_iommu_notifier** (MemoryRegion *mr, IOMMUNotifier *n)
unregister a notifier for changes to IOMMU translation entries.

Parameters

MemoryRegion *mr the memory region which was observed and for which notify_stopped() needs to be called

IOMMUNotifier *n the notifier to be removed.

int **memory_region_iommu_get_attr** (IOMMUMemoryRegion *iommu_mr, enum IOMMUMemoryRegionAttr attr, void *data)
return an IOMMU attr if get_attr() is defined on the IOMMU.

Parameters

IOMMUMemoryRegion *iommu_mr the memory region

enum IOMMUMemoryRegionAttr attr the requested attribute

void *data a pointer to the requested attribute data

Description

Returns 0 on success, or a negative errno otherwise. In particular, -EINVAL indicates that the IOMMU does not support the requested attribute.

int **memory_region_iommu_attrs_to_index** (IOMMUMemoryRegion *iommu_mr, MemTxAttrs attrs)
return the IOMMU index to use for translations with the given memory transaction attributes.

Parameters

IOMMUMemoryRegion *iommu_mr the memory region

MemTxAttrs attrs the memory transaction attributes

int **memory_region_iommu_num_indexes** (IOMMUMemoryRegion *iommu_mr)
return the total number of IOMMU indexes that this IOMMU supports.

Parameters

IOMMUMemoryRegion *iommu_mr the memory region

int **memory_region_iommu_set_page_size_mask** (IOMMUMemoryRegion *iommu_mr, uint64_t page_size_mask, Error **errp)
set the supported page sizes for a given IOMMU memory region

Parameters

IOMMUMemoryRegion *iommu_mr IOMMU memory region

uint64_t page_size_mask supported page size mask

Error **errp pointer to Error*, to store an error if it happens.

const char * **memory_region_name** (const MemoryRegion *mr)
get a memory region's name

Parameters

const MemoryRegion *mr the memory region being queried

Description

Returns the string that was used to initialize the memory region.

bool **memory_region_is_logging** (MemoryRegion **mr*, uint8_t *client*)
 return whether a memory region is logging writes

Parameters

MemoryRegion *mr the memory region being queried

uint8_t client the client being queried

Description

Returns `true` if the memory region is logging writes for the given client

uint8_t **memory_region_get_dirty_log_mask** (MemoryRegion **mr*)
 return the clients for which a memory region is logging writes.

Parameters

MemoryRegion *mr the memory region being queried

Description

Returns a bitmap of clients, in which the `DIRTY_MEMORY_*` constants are the bit indices.

bool **memory_region_is_rom** (MemoryRegion **mr*)
 check whether a memory region is ROM

Parameters

MemoryRegion *mr the memory region being queried

Description

Returns `true` if a memory region is read-only memory.

bool **memory_region_is_nonvolatile** (MemoryRegion **mr*)
 check whether a memory region is non-volatile

Parameters

MemoryRegion *mr the memory region being queried

Description

Returns `true` is a memory region is non-volatile memory.

int **memory_region_get_fd** (MemoryRegion **mr*)
 Get a file descriptor backing a RAM memory region.

Parameters

MemoryRegion *mr the RAM or alias memory region being queried.

Description

Returns a file descriptor backing a file-based RAM memory region, or -1 if the region is not a file-based RAM memory region.

MemoryRegion * **memory_region_from_host** (void **ptr*, ram_addr_t **offset*)
 Convert a pointer into a RAM memory region and an offset within it.

Parameters

void *ptr the host pointer to be converted

ram_addr_t *offset the offset within memory region

Description

Given a host pointer inside a RAM memory region (created with `memory_region_init_ram()` or `memory_region_init_ram_ptr()`), return the `MemoryRegion` and the offset within it.

Use with care; by the time this function returns, the returned pointer is not protected by RCU anymore. If the caller is not within an RCU critical section and does not hold the `iothread` lock, it must have other means of protecting the pointer, such as a reference to the region that includes the incoming `ram_addr_t`.

void ***memory_region_get_ram_ptr** (`MemoryRegion *mr`)
Get a pointer into a RAM memory region.

Parameters

MemoryRegion *mr the memory region being queried.

Description

Returns a host pointer to a RAM memory region (created with `memory_region_init_ram()` or `memory_region_init_ram_ptr()`).

Use with care; by the time this function returns, the returned pointer is not protected by RCU anymore. If the caller is not within an RCU critical section and does not hold the `iothread` lock, it must have other means of protecting the pointer, such as a reference to the region that includes the incoming `ram_addr_t`.

void **memory_region_msync** (`MemoryRegion *mr`, `hwaddr addr`, `hwaddr size`)
Synchronize selected address range of a memory mapped region

Parameters

MemoryRegion *mr the memory region to be msync

hwaddr addr the initial address of the range to be sync

hwaddr size the size of the range to be sync

void **memory_region_writeback** (`MemoryRegion *mr`, `hwaddr addr`, `hwaddr size`)
Trigger cache writeback for selected address range

Parameters

MemoryRegion *mr the memory region to be updated

hwaddr addr the initial address of the range to be written back

hwaddr size the size of the range to be written back

void **memory_region_set_log** (`MemoryRegion *mr`, `bool log`, `unsigned client`)
Turn dirty logging on or off for a region.

Parameters

MemoryRegion *mr the memory region being updated.

bool log whether dirty logging is to be enabled or disabled.

unsigned client the user of the logging information; `DIRTY_MEMORY_VGA` only.

Description

Turns dirty logging on or off for a specified client (display, migration). Only meaningful for RAM regions.

void **memory_region_set_dirty** (`MemoryRegion *mr`, `hwaddr addr`, `hwaddr size`)
Mark a range of bytes as dirty in a memory region.

Parameters

MemoryRegion *mr the memory region being dirtied.

hwaddr addr the address (relative to the start of the region) being dirtied.

hwaddr size size of the range being dirtied.

Description

Marks a range of bytes as dirty, after it has been dirtied outside guest code.

void **memory_region_clear_dirty_bitmap** (MemoryRegion *mr, hwaddr start, hwaddr len)
clear dirty bitmap for memory range

Parameters

MemoryRegion *mr the memory region to clear the dirty log upon

hwaddr start start address offset within the memory region

hwaddr len length of the memory region to clear dirty bitmap

Description

This function is called when the caller wants to clear the remote dirty bitmap of a memory range within the memory region. This can be used by e.g. KVM to manually clear dirty log when KVM_CAP_MANUAL_DIRTY_LOG_PROTECT is declared support by the host kernel.

DirtyBitmapSnapshot * **memory_region_snapshot_and_clear_dirty** (MemoryRegion *mr,
hwaddr addr, hwaddr size,
unsigned client)

Get a snapshot of the dirty bitmap and clear it.

Parameters

MemoryRegion *mr the memory region being queried.

hwaddr addr the address (relative to the start of the region) being queried.

hwaddr size the size of the range being queried.

unsigned client the user of the logging information; typically DIRTY_MEMORY_VGA.

Description

Creates a snapshot of the dirty bitmap, clears the dirty bitmap and returns the snapshot. The snapshot can then be used to query dirty status, using memory_region_snapshot_get_dirty. Snapshotting allows querying the same page multiple times, which is especially useful for display updates where the scanlines often are not page aligned.

The dirty bitmap region which gets copied into the snapshot (and cleared afterwards) can be larger than requested. The boundaries are rounded up/down so complete bitmap longs (covering 64 pages on 64bit hosts) can be copied over into the bitmap snapshot. Which isn't a problem for display updates as the extra pages are outside the visible area, and in case the visible area changes a full display redraw is due anyway. Should other use cases for this function emerge we might have to revisit this implementation detail.

Use g_free to release DirtyBitmapSnapshot.

bool **memory_region_snapshot_get_dirty** (MemoryRegion *mr, DirtyBitmapSnapshot *snap,
hwaddr addr, hwaddr size)
Check whether a range of bytes is dirty in the specified dirty bitmap snapshot.

Parameters

MemoryRegion *mr the memory region being queried.

DirtyBitmapSnapshot *snap the dirty bitmap snapshot

hwaddr addr the address (relative to the start of the region) being queried.

hwaddr size the size of the range being queried.

void **memory_region_reset_dirty** (MemoryRegion **mr*, hwaddr *addr*, hwaddr *size*, unsigned *client*)
Mark a range of pages as clean, for a specified client.

Parameters

MemoryRegion *mr the region being updated.

hwaddr addr the start of the subrange being cleaned.

hwaddr size the size of the subrange being cleaned.

unsigned client the user of the logging information; DIRTY_MEMORY_MIGRATION or DIRTY_MEMORY_VGA.

Description

Marks a range of pages as no longer dirty.

void **memory_region_flush_rom_device** (MemoryRegion **mr*, hwaddr *addr*, hwaddr *size*)
Mark a range of pages dirty and invalidate TBs (for self-modifying code).

Parameters

MemoryRegion *mr the region being flushed.

hwaddr addr the start, relative to the start of the region, of the range being flushed.

hwaddr size the size, in bytes, of the range being flushed.

Description

The MemoryRegionOps->write() callback of a ROM device must use this function to mark byte ranges that have been modified internally, such as by directly accessing the memory returned by memory_region_get_ram_ptr().

This function marks the range dirty and invalidates TBs so that TCG can detect self-modifying code.

void **memory_region_set_readonly** (MemoryRegion **mr*, bool *readonly*)
Turn a memory region read-only (or read-write)

Parameters

MemoryRegion *mr the region being updated.

bool readonly whether the region is to be ROM or RAM.

Description

Allows a memory region to be marked as read-only (turning it into a ROM). only useful on RAM regions.

void **memory_region_set_nonvolatile** (MemoryRegion **mr*, bool *nonvolatile*)
Turn a memory region non-volatile

Parameters

MemoryRegion *mr the region being updated.

bool nonvolatile whether the region is to be non-volatile.

Description

Allows a memory region to be marked as non-volatile. only useful on RAM regions.

void **memory_region_rom_device_set_romd** (MemoryRegion **mr*, bool *romd_mode*)
enable/disable ROMD mode

Parameters

MemoryRegion *mr the memory region to be updated

bool romd_mode `true` to put the region into ROMD mode

Description

Allows a ROM device (initialized with `memory_region_init_rom_device()` to set to ROMD mode (default) or MMIO mode. When it is in ROMD mode, the device is mapped to guest memory and satisfies read access directly. When in MMIO mode, reads are forwarded to the `MemoryRegion.read` function. Writes are always handled by the `MemoryRegion.write` function.

void **memory_region_set_coalescing** (MemoryRegion *mr)
 Enable memory coalescing for the region.

Parameters

MemoryRegion *mr the memory region to be write coalesced

Description

Enabled writes to a region to be queued for later processing. MMIO ->write callbacks may be delayed until a non-coalesced MMIO is issued. Only useful for IO regions. Roughly similar to write-combining hardware.

void **memory_region_add_coalescing** (MemoryRegion *mr, hwaddr offset, uint64_t size)
 Enable memory coalescing for a sub-range of a region.

Parameters

MemoryRegion *mr the memory region to be updated.

hwaddr offset the start of the range within the region to be coalesced.

uint64_t size the size of the subrange to be coalesced.

Description

Like `memory_region_set_coalescing()`, but works on a sub-range of a region. Multiple calls can be issued coalesced disjoint ranges.

void **memory_region_clear_coalescing** (MemoryRegion *mr)
 Disable MMIO coalescing for the region.

Parameters

MemoryRegion *mr the memory region to be updated.

Description

Disables any coalescing caused by `memory_region_set_coalescing()` or `memory_region_add_coalescing()`. Roughly equivalent to uncacheble memory hardware.

void **memory_region_set_flush_coalesced** (MemoryRegion *mr)
 Enforce memory coalescing flush before accesses.

Parameters

MemoryRegion *mr the memory region to be updated.

Description

Ensure that pending coalesced MMIO request are flushed before the memory region is accessed. This property is automatically enabled for all regions passed to `memory_region_set_coalescing()` and `memory_region_add_coalescing()`.

void **memory_region_clear_flush_coalesced** (MemoryRegion *mr)
 Disable memory coalescing flush before accesses.

Parameters

MemoryRegion *mr the memory region to be updated.

Description

Clear the automatic coalesced MMIO flushing enabled via `memory_region_set_flush_coalesced`. Note that this service has no effect on memory regions that have MMIO coalescing enabled for themselves. For them, automatic flushing will stop once coalescing is disabled.

void **memory_region_add_eventfd** (MemoryRegion *mr, hwaddr addr, unsigned size, bool match_data, uint64_t data, EventNotifier *e)

Request an eventfd to be triggered when a word is written to a location.

Parameters

MemoryRegion *mr the memory region being updated.

hwaddr addr the address within **mr** that is to be monitored

unsigned size the size of the access to trigger the eventfd

bool match_data whether to match against **data**, instead of just **addr**

uint64_t data the data to match against the guest write

EventNotifier *e event notifier to be triggered when **addr**, **size**, and **data** all match.

Description

Marks a word in an IO region (initialized with `memory_region_init_io()`) as a trigger for an eventfd event. The I/O callback will not be called. The caller must be prepared to handle failure (that is, take the required action if the callback `_is_called`).

void **memory_region_del_eventfd** (MemoryRegion *mr, hwaddr addr, unsigned size, bool match_data, uint64_t data, EventNotifier *e)

Cancel an eventfd.

Parameters

MemoryRegion *mr the memory region being updated.

hwaddr addr the address within **mr** that is to be monitored

unsigned size the size of the access to trigger the eventfd

bool match_data whether to match against **data**, instead of just **addr**

uint64_t data the data to match against the guest write

EventNotifier *e event notifier to be triggered when **addr**, **size**, and **data** all match.

Description

Cancels an eventfd trigger requested by a previous `memory_region_add_eventfd()` call.

void **memory_region_add_subregion** (MemoryRegion *mr, hwaddr offset, MemoryRegion *subregion)

Add a subregion to a container.

Parameters

MemoryRegion *mr the region to contain the new subregion; must be a container initialized with `memory_region_init()`.

hwaddr offset the offset relative to **mr** where **subregion** is added.

MemoryRegion *subregion the subregion to be added.

Description

Adds a subregion at **offset**. The subregion may not overlap with other subregions (except for those explicitly marked as overlapping). A region may only be added once as a subregion (unless removed with `memory_region_del_subregion()`); use `memory_region_init_alias()` if you want a region to be a subregion in multiple locations.

```
void memory_region_add_subregion_overlap (MemoryRegion *mr, hwaddr offset, MemoryRegion *subregion, int priority)
```

Add a subregion to a container with overlap.

Parameters

MemoryRegion *mr the region to contain the new subregion; must be a container initialized with `memory_region_init()`.

hwaddr offset the offset relative to **mr** where **subregion** is added.

MemoryRegion *subregion the subregion to be added.

int priority used for resolving overlaps; highest priority wins.

Description

Adds a subregion at **offset**. The subregion may overlap with other subregions. Conflicts are resolved by having a higher **priority** hide a lower **priority**. Subregions without priority are taken as **priority 0**. A region may only be added once as a subregion (unless removed with `memory_region_del_subregion()`); use `memory_region_init_alias()` if you want a region to be a subregion in multiple locations.

```
ram_addr_t memory_region_get_ram_addr (MemoryRegion *mr)
```

Get the ram address associated with a memory region

Parameters

MemoryRegion *mr the region to be queried

```
void memory_region_del_subregion (MemoryRegion *mr, MemoryRegion *subregion)
```

Remove a subregion.

Parameters

MemoryRegion *mr the container to be updated.

MemoryRegion *subregion the region being removed; must be a current subregion of **mr**.

Description

Removes a subregion from its container.

```
bool memory_region_present (MemoryRegion *container, hwaddr addr)
```

checks if an address relative to a **container** translates into `MemoryRegion` within **container**

Parameters

MemoryRegion *container a `MemoryRegion` within which **addr** is a relative address

hwaddr addr the area within **container** to be searched

Description

Answer whether a `MemoryRegion` within **container** covers the address **addr**.

```
bool memory_region_is_mapped (MemoryRegion *mr)
```

returns true if `MemoryRegion` is mapped into any address space.

Parameters

MemoryRegion *mr a `MemoryRegion` which should be checked if it's mapped

MemoryRegionSection **memory_region_find** (MemoryRegion **mr*, hwaddr *addr*, uint64_t *size*)
translate an address/size relative to a MemoryRegion into a *MemoryRegionSection*.

Parameters

MemoryRegion *mr a MemoryRegion within which **addr** is a relative address

hwaddr addr start of the area within **as** to be searched

uint64_t size size of the area to be searched

Description

Locates the first MemoryRegion within **mr** that overlaps the range given by **addr** and **size**.

Returns a *MemoryRegionSection* that describes a contiguous overlap. It will have the following characteristics:
- **size** = 0 iff no overlap was found - **mr** is non-NULL iff an overlap was found

Remember that in the return value the **offset_within_region** is relative to the returned region (in the **mr** field), not to the **mr** argument.

Similarly, the **offset_within_address_space** is relative to the address space that contains both regions, the passed and the returned one. However, in the special case where the **mr** argument has no container (and thus is the root of the address space), the following will hold: - **offset_within_address_space** >= **addr** - **offset_within_address_space** + **size** <= **addr** + **size**

void **memory_global_dirty_log_sync** (void)
synchronize the dirty log for all memory

Parameters

void no arguments

Description

Synchronizes the dirty page log for all address spaces.

void **memory_global_after_dirty_log_sync** (void)
synchronize the dirty log for all memory

Parameters

void no arguments

Description

Synchronizes the vCPUs with a thread that is reading the dirty bitmap. This function must be called after the dirty log bitmap is cleared, and before dirty guest memory pages are read. If you are using DirtyBitmapSnapshot, **memory_region_snapshot_and_clear_dirty()** takes care of doing this.

void **memory_region_transaction_begin** (void)
Start a transaction.

Parameters

void no arguments

Description

During a transaction, changes will be accumulated and made visible only when the transaction ends (is committed).

void **memory_region_transaction_commit** (void)
Commit a transaction and make changes visible to the guest.

Parameters

void no arguments

void **memory_listener_register** (*MemoryListener* *listener, *AddressSpace* *filter)
 register callbacks to be called when memory sections are mapped or unmapped into an address space

Parameters

MemoryListener *listener an object containing the callbacks to be called

AddressSpace *filter if non-NULL, only regions in this address space will be observed

void **memory_listener_unregister** (*MemoryListener* *listener)
 undo the effect of memory_listener_register()

Parameters

MemoryListener *listener an object containing the callbacks to be removed

void **memory_global_dirty_log_start** (void)
 begin dirty logging for all regions

Parameters

void no arguments

void **memory_global_dirty_log_stop** (void)
 end dirty logging for all regions

Parameters

void no arguments

MemTxResult **memory_region_dispatch_read** (MemoryRegion *mr, hwaddr addr, uint64_t *pval,
 MemOp op, MemTxAttrs attrs)
 perform a read directly to the specified MemoryRegion.

Parameters

MemoryRegion *mr MemoryRegion to access

hwaddr addr address within that region

uint64_t *pval pointer to uint64_t which the data is written to

MemOp op size, sign, and endianness of the memory operation

MemTxAttrs attrs memory transaction attributes to use for the access

MemTxResult **memory_region_dispatch_write** (MemoryRegion *mr, hwaddr addr, uint64_t data,
 MemOp op, MemTxAttrs attrs)
 perform a write directly to the specified MemoryRegion.

Parameters

MemoryRegion *mr MemoryRegion to access

hwaddr addr address within that region

uint64_t data data to write

MemOp op size, sign, and endianness of the memory operation

MemTxAttrs attrs memory transaction attributes to use for the access

void **address_space_init** (*AddressSpace* *as, MemoryRegion *root, const char *name)
 initializes an address space

Parameters

AddressSpace *as an uninitialized *AddressSpace*

MemoryRegion *root a `MemoryRegion` that routes addresses for the address space
const char *name an address space name. The name is only used for debugging output.
void **address_space_destroy** (*AddressSpace *as*)
destroy an address space

Parameters

AddressSpace *as address space to be destroyed

Description

Releases all resources associated with an address space. After an address space is destroyed, its root memory region (given by `address_space_init()`) may be destroyed as well.

void **address_space_remove_listeners** (*AddressSpace *as*)
unregister all listeners of an address space

Parameters

AddressSpace *as an initialized *AddressSpace*

Description

Removes all callbacks previously registered with `memory_listener_register()` for **as**.

MemTxResult **address_space_rw** (*AddressSpace *as*, hwaddr *addr*, MemTxAttrs *attrs*, void **buf*,
hwaddr *len*, bool *is_write*)
read from or write to an address space.

Parameters

AddressSpace *as *AddressSpace* to be accessed

hwaddr addr address within that address space

MemTxAttrs attrs memory transaction attributes

void *buf buffer with the data transferred

hwaddr len the number of bytes to read or write

bool is_write indicates the transfer direction

Description

Return a `MemTxResult` indicating whether the operation succeeded or failed (eg unassigned memory, device rejected the transaction, IOMMU fault).

MemTxResult **address_space_write** (*AddressSpace *as*, hwaddr *addr*, MemTxAttrs *attrs*, const
void **buf*, hwaddr *len*)
write to address space.

Parameters

AddressSpace *as *AddressSpace* to be accessed

hwaddr addr address within that address space

MemTxAttrs attrs memory transaction attributes

const void *buf buffer with the data transferred

hwaddr len the number of bytes to write

Description

Return a MemTxResult indicating whether the operation succeeded or failed (eg unassigned memory, device rejected the transaction, IOMMU fault).

MemTxResult **address_space_write_rom** (*AddressSpace* *as, hwaddr addr, MemTxAttrs attrs, const void *buf, hwaddr len)
 write to address space, including ROM.

Parameters

AddressSpace *as *AddressSpace* to be accessed

hwaddr addr address within that address space

MemTxAttrs attrs memory transaction attributes

const void *buf buffer with the data transferred

hwaddr len the number of bytes to write

Description

This function writes to the specified address space, but will write data to both ROM and RAM. This is used for non-guest writes like writes from the gdb debug stub or initial loading of ROM contents.

Note that portions of the write which attempt to write data to a device will be silently ignored – only real RAM and ROM will be written to.

Return a MemTxResult indicating whether the operation succeeded or failed (eg unassigned memory, device rejected the transaction, IOMMU fault).

void **address_space_cache_invalidate** (MemoryRegionCache *cache, hwaddr addr, hwaddr access_len)
 complete a write to a MemoryRegionCache

Parameters

MemoryRegionCache *cache The MemoryRegionCache to operate on.

hwaddr addr The first physical address that was written, relative to the address that was passed to **address_space_cache_init**.

hwaddr access_len The number of bytes that were written starting at **addr**.

void **address_space_cache_destroy** (MemoryRegionCache *cache)
 free a MemoryRegionCache

Parameters

MemoryRegionCache *cache The MemoryRegionCache whose memory should be released.

MemTxResult **address_space_read** (*AddressSpace* *as, hwaddr addr, MemTxAttrs attrs, void *buf, hwaddr len)
 read from an address space.

Parameters

AddressSpace *as *AddressSpace* to be accessed

hwaddr addr address within that address space

MemTxAttrs attrs memory transaction attributes

void *buf buffer with the data transferred

hwaddr len length of the data transferred

Description

Return a MemTxResult indicating whether the operation succeeded or failed (eg unassigned memory, device rejected the transaction, IOMMU fault). Called within RCU critical section.

MemTxResult **address_space_read_cached** (MemoryRegionCache **cache*, hwaddr *addr*, void **buf*, hwaddr *len*)
read from a cached RAM region

Parameters

MemoryRegionCache *cache Cached region to be addressed

hwaddr addr address relative to the base of the RAM region

void *buf buffer with the data transferred

hwaddr len length of the data transferred

MemTxResult **address_space_write_cached** (MemoryRegionCache **cache*, hwaddr *addr*, const void **buf*, hwaddr *len*)
write to a cached RAM region

Parameters

MemoryRegionCache *cache Cached region to be addressed

hwaddr addr address relative to the base of the RAM region

const void *buf buffer with the data transferred

hwaddr len length of the data transferred

6.8 Migration

QEMU has code to load/save the state of the guest that it is running. These are two complementary operations. Saving the state just does that, saves the state for each device that the guest is running. Restoring a guest is just the opposite operation: we need to load the state of each device.

For this to work, QEMU has to be launched with the same arguments the two times. I.e. it can only restore the state in one guest that has the same devices that the one it was saved (this last requirement can be relaxed a bit, but for now we can consider that configuration has to be exactly the same).

Once that we are able to save/restore a guest, a new functionality is requested: migration. This means that QEMU is able to start in one machine and being “migrated” to another machine. I.e. being moved to another machine.

Next was the “live migration” functionality. This is important because some guests run with a lot of state (specially RAM), and it can take a while to move all state from one machine to another. Live migration allows the guest to continue running while the state is transferred. Only while the last part of the state is transferred has the guest to be stopped. Typically the time that the guest is unresponsive during live migration is the low hundred of milliseconds (notice that this depends on a lot of things).

6.8.1 Transports

The migration stream is normally just a byte stream that can be passed over any transport.

- tcp migration: do the migration using tcp sockets
- unix migration: do the migration using unix sockets
- exec migration: do the migration using the stdin/stdout through a process.

- fd migration: do the migration using a file descriptor that is passed to QEMU. QEMU doesn't care how this file descriptor is opened.

In addition, support is included for migration using RDMA, which transports the page data using RDMA, where the hardware takes care of transporting the pages, and the load on the CPU is much lower. While the internals of RDMA migration are a bit different, this isn't really visible outside the RAM migration code.

All these migration protocols use the same infrastructure to save/restore state devices. This infrastructure is shared with the savevm/loadvm functionality.

6.8.2 Debugging

The migration stream can be analyzed thanks to `scripts/analyze-migration.py`.

Example usage:

```
$ qemu-system-x86_64 -display none -monitor stdio
(qemu) migrate "exec:cat > mig"
(qemu) q
$ ./scripts/analyze-migration.py -f mig
{
  "ram (3)": {
    "section sizes": {
      "pc.ram": "0x0000000008000000",
      ...
    }
  }
}
```

See also `analyze-migration.py -h` help for more options.

6.8.3 Common infrastructure

The files, sockets or fd's that carry the migration stream are abstracted by the `QEMUFile` type (see `migration/qemu-file.h`). In most cases this is connected to a subtype of `QIOChannel` (see `io/`).

6.8.4 Saving the state of one device

For most devices, the state is saved in a single call to the migration infrastructure; these are *non-iterative* devices. The data for these devices is sent at the end of precopy migration, when the CPUs are paused. There are also *iterative* devices, which contain a very large amount of data (e.g. RAM or large tables). See the iterative device section below.

General advice for device developers

- The migration state saved should reflect the device being modelled rather than the way your implementation works. That way if you change the implementation later the migration stream will stay compatible. That model may include internal state that's not directly visible in a register.
- When saving a migration stream the device code may walk and check the state of the device. These checks might fail in various ways (e.g. discovering internal state is corrupt or that the guest has done something bad). Consider carefully before asserting/aborting at this point, since the normal response from users is that *migration broke their VM* since it had apparently been running fine until then. In these error cases, the device should log a message indicating the cause of error, and should consider putting the device into an error state, allowing the rest of the VM to continue execution.

- The migration might happen at an inconvenient point, e.g. right in the middle of the guest reprogramming the device, during guest reboot or shutdown or while the device is waiting for external IO. It's strongly preferred that migrations do not fail in this situation, since in the cloud environment migrations might happen automatically to VMs that the administrator doesn't directly control.
- If you do need to fail a migration, ensure that sufficient information is logged to identify what went wrong.
- The destination should treat an incoming migration stream as hostile (which we do to varying degrees in the existing code). Check that offsets into buffers and the like can't cause overruns. Fail the incoming migration in the case of a corrupted stream like this.
- Take care with internal device state or behaviour that might become migration version dependent. For example, the order of PCI capabilities is required to stay constant across migration. Another example would be that a special case handled by subsections (see below) might become much more common if a default behaviour is changed.
- The state of the source should not be changed or destroyed by the outgoing migration. Migrations timing out or being failed by higher levels of management, or failures of the destination host are not unusual, and in that case the VM is restarted on the source. Note that the management layer can validly revert the migration even though the QEMU level of migration has succeeded as long as it does it before starting execution on the destination.
- Buses and devices should be able to explicitly specify addresses when instantiated, and management tools should use those. For example, when hot adding USB devices it's important to specify the ports and addresses, since implicit ordering based on the command line order may be different on the destination. This can result in the device state being loaded into the wrong device.

VMState

Most device data can be described using the `VMSTATE` macros (mostly defined in `include/migration/vmstate.h`).

An example (from `hw/input/pckbd.c`)

```
static const VMStateDescription vmstate_kbd = {
    .name = "pckbd",
    .version_id = 3,
    .minimum_version_id = 3,
    .fields = (VMStateField[]) {
        VMSTATE_UINT8(write_cmd, KBDState),
        VMSTATE_UINT8(status, KBDState),
        VMSTATE_UINT8(mode, KBDState),
        VMSTATE_UINT8(pending, KBDState),
        VMSTATE_END_OF_LIST()
    }
};
```

We are declaring the state with name "pckbd". The `version_id` is 3, and the fields are 4 `uint8_t` in a `KBDState` structure. We registered this with:

```
vmstate_register(NULL, 0, &vmstate_kbd, s);
```

For devices that are *qdev* based, we can register the device in the class init function:

```
dc->vmsd = &vmstate_kbd_isa;
```

The `VMState` macros take care of ensuring that the device data section is formatted portably (normally big endian) and make some compile time checks against the types of the fields in the structures.

VMState macros can include other VMStateDescriptions to store substructures (see `VMSTATE_STRUCT__`), arrays (`VMSTATE_ARRAY__`) and variable length arrays (`VMSTATE_VARRAY__`). Various other macros exist for special cases.

Note that the format on the wire is still very raw; i.e. a `VMSTATE_UINT32` ends up with a 4 byte bigendian representation on the wire; in the future it might be possible to use a more structured format.

Legacy way

This way is going to disappear as soon as all current users are ported to VMSTATE; although converting existing code can be tricky, and thus ‘soon’ is relative.

Each device has to register two functions, one to save the state and another to load the state back.

```
int register_savevm_live(const char *idstr,
                        int instance_id,
                        int version_id,
                        SaveVMHandlers *ops,
                        void *opaque);
```

Two functions in the `ops` structure are the `save_state` and `load_state` functions. Notice that `load_state` receives a `version_id` parameter to know what state format is receiving. `save_state` doesn’t have a `version_id` parameter because it always uses the latest version.

Note that because the VMState macros still save the data in a raw format, in many cases it’s possible to replace legacy code with a carefully constructed VMState description that matches the byte layout of the existing code.

Changing migration data structures

When we migrate a device, we save/load the state as a series of fields. Sometimes, due to bugs or new functionality, we need to change the state to store more/different information. Changing the migration state saved for a device can break migration compatibility unless care is taken to use the appropriate techniques. In general QEMU tries to maintain forward migration compatibility (i.e. migrating from QEMU $n \rightarrow n+1$) and there are users who benefit from backward compatibility as well.

Subsections

The most common structure change is adding new data, e.g. when adding a newer form of device, or adding that state that you previously forgot to migrate. This is best solved using a subsection.

A subsection is “like” a device vmstate, but with a particularity, it has a Boolean function that tells if that values are needed to be sent or not. If this functions returns false, the subsection is not sent. Subsections have a unique name, that is looked for on the receiving side.

On the receiving side, if we found a subsection for a device that we don’t understand, we just fail the migration. If we understand all the subsections, then we load the state with success. There’s no check that a subsection is loaded, so a newer QEMU that knows about a subsection can (with care) load a stream from an older QEMU that didn’t send the subsection.

If the new data is only needed in a rare case, then the subsection can be made conditional on that case and the migration will still succeed to older QEMUs in most cases. This is OK for data that’s critical, but in some use cases it’s preferred that the migration should succeed even with the data missing. To support this the subsection can be connected to a device property and from there to a versioned machine type.

The ‘pre_load’ and ‘post_load’ functions on subsections are only called if the subsection is loaded.

One important note is that the outer `post_load()` function is called “after” loading all subsections, because a newer subsection could change the same value that it uses. A flag, and the combination of outer `pre_load` and `post_load` can be used to detect whether a subsection was loaded, and to fall back on default behaviour when the subsection isn’t present.

Example:

```
static bool ide_drive_pio_state_needed(void *opaque)
{
    IDEState *s = opaque;

    return ((s->status & DRQ_STAT) != 0)
        || (s->bus->error_status & BM_STATUS_PIO_RETRY);
}

const VMStateDescription vmstate_ide_drive_pio_state = {
    .name = "ide_drive/pio_state",
    .version_id = 1,
    .minimum_version_id = 1,
    .pre_save = ide_drive_pio_pre_save,
    .post_load = ide_drive_pio_post_load,
    .needed = ide_drive_pio_state_needed,
    .fields = (VMStateField[]) {
        VMSTATE_INT32(req_nb_sectors, IDEState),
        VMSTATE_VARRAY_INT32(io_buffer, IDEState, io_buffer_total_len, 1,
                             vmstate_info_uint8, uint8_t),
        VMSTATE_INT32(cur_io_buffer_offset, IDEState),
        VMSTATE_INT32(cur_io_buffer_len, IDEState),
        VMSTATE_UINT8(end_transfer_fn_idx, IDEState),
        VMSTATE_INT32(elementary_transfer_size, IDEState),
        VMSTATE_INT32(packet_transfer_size, IDEState),
        VMSTATE_END_OF_LIST()
    }
};

const VMStateDescription vmstate_ide_drive = {
    .name = "ide_drive",
    .version_id = 3,
    .minimum_version_id = 0,
    .post_load = ide_drive_post_load,
    .fields = (VMStateField[]) {
        .... several fields ....
        VMSTATE_END_OF_LIST()
    },
    .subsections = (const VMStateDescription*[]) {
        &vmstate_ide_drive_pio_state,
        NULL
    }
};
```

Here we have a subsection for the pio state. We only need to save/send this state when we are in the middle of a pio operation (that is what `ide_drive_pio_state_needed()` checks). If `DRQ_STAT` is not enabled, the values on that fields are garbage and don’t need to be sent.

Connecting subsections to properties

Using a condition function that checks a ‘property’ to determine whether to send a subsection allows backward migration compatibility when new subsections are added, especially when combined with versioned machine types.

For example:

- a) Add a new property using `DEFINE_PROP_BOOL` - e.g. `support-foo` and default it to `true`.
- b) Add an entry to the `hw_compat_` for the previous version that sets the property to `false`.
- c) Add a static bool `support_foo` function that tests the property.
- d) Add a subsection with a `.needed` set to the `support_foo` function
- e) (potentially) Add an outer `pre_load` that sets up a default value for 'foo' to be used if the subsection isn't loaded.

Now that subsection will not be generated when using an older machine type and the migration stream will be accepted by older QEMU versions.

Not sending existing elements

Sometimes members of the VMState are no longer needed:

- removing them will break migration compatibility
- making them version dependent and bumping the version will break backward migration compatibility.

Adding a dummy field into the migration stream is normally the best way to preserve compatibility.

If the field really does need to be removed then:

- a) Add a new property/compatibility/function in the same way for subsections above.
- b) replace the `VMSTATE` macro with the `_TEST` version of the macro, e.g.:

```
VMSTATE_UINT32(foo, barstruct)
```

becomes

```
VMSTATE_UINT32_TEST(foo, barstruct, pre_version_baz)
```

Sometime in the future when we no longer care about the ancient versions these can be killed off. Note that for backward compatibility it's important to fill in the structure with data that the destination will understand.

Any difference in the predicates on the source and destination will end up with different fields being enabled and data being loaded into the wrong fields; for this reason conditional fields like this are very fragile.

Versions

Version numbers are intended for major incompatible changes to the migration of a device, and using them breaks backward-migration compatibility; in general most changes can be made by adding Subsections (see above) or `_TEST` macros (see above) which won't break compatibility.

Each version is associated with a series of fields saved. The `save_state` always saves the state as the newer version. But `load_state` sometimes is able to load state from an older version.

You can see that there are several version fields:

- `version_id`: the maximum `version_id` supported by VMState for that device.
- `minimum_version_id`: the minimum `version_id` that VMState is able to understand for that device.
- `minimum_version_id_old`: For devices that were not able to port to vmstate, we can assign a function that knows how to read this old state. This field is ignored if there is no `load_state_old` handler.

VMState is able to read versions from `minimum_version_id` to `version_id`. And the function `load_state_old()` (if present) is able to load state from `minimum_version_id_old` to `minimum_version_id`. This function is deprecated and will be removed when no more users are left.

There are `_V` forms of many `VMSTATE_` macros to load fields for version dependent fields, e.g.

```
VMSTATE_UINT16_V(ip_id, Slirp, 2),
```

only loads that field for versions 2 and newer.

Saving state will always create a section with the `'version_id'` value and thus can't be loaded by any older QEMU.

Massaging functions

Sometimes, it is not enough to be able to save the state directly from one structure, we need to fill the correct values there. One example is when we are using kvm. Before saving the cpu state, we need to ask kvm to copy to QEMU the state that it is using. And the opposite when we are loading the state, we need a way to tell kvm to load the state for the cpu that we have just loaded from the QEMUFile.

The functions to do that are inside a `vmstate` definition, and are called:

- `int (*pre_load)(void *opaque);`

This function is called before we load the state of one device.

- `int (*post_load)(void *opaque, int version_id);`

This function is called after we load the state of one device.

- `int (*pre_save)(void *opaque);`

This function is called before we save the state of one device.

- `int (*post_save)(void *opaque);`

This function is called after we save the state of one device (even upon failure, unless the call to `pre_save` returned an error).

Example: You can look at `hpet.c`, that uses the first three functions to massage the state that is transferred.

The `VMSTATE_WITH_TMP` macro may be useful when the migration data doesn't match the stored device data well; it allows an intermediate temporary structure to be populated with migration data and then transferred to the main structure.

If you use memory API functions that update memory layout outside initialization (i.e., in response to a guest action), this is a strong indication that you need to call these functions in a *post_load* callback. Examples of such memory API functions are:

- `memory_region_add_subregion()`
- `memory_region_del_subregion()`
- `memory_region_set_readonly()`
- `memory_region_set_nonvolatile()`
- `memory_region_set_enabled()`
- `memory_region_set_address()`
- `memory_region_set_alias_offset()`

Iterative device migration

Some devices, such as RAM, Block storage or certain platform devices, have large amounts of data that would mean that the CPUs would be paused for too long if they were sent in one section. For these devices an *iterative* approach is taken.

The iterative devices generally don't use VMState macros (although it may be possible in some cases) and instead use `qemu_put_*/qemu_get_*` macros to read/write data to the stream. Specialist versions exist for high bandwidth IO.

An iterative device must provide:

- A `save_setup` function that initialises the data structures and transmits a first section containing information on the device. In the case of RAM this transmits a list of RAMBlocks and sizes.
- A `load_setup` function that initialises the data structures on the destination.
- A `save_live_pending` function that is called repeatedly and must indicate how much more data the iterative data must save. The core migration code will use this to determine when to pause the CPUs and complete the migration.
- A `save_live_iterate` function (called after `save_live_pending` when there is significant data still to be sent). It should send a chunk of data until the point that stream bandwidth limits tell it to stop. Each call generates one section.
- A `save_live_complete_precopy` function that must transmit the last section for the device containing any remaining data.
- A `load_state` function used to load sections generated by any of the save functions that generate sections.
- `cleanup` functions for both save and load that are called at the end of migration.

Note that the contents of the sections for iterative migration tend to be open-coded by the devices; care should be taken in parsing the results and structuring the stream to make them easy to validate.

Device ordering

There are cases in which the ordering of device loading matters; for example in some systems where a device may assert an interrupt during loading, if the interrupt controller is loaded later then it might lose the state.

Some ordering is implicitly provided by the order in which the machine definition creates devices, however this is somewhat fragile.

The `MigrationPriority` enum provides a means of explicitly enforcing ordering. Numerically higher priorities are loaded earlier. The priority is set by setting the `priority` field of the top level `VMStateDescription` for the device.

6.8.5 Stream structure

The stream tries to be word and endian agnostic, allowing migration between hosts of different characteristics running the same VM.

- Header
 - Magic
 - Version
 - VM configuration section
 - * Machine type

* Target page bits

- List of sections Each section contains a device, or one iteration of a device save.
 - section type
 - section id
 - ID string (First section of each device)
 - instance id (First section of each device)
 - version id (First section of each device)
 - <device data>
 - Footer mark
- EOF mark
- VM Description structure Consisting of a JSON description of the contents for analysis only

The `device data` in each section consists of the data produced by the code described above. For non-iterative devices they have a single section; iterative devices have an initial and last section and a set of parts in between. Note that there is very little checking by the common code of the integrity of the `device data` contents, that's up to the devices themselves. The `footer mark` provides a little bit of protection for the case where the receiving side reads more or less data than expected.

The `ID string` is normally unique, having been formed from a bus name and device address, PCI devices and storage devices hung off PCI controllers fit this pattern well. Some devices are fixed single instances (e.g. "pc-ram"). Others (especially either older devices or system devices which for some reason don't have a bus concept) make use of the `instance id` for otherwise identically named devices.

Return path

Only a unidirectional stream is required for normal migration, however a `return path` can be created when bidirectional communication is desired. This is primarily used by postcopy, but is also used to return a success flag to the source at the end of migration.

`qemu_file_get_return_path(QEMUFile* fwdpath)` gives the `QEMUFile*` for the return path.

Source side

Forward path - written by migration thread Return path - opened by main thread, read by return-path thread

Destination side

Forward path - read by main thread Return path - opened by main thread, written by main thread AND postcopy thread (protected by `rp_mutex`)

6.8.6 Postcopy

'Postcopy' migration is a way to deal with migrations that refuse to converge (or take too long to converge) its plus side is that there is an upper bound on the amount of migration traffic and time it takes, the down side is that during the postcopy phase, a failure of *either* side or the network connection causes the guest to be lost.

In postcopy the destination CPUs are started before all the memory has been transferred, and accesses to pages that are yet to be transferred cause a fault that's translated by QEMU into a request to the source QEMU.

Postcopy can be combined with precopy (i.e. normal migration) so that if precopy doesn't finish in a given time the switch is made to postcopy.

Enabling postcopy

To enable postcopy, issue this command on the monitor (both source and destination) prior to the start of migration:

```
migrate_set_capability postcopy-ram on
```

The normal commands are then used to start a migration, which is still started in precopy mode. Issuing:

```
migrate_start_postcopy
```

will now cause the transition from precopy to postcopy. It can be issued immediately after migration is started or any time later on. Issuing it after the end of a migration is harmless.

Blocktime is a postcopy live migration metric, intended to show how long the vCPU was in state of interruptible sleep due to pagefault. That metric is calculated both for all vCPUs as overlapped value, and separately for each vCPU. These values are calculated on destination side. To enable postcopy blocktime calculation, enter following command on destination monitor:

```
migrate_set_capability postcopy-blocktime on
```

Postcopy blocktime can be retrieved by query-migrate qmp command. postcopy-blocktime value of qmp command will show overlapped blocking time for all vCPU, postcopy-vcpu-blocktime will show list of blocking time per vCPU.

Note: During the postcopy phase, the bandwidth limits set using `migrate_set_speed` is ignored (to avoid delaying requested pages that the destination is waiting for).

Postcopy device transfer

Loading of device data may cause the device emulation to access guest RAM that may trigger faults that have to be resolved by the source, as such the migration stream has to be able to respond with page data *during* the device load, and hence the device data has to be read from the stream completely before the device load begins to free the stream up. This is achieved by ‘packaging’ the device data into a blob that’s read in one go.

Source behaviour

Until postcopy is entered the migration stream is identical to normal precopy, except for the addition of a ‘postcopy advise’ command at the beginning, to tell the destination that postcopy might happen. When postcopy starts the source sends the page discard data and then forms the ‘package’ containing:

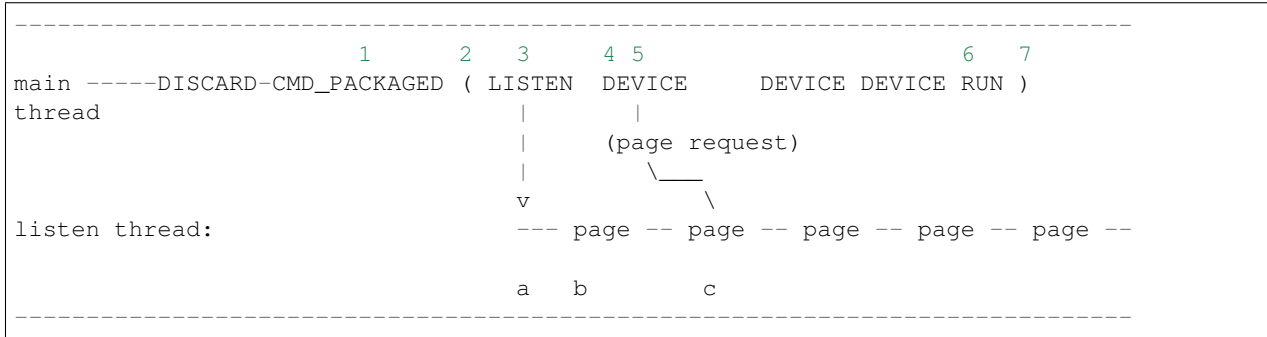
- Command: ‘postcopy listen’
- The device state
 - A series of sections, identical to the precopy streams device state stream containing everything except postcopyable devices (i.e. RAM)
- Command: ‘postcopy run’

The ‘package’ is sent as the data part of a Command: `CMD_PACKAGED`, and the contents are formatted in the same way as the main migration stream.

During postcopy the source scans the list of dirty pages and sends them to the destination without being requested (in much the same way as precopy), however when a page request is received from the destination, the dirty page scanning restarts from the requested location. This causes requested pages to be sent quickly, and also causes pages directly after the requested page to be sent quickly in the hope that those pages are likely to be used by the destination soon.

Destination behaviour

Initially the destination looks the same as precopy, with a single thread reading the migration stream; the ‘postcopy advise’ and ‘discard’ commands are processed to change the way RAM is managed, but don’t affect the stream processing.



- On receipt of CMD_PACKAGED (1)

All the data associated with the package - the (...) section in the diagram - is read into memory, and the main thread recurses into qemu_loadvm_state_main to process the contents of the package (2) which contains commands (3,6) and devices (4...)
- On receipt of ‘postcopy listen’ - 3 -(i.e. the 1st command in the package)

a new thread (a) is started that takes over servicing the migration stream, while the main thread carries on loading the package. It loads normal background page data (b) but if during a device load a fault happens (5) the returned page (c) is loaded by the listen thread allowing the main threads device load to carry on.
- The last thing in the CMD_PACKAGED is a ‘RUN’ command (6)

letting the destination CPUs start running. At the end of the CMD_PACKAGED (7) the main thread returns to normal running behaviour and is no longer used by migration, while the listen thread carries on servicing page data until the end of migration.

Postcopy states

Postcopy moves through a series of states (see postcopy_state) from ADVISE->DISCARD->LISTEN->RUNNING->END

- Advise

Set at the start of migration if postcopy is enabled, even if it hasn’t had the start command; here the destination checks that its OS has the support needed for postcopy, and performs setup to ensure the RAM mappings are suitable for later postcopy. The destination will fail early in migration at this point if the required OS support is not present. (Triggered by reception of POSTCOPY_ADVISE command)
- Discard

Entered on receipt of the first ‘discard’ command; prior to the first Discard being performed, hugepages are switched off (using madvise) to ensure that no new huge pages are created during the postcopy phase, and to cause any huge pages that have discards on them to be broken.
- Listen

The first command in the package, `POSTCOPY_LISTEN`, switches the destination state to Listen, and starts a new thread (the ‘listen thread’) which takes over the job of receiving pages off the migration stream, while the main thread carries on processing the blob. With this thread able to process page reception, the destination now ‘sensitises’ the RAM to detect any access to missing pages (on Linux using the ‘userfault’ system).

- Running

`POSTCOPY_RUN` causes the destination to synchronise all state and start the CPUs and IO devices running. The main thread now finishes processing the migration package and now carries on as it would for normal precopy migration (although it can’t do the cleanup it would do as it finishes a normal migration).

- End

The listen thread can now quit, and perform the cleanup of migration state, the migration is now complete.

Source side page maps

The source side keeps two bitmaps during postcopy; ‘the migration bitmap’ and ‘unsent map’. The ‘migration bitmap’ is basically the same as in the precopy case, and holds a bit to indicate that page is ‘dirty’ - i.e. needs sending. During the precopy phase this is updated as the CPU dirties pages, however during postcopy the CPUs are stopped and nothing should dirty anything any more.

The ‘unsent map’ is used for the transition to postcopy. It is a bitmap that has a bit cleared whenever a page is sent to the destination, however during the transition to postcopy mode it is combined with the migration bitmap to form a set of pages that:

- Have been sent but then redirtied (which must be discarded)
- Have not yet been sent - which also must be discarded to cause any transparent huge pages built during precopy to be broken.

Note that the contents of the unsentmap are sacrificed during the calculation of the discard set and thus aren’t valid once in postcopy. The dirtymap is still valid and is used to ensure that no page is sent more than once. Any request for a page that has already been sent is ignored. Duplicate requests such as this can happen as a page is sent at about the same time the destination accesses it.

Postcopy with hugepages

Postcopy now works with hugetlbfs backed memory:

- The linux kernel on the destination must support userfault on hugepages.
- The huge-page configuration on the source and destination VMs must be identical; i.e. RAMBlocks on both sides must use the same page size.
- Note that `-mem-path /dev/hugepages` will fall back to allocating normal RAM if it doesn’t have enough hugepages, triggering (b) to fail. Using `-mem-prealloc` enforces the allocation using hugepages.
- Care should be taken with the size of hugepage used; postcopy with 2MB hugepages works well, however 1GB hugepages are likely to be problematic since it takes ~1 second to transfer a 1GB hugepage across a 10Gbps link, and until the full page is transferred the destination thread is blocked.

Postcopy with shared memory

Postcopy migration with shared memory needs explicit support from the other processes that share memory and from QEMU. There are restrictions on the type of memory that userfault can support shared.

The Linux kernel userfault support works on `/dev/shm` memory and on `hugetlbfs` (although the kernel doesn't provide an equivalent to `madvise(MADV_DONTNEED)` for `hugetlbfs` which may be a problem in some configurations).

The `vhost-user` code in QEMU supports clients that have Postcopy support, and the `vhost-user-bridge` (in `tests/`) and the DPDK package have changes to support postcopy.

The client needs to open a `userfaultfd` and register the areas of memory that it maps with `userfault`. The client must then pass the `userfaultfd` back to QEMU together with a mapping table that allows fault addresses in the clients address space to be converted back to `RAMBlock`/offsets. The client's `userfaultfd` is added to the postcopy fault-thread and page requests are made on behalf of the client by QEMU. QEMU performs 'wake' operations on the client's `userfaultfd` to allow it to continue after a page has arrived.

Note:

There are two future improvements that would be nice:

- a) Some way to make QEMU ignorant of the addresses in the clients address space
 - b) Avoiding the need for QEMU to perform `ufd-wake` calls after the pages have arrived
-

Retro-fitting postcopy to existing clients is possible:

- a) A mechanism is needed for the registration with `userfault` as above, and the registration needs to be co-ordinated with the phases of postcopy. In `vhost-user` extra messages are added to the existing control channel.
- b) Any thread that can block due to guest memory accesses must be identified and the implication understood; for example if the guest memory access is made while holding a lock then all other threads waiting for that lock will also be blocked.

6.8.7 Firmware

Migration migrates the copies of RAM and ROM, and thus when running on the destination it includes the firmware from the source. Even after resetting a VM, the old firmware is used. Only once QEMU has been restarted is the new firmware in use.

- Changes in firmware size can cause changes in the required `RAMBlock` size to hold the firmware and thus migration can fail. In practice it's best to pad firmware images to convenient powers of 2 with plenty of space for growth.
- Care should be taken with device emulation code so that newer emulation code can work with older firmware to allow forward migration.
- Care should be taken with newer firmware so that backward migration to older systems with older device emulation code will work.

In some cases it may be best to tie specific firmware versions to specific versioned machine types to cut down on the combinations that will need support. This is also useful when newer versions of firmware outgrow the padding.

6.9 Atomic operations in QEMU

CPUs perform independent memory operations effectively in random order. but this can be a problem for CPU-CPU interaction (including interactions between QEMU and the guest). Multi-threaded programs use various tools to instruct the compiler and the CPU to restrict the order to something that is consistent with the expectations of the programmer.

The most basic tool is locking. Mutexes, condition variables and semaphores are used in QEMU, and should be the default approach to synchronization. Anything else is considerably harder, but it's also justified more often than one would like; the most performance-critical parts of QEMU in particular require a very low level approach to concurrency, involving memory barriers and atomic operations. The semantics of concurrent memory accesses are governed by the C11 memory model.

QEMU provides a header, `qemu/atomic.h`, which wraps C11 atomics to provide better portability and a less verbose syntax. `qemu/atomic.h` provides macros that fall in three camps:

- compiler barriers: `barrier()`;
- weak atomic access and manual memory barriers: `qatomic_read()`, `qatomic_set()`, `smp_rmb()`, `smp_wmb()`, `smp_mb()`, `smp_mb_acquire()`, `smp_mb_release()`, `smp_read_barrier_depends()`;
- sequentially consistent atomic access: everything else.

In general, use of `qemu/atomic.h` should be wrapped with more easily used data structures (e.g. the lock-free singly-linked list operations `QSLIST_INSERT_HEAD_ATOMIC` and `QSLIST_MOVE_ATOMIC`) or synchronization primitives (such as `RCU`, `QemuEvent` or `QemuLockCnt`). Bare use of atomic operations and memory barriers should be limited to inter-thread checking of flags and documented thoroughly.

6.9.1 Compiler memory barrier

`barrier()` prevents the compiler from moving the memory accesses on either side of it to the other side. The compiler barrier has no direct effect on the CPU, which may then reorder things however it wishes.

`barrier()` is mostly used within `qemu/atomic.h` itself. On some architectures, CPU guarantees are strong enough that blocking compiler optimizations already ensures the correct order of execution. In this case, `qemu/atomic.h` will reduce stronger memory barriers to simple compiler barriers.

Still, `barrier()` can be useful when writing code that can be interrupted by signal handlers.

6.9.2 Sequentially consistent atomic access

Most of the operations in the `qemu/atomic.h` header ensure *sequential consistency*, where “the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program”.

`qemu/atomic.h` provides the following set of atomic read-modify-write operations:

```
void qatomic_inc(ptr)
void qatomic_dec(ptr)
void qatomic_add(ptr, val)
void qatomic_sub(ptr, val)
void qatomic_and(ptr, val)
void qatomic_or(ptr, val)

typedef(*ptr) qatomic_fetch_inc(ptr)
```

(continues on next page)

(continued from previous page)

```

typeof(*ptr) qatomic_fetch_dec(ptr)
typeof(*ptr) qatomic_fetch_add(ptr, val)
typeof(*ptr) qatomic_fetch_sub(ptr, val)
typeof(*ptr) qatomic_fetch_and(ptr, val)
typeof(*ptr) qatomic_fetch_or(ptr, val)
typeof(*ptr) qatomic_fetch_xor(ptr, val)
typeof(*ptr) qatomic_fetch_inc_nonzero(ptr)
typeof(*ptr) qatomic_xchg(ptr, val)
typeof(*ptr) qatomic_cmpxchg(ptr, old, new)

```

all of which return the old value of `*ptr`. These operations are polymorphic; they operate on any type that is as wide as a pointer or smaller.

Similar operations return the new value of `*ptr`:

```

typeof(*ptr) qatomic_inc_fetch(ptr)
typeof(*ptr) qatomic_dec_fetch(ptr)
typeof(*ptr) qatomic_add_fetch(ptr, val)
typeof(*ptr) qatomic_sub_fetch(ptr, val)
typeof(*ptr) qatomic_and_fetch(ptr, val)
typeof(*ptr) qatomic_or_fetch(ptr, val)
typeof(*ptr) qatomic_xor_fetch(ptr, val)

```

`qemu/atomic.h` also provides loads and stores that cannot be reordered with each other:

```

typeof(*ptr) qatomic_mb_read(ptr)
void          qatomic_mb_set(ptr, val)

```

However these do not provide sequential consistency and, in particular, they do not participate in the total ordering enforced by sequentially-consistent operations. For this reason they are deprecated. They should instead be replaced with any of the following (ordered from easiest to hardest):

- accesses inside a mutex or spinlock
- lightweight synchronization primitives such as `QemuEvent`
- RCU operations (`qatomic_rcu_read`, `qatomic_rcu_set`) when publishing or accessing a new version of a data structure
- other atomic accesses: `qatomic_read` and `qatomic_load_acquire` for loads, `qatomic_set` and `qatomic_store_release` for stores, `smp_mb` to forbid reordering subsequent loads before a store.

6.9.3 Weak atomic access and manual memory barriers

Compared to sequentially consistent atomic access, programming with weaker consistency models can be considerably more complicated. The only guarantees that you can rely upon in this case are:

- atomic accesses will not cause data races (and hence undefined behavior); ordinary accesses instead cause data races if they are concurrent with other accesses of which at least one is a write. In order to ensure this, the compiler will not optimize accesses out of existence, create unsolicited accesses, or perform other similar optimizations.
- acquire operations will appear to happen, with respect to the other components of the system, before all the LOAD or STORE operations specified afterwards.
- release operations will appear to happen, with respect to the other components of the system, after all the LOAD or STORE operations specified before.

- release operations will *synchronize with* acquire operations; see *Acquire/release pairing and the synchronizes-with relation* for a detailed explanation.

When using this model, variables are accessed with:

- `qatomic_read()` and `qatomic_set()`; these prevent the compiler from optimizing accesses out of existence and creating unsolicited accesses, but do not otherwise impose any ordering on loads and stores: both the compiler and the processor are free to reorder them.
- `qatomic_load_acquire()`, which guarantees the LOAD to appear to happen, with respect to the other components of the system, before all the LOAD or STORE operations specified afterwards. Operations coming before `qatomic_load_acquire()` can still be reordered after it.
- `qatomic_store_release()`, which guarantees the STORE to appear to happen, with respect to the other components of the system, after all the LOAD or STORE operations specified before. Operations coming after `qatomic_store_release()` can still be reordered before it.

Restrictions to the ordering of accesses can also be specified using the memory barrier macros: `smp_rmb()`, `smp_wmb()`, `smp_mb()`, `smp_mb_acquire()`, `smp_mb_release()`, `smp_read_barrier_depends()`.

Memory barriers control the order of references to shared memory. They come in six kinds:

- `smp_rmb()` guarantees that all the LOAD operations specified before the barrier will appear to happen before all the LOAD operations specified after the barrier with respect to the other components of the system.

In other words, `smp_rmb()` puts a partial ordering on loads, but is not required to have any effect on stores.

- `smp_wmb()` guarantees that all the STORE operations specified before the barrier will appear to happen before all the STORE operations specified after the barrier with respect to the other components of the system.

In other words, `smp_wmb()` puts a partial ordering on stores, but is not required to have any effect on loads.

- `smp_mb_acquire()` guarantees that all the LOAD operations specified before the barrier will appear to happen before all the LOAD or STORE operations specified after the barrier with respect to the other components of the system.
- `smp_mb_release()` guarantees that all the STORE operations specified *after* the barrier will appear to happen after all the LOAD or STORE operations specified *before* the barrier with respect to the other components of the system.
- `smp_mb()` guarantees that all the LOAD and STORE operations specified before the barrier will appear to happen before all the LOAD and STORE operations specified after the barrier with respect to the other components of the system.

`smp_mb()` puts a partial ordering on both loads and stores. It is stronger than both a read and a write memory barrier; it implies both `smp_mb_acquire()` and `smp_mb_release()`, but it also prevents STOREs coming before the barrier from overtaking LOADs coming after the barrier and vice versa.

- `smp_read_barrier_depends()` is a weaker kind of read barrier. On most processors, whenever two loads are performed such that the second depends on the result of the first (e.g., the first load retrieves the address to which the second load will be directed), the processor will guarantee that the first LOAD will appear to happen before the second with respect to the other components of the system. However, this is not always true—for example, it was not true on Alpha processors. Whenever this kind of access happens to shared memory (that is not protected by a lock), a read barrier is needed, and `smp_read_barrier_depends()` can be used instead of `smp_rmb()`.

Note that the first load really has to have a `_data_` dependency and not a control dependency. If the address for the second load is dependent on the first load, but the dependency is through a conditional rather than actually loading the address itself, then it's a `_control_` dependency and a full read barrier or better is required.

Memory barriers and `qatomic_load_acquire/qatomic_store_release` are mostly used when a data structure has one thread that is always a writer and one thread that is always a reader:

thread 1	thread 2
<pre>qatomic_store_release(&a, x); qatomic_store_release(&b, y);</pre>	<pre>y = qatomic_load_acquire(&b); x = qatomic_load_acquire(&a);</pre>

In this case, correctness is easy to check for using the “pairing” trick that is explained below.

Sometimes, a thread is accessing many variables that are otherwise unrelated to each other (for example because, apart from the current thread, exactly one other thread will read or write each of these variables). In this case, it is possible to “hoist” the barriers outside a loop. For example:

before	after
<pre>n = 0; for (i = 0; i < 10; i++) n += qatomic_load_acquire(&a[i]);</pre>	<pre>n = 0; for (i = 0; i < 10; i++) n += qatomic_read(&a[i]); smp_mb_acquire();</pre>
<pre>for (i = 0; i < 10; i++) qatomic_store_release(&a[i], ↪false);</pre>	<pre>smp_mb_release(); for (i = 0; i < 10; i++) qatomic_set(&a[i], false);</pre>

Splitting a loop can also be useful to reduce the number of barriers:

before	after
<pre>n = 0; for (i = 0; i < 10; i++) { qatomic_store_release(&a[i], ↪false); smp_mb(); n += qatomic_read(&b[i]); }</pre>	<pre>smp_mb_release(); for (i = 0; i < 10; i++) qatomic_set(&a[i], false); smb_mb(); n = 0; for (i = 0; i < 10; i++) n += qatomic_read(&b[i]);</pre>

In this case, a `smp_mb_release()` is also replaced with a (possibly cheaper, and clearer as well) `smp_wmb()`:

before	after
<pre>for (i = 0; i < 10; i++) { qatomic_store_release(&a[i], ↪false); qatomic_store_release(&b[i], ↪false); }</pre>	<pre>smp_mb_release(); for (i = 0; i < 10; i++) qatomic_set(&a[i], false); smb_wmb(); for (i = 0; i < 10; i++) qatomic_set(&b[i], false);</pre>

Acquire/release pairing and the *synchronizes-with* relation

Atomic operations other than `qatomic_set()` and `qatomic_read()` have either *acquire* or *release* semantics¹. This has two effects:

- within a thread, they are ordered either before subsequent operations (for acquire) or after previous operations (for release).
- if a release operation in one thread *synchronizes with* an acquire operation in another thread, the ordering constraints propagates from the first to the second thread. That is, everything before the release operation in the first thread is guaranteed to *happen before* everything after the acquire operation in the second thread.

The concept of acquire and release semantics is not exclusive to atomic operations; almost all higher-level synchronization primitives also have acquire or release semantics. For example:

- `pthread_mutex_lock` has acquire semantics, `pthread_mutex_unlock` has release semantics and synchronizes with a `pthread_mutex_lock` for the same mutex.
- `pthread_cond_signal` and `pthread_cond_broadcast` have release semantics; `pthread_cond_wait` has both release semantics (synchronizing with `pthread_mutex_lock`) and acquire semantics (synchronizing with `pthread_mutex_unlock` and signaling of the condition variable).
- `pthread_create` has release semantics and synchronizes with the start of the new thread; `pthread_join` has acquire semantics and synchronizes with the exiting of the thread.
- `qemu_event_set` has release semantics, `qemu_event_wait` has acquire semantics.

For example, in the following example there are no atomic accesses, but still thread 2 is relying on the *synchronizes-with* relation between `pthread_exit` (release) and `pthread_join` (acquire):

thread 1	thread 2
<pre>*a = 1; pthread_exit(a);</pre>	<pre>pthread_join(thread1, &a); x = *a;</pre>

Synchronization between threads basically descends from this pairing of a release operation and an acquire operation. Therefore, atomic operations other than `qatomic_set()` and `qatomic_read()` will almost always be paired with another operation of the opposite kind: an acquire operation will pair with a release operation and vice versa. This rule of thumb is extremely useful; in the case of QEMU, however, note that the other operation may actually be in a driver that runs in the guest!

`smp_read_barrier_depends()`, `smp_rmb()`, `smp_mb_acquire()`, `qatomic_load_acquire()` and `qatomic_rcu_read()` all count as acquire operations. `smp_wmb()`, `smp_mb_release()`, `qatomic_store_release()` and `qatomic_rcu_set()` all count as release operations. `smp_mb()` counts as both acquire and release, therefore it can pair with any other atomic operation. Here is an example:

thread 1	thread 2
<pre>qatomic_set(&a, 1); smp_wmb(); qatomic_set(&b, 2);</pre>	<pre>x = qatomic_read(&b); smp_rmb(); y = qatomic_read(&a);</pre>

¹ Read-modify-write operations can have both—acquire applies to the read part, and release to the write.

Note that a load-store pair only counts if the two operations access the same variable: that is, a store-release on a variable `x` *synchronizes with* a load-acquire on a variable `x`, while a release barrier synchronizes with any acquire operation. The following example shows correct synchronization:

thread 1	thread 2
<pre>qatomic_set(&a, 1); qatomic_store_release(&b, 2);</pre>	<pre>x = qatomic_load_acquire(&b); y = qatomic_read(&a);</pre>

Acquire and release semantics of higher-level primitives can also be relied upon for the purpose of establishing the *synchronizes with* relation.

Note that the “writing” thread is accessing the variables in the opposite order as the “reading” thread. This is expected: stores before a release operation will normally match the loads after the acquire operation, and vice versa. In fact, this happened already in the `pthread_exit/pthread_join` example above.

Finally, this more complex example has more than two accesses and data dependency barriers. It also does not use atomic accesses whenever there cannot be a data race:

thread 1	thread 2
<pre>b[2] = 1; smp_wmb(); x->i = 2; smp_wmb(); qatomic_set(&a, x);</pre>	<pre>x = qatomic_read(&a); smp_read_barrier_depends(); y = x->i; smp_read_barrier_depends(); z = b[y];</pre>

6.9.4 Comparison with Linux kernel primitives

Here is a list of differences between Linux kernel atomic operations and memory barriers, and the equivalents in QEMU:

- atomic operations in Linux are always on a 32-bit int type and use a boxed `atomic_t` type; atomic operations in QEMU are polymorphic and use normal C types.
- Originally, `atomic_read` and `atomic_set` in Linux gave no guarantee at all. Linux 4.1 updated them to implement volatile semantics via `ACCESS_ONCE` (or the more recent `READ/WRITE_ONCE`).

QEMU’s `qatomic_read` and `qatomic_set` implement C11 atomic relaxed semantics if the compiler supports it, and volatile semantics otherwise. Both semantics prevent the compiler from doing certain transformations; the difference is that atomic accesses are guaranteed to be atomic, while volatile accesses aren’t. Thus, in the volatile case we just cross our fingers hoping that the compiler will generate atomic accesses, since we assume the variables passed are machine-word sized and properly aligned.

No barriers are implied by `qatomic_read` and `qatomic_set` in either Linux or QEMU.

- atomic read-modify-write operations in Linux are of three kinds:

<code>atomic_OP</code>	returns void
<code>atomic_OP_return</code>	returns new value of the variable
<code>atomic_fetch_OP</code>	returns the old value of the variable
<code>atomic_cmpxchg</code>	returns the old value of the variable

In QEMU, the second kind is named `atomic_OP_fetch`.

- different atomic read-modify-write operations in Linux imply a different set of memory barriers; in QEMU, all of them enforce sequential consistency.
- in QEMU, `qatomic_read()` and `qatomic_set()` do not participate in the total ordering enforced by sequentially-consistent operations. This is because QEMU uses the C11 memory model. The following example is correct in Linux but not in QEMU:

Linux (correct)	QEMU (incorrect)
<pre>a = atomic_fetch_add(&x, 2); b = READ_ONCE(&y);</pre>	<pre>a = qatomic_fetch_add(&x, 2); b = qatomic_read(&y);</pre>

because the read of `y` can be moved (by either the processor or the compiler) before the write of `x`.

Fixing this requires an `smp_mb()` memory barrier between the write of `x` and the read of `y`. In the common case where only one thread writes `x`, it is also possible to write it like this:

QEMU (correct)
<pre>a = qatomic_read(&x); qatomic_set(&x, a + 2); smp_mb(); b = qatomic_read(&y);</pre>

6.9.5 Sources

- `Documentation/memory-barriers.txt` from the Linux kernel

6.10 QEMU and the stable process

6.10.1 QEMU stable releases

QEMU stable releases are based upon the last released QEMU version and marked by an additional version number, e.g. 2.10.1. Occasionally, a four-number version is released, if a single urgent fix needs to go on top.

Usually, stable releases are only provided for the last major QEMU release. For example, when QEMU 2.11.0 is released, 2.11.x or 2.11.x.y stable releases are produced only until QEMU 2.12.0 is released, at which point the stable process moves to producing 2.12.x/2.12.x.y releases.

6.10.2 What should go into a stable release?

Generally, the following patches are considered stable material:

- Patches that fix severe issues, like fixes for CVEs
- Patches that fix regressions

If you think the patch would be important for users of the current release (or for a distribution picking fixes), it is usually a good candidate for stable.

6.10.3 How to get a patch into QEMU stable

There are various ways to get a patch into stable:

- Preferred: Make sure that the stable maintainers are on copy when you send the patch by adding

```
Cc: qemu-stable@nongnu.org
```

to the patch description. By default, this will send a copy of the patch to `qemu-stable@nongnu.org` if you use `git send-email`, which is where patches that are stable candidates are tracked by the maintainers.

- You can also reply to a patch and put `qemu-stable@nongnu.org` on copy directly in your mail client if you think a previously submitted patch should be considered for a stable release.
- If a maintainer judges the patch appropriate for stable later on (or you notify them), they will add the same line to the patch, meaning that the stable maintainers will be on copy on the maintainer's pull request.
- If you judge an already merged patch suitable for stable, send a mail (preferably as a reply to the most recent patch submission) to `qemu-stable@nongnu.org` along with `qemu-devel@nongnu.org` and appropriate other people (like the patch author or the relevant maintainer) on copy.

6.10.4 Stable release process

When the stable maintainers prepare a new stable release, they will prepare a git branch with a release candidate and send the patches out to `qemu-devel@nongnu.org` for review. If any of your patches are included, please verify that they look fine, especially if the maintainer had to tweak the patch as part of back-porting things across branches. You may also nominate other patches that you think are suitable for inclusion. After review is complete (may involve more release candidates), a new stable release is made available.

6.11 QTest Device Emulation Testing Framework

QTest is a device emulation testing framework. It can be very useful to test device models; it could also control certain aspects of QEMU (such as virtual clock stepping), with a special purpose “qtest” protocol. Refer to [QTest Protocol](#) for more details of the protocol.

QTest cases can be executed with

```
make check-qtest
```

The QTest library is implemented by `tests/qtest/libqtest.c` and the API is defined in `tests/qtest/libqtest.h`.

Consider adding a new QTest case when you are introducing a new virtual hardware, or extending one if you are adding functionalities to an existing virtual device.

On top of `libqtest`, a higher level library, `libqos`, was created to encapsulate common tasks of device drivers, such as memory management and communicating with system buses or devices. Many virtual device tests use `libqos` instead of directly calling into `libqtest`.

Steps to add a new QTest case are:

1. Create a new source file for the test. (More than one file can be added as necessary.) For example, `tests/qtest/foo-test.c`.
2. Write the test code with the `glib` and `libqtest/libqos` API. See also existing tests and the library headers for reference.

3. Register the new test in `tests/qtest/meson.build`. Add the test executable name to an appropriate `qtests_*` variable. There is one variable per architecture, plus `qtests_generic` for tests that can be run for all architectures. For example:

```
qtests_generic = [
    ...
    'foo-test',
    ...
]
```

4. If the test has more than one source file or needs to be linked with any dependency other than `qemuutil` and `qos`, list them in the `qtests` dictionary. For example a test that needs to use the `QIO` library will have an entry like:

```
{
    ...
    'foo-test': [io],
    ...
}
```

Debugging a QTest failure is slightly harder than the unit test because the tests look up QEMU program names in the environment variables, such as `QTEST_QEMU_BINARY` and `QTEST_QEMU_IMG`, and also because it is not easy to attach `gdb` to the QEMU process spawned from the test. But manual invoking and using `gdb` on the test is still simple to do: find out the actual command from the output of

```
make check-qtest V=1
```

which you can run manually.

6.11.1 QTest Protocol

Line based protocol, request/response based. Server can send async messages so clients should always handle many async messages before the response comes in.

Valid requests

Clock management:

The `qtest` client is completely in charge of the `QEMU_CLOCK_VIRTUAL`. `qtest` commands let you adjust the value of the clock (monotonically). All the commands return the current value of the clock in nanoseconds.

```
> clock_step
< OK VALUE
```

Advance the clock to the next deadline. Useful when waiting for asynchronous events.

```
> clock_step NS
< OK VALUE
```

Advance the clock by `NS` nanoseconds.

```
> clock_set NS
< OK VALUE
```

Advance the clock to `NS` nanoseconds (do nothing if it's already past).

PIO and memory access:

```
> outb ADDR VALUE  
< OK
```

```
> outw ADDR VALUE  
< OK
```

```
> outl ADDR VALUE  
< OK
```

```
> inb ADDR  
< OK VALUE
```

```
> inw ADDR  
< OK VALUE
```

```
> inl ADDR  
< OK VALUE
```

```
> writeb ADDR VALUE  
< OK
```

```
> writew ADDR VALUE  
< OK
```

```
> writel ADDR VALUE  
< OK
```

```
> writeq ADDR VALUE  
< OK
```

```
> readb ADDR  
< OK VALUE
```

```
> readw ADDR  
< OK VALUE
```

```
> readl ADDR  
< OK VALUE
```

```
> readq ADDR  
< OK VALUE
```

```
> read ADDR SIZE  
< OK DATA
```

```
> write ADDR SIZE DATA  
< OK
```

```
> b64read ADDR SIZE
< OK B64_DATA
```

```
> b64write ADDR SIZE B64_DATA
< OK
```

```
> memset ADDR SIZE VALUE
< OK
```

ADDR, SIZE, VALUE are all integers parsed with `strtoul()` with a base of 0. For ‘memset’ a zero size is permitted and does nothing.

DATA is an arbitrarily long hex number prefixed with ‘0x’. If it’s smaller than the expected size, the value will be zero filled at the end of the data sequence.

B64_DATA is an arbitrarily long base64 encoded string. If the sizes do not match, the data will be truncated.

IRQ management:

```
> irq_intercept_in QOM-PATH
< OK
```

```
> irq_intercept_out QOM-PATH
< OK
```

Attach to the gpio-in (resp. gpio-out) pins exported by the device at QOM-PATH. When the pin is triggered, one of the following async messages will be printed to the QTest stream:

```
IRQ raise NUM
IRQ lower NUM
```

where NUM is an IRQ number. For the PC, interrupts can be intercepted simply with “irq_intercept_in ioapic” (note that IRQ0 comes out with NUM=0 even though it is remapped to GSI 2).

Setting interrupt level:

```
> set_irq_in QOM-PATH NAME NUM LEVEL
< OK
```

where NAME is the name of the irq/gpio list, NUM is an IRQ number and LEVEL is a signed integer IRQ level.

Forcibly set the given interrupt pin to the given level.

6.11.2 libqtest API reference

QTestState * **qtest_initf** (const char **fmt*, ...)

Parameters

const char *fmt Format for creating other arguments to pass to QEMU, formatted like `sprintf()`.

... variable arguments

Description

Convenience wrapper around `qtest_init()`.

Return

`QTestState` instance.

`QTestState * qtest_vinitf` (const char **fmt*, va_list *ap*)

Parameters

const char *fmt Format for creating other arguments to pass to QEMU, formatted like `vsprintf()`.

va_list ap Format arguments.

Description

Convenience wrapper around `qtest_init()`.

Return

`QTestState` instance.

`QTestState * qtest_init` (const char **extra_args*)

Parameters

const char *extra_args other arguments to pass to QEMU. CAUTION: these arguments are subject to word splitting and shell evaluation.

Return

`QTestState` instance.

`QTestState * qtest_init_without_qmp_handshake` (const char **extra_args*)

Parameters

const char *extra_args other arguments to pass to QEMU. CAUTION: these arguments are subject to word splitting and shell evaluation.

Return

`QTestState` instance.

`QTestState * qtest_init_with_serial` (const char **extra_args*, int **sock_fd*)

Parameters

const char *extra_args other arguments to pass to QEMU. CAUTION: these arguments are subject to word splitting and shell evaluation.

int *sock_fd pointer to store the socket file descriptor for connection with serial.

Return

`QTestState` instance.

void `qtest_quit` (`QTestState *s`)

Parameters

QTestState *s `QTestState` instance to operate on.

Description

Shut down the QEMU process associated to `s`.

`QDict * qtest_qmp_fds` (`QTestState *s`, int **fds*, size_t *fds_num*, const char **fmt*, ...)

Parameters

QTestState *s QTestState instance to operate on.

int *fds array of file descriptors

size_t fds_num number of elements in **fds**

const char *fmt QMP message to send to qemu, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'. See

... variable arguments

Description

Sends a QMP message to QEMU with **fds** and returns the response.

QDict * **qtest_qmp** (QTestState *s, const char *fmt, ...)

Parameters

QTestState *s QTestState instance to operate on.

const char *fmt QMP message to send to qemu, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'. See

... variable arguments

Description

Sends a QMP message to QEMU and returns the response.

void **qtest_qmp_send** (QTestState *s, const char *fmt, ...)

Parameters

QTestState *s QTestState instance to operate on.

const char *fmt QMP message to send to qemu, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'. See

... variable arguments

Description

Sends a QMP message to QEMU and leaves the response in the stream.

void **qtest_qmp_send_raw** (QTestState *s, const char *fmt, ...)

Parameters

QTestState *s QTestState instance to operate on.

const char *fmt text to send, formatted like `sprintf()`

... variable arguments

Description

Sends text to the QMP monitor verbatim. Need not be valid JSON; this is useful for negative tests.

QDict * **qtest_vqmp_fds** (QTestState *s, int *fds, size_t fds_num, const char *fmt, va_list ap)

Parameters

QTestState *s QTestState instance to operate on.

int *fds array of file descriptors

size_t fds_num number of elements in **fds**

const char *fmt QMP message to send to QEMU, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'.
va_list ap QMP message arguments

Description

Sends a QMP message to QEMU with `fds` and returns the response.

QDict * **qtest_vqmp** (QTestState *s, const char *fmt, va_list ap)

Parameters

QTestState *s QTestState instance to operate on.

const char *fmt QMP message to send to QEMU, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'.
va_list ap QMP message arguments

Description

Sends a QMP message to QEMU and returns the response.

void **qtest_qmp_vsend_fds** (QTestState *s, int *fds, size_t fds_num, const char *fmt, va_list ap)

Parameters

QTestState *s QTestState instance to operate on.

int *fds array of file descriptors

size_t fds_num number of elements in `fds`

const char *fmt QMP message to send to QEMU, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'.
va_list ap QMP message arguments

Description

Sends a QMP message to QEMU and leaves the response in the stream.

void **qtest_qmp_vsend** (QTestState *s, const char *fmt, va_list ap)

Parameters

QTestState *s QTestState instance to operate on.

const char *fmt QMP message to send to QEMU, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'.
va_list ap QMP message arguments

Description

Sends a QMP message to QEMU and leaves the response in the stream.

QDict * **qtest_qmp_receive_dict** (QTestState *s)

Parameters

QTestState *s QTestState instance to operate on.

Description

Reads a QMP message from QEMU and returns the response.

QDict * **qtest_qmp_receive** (QTestState *s)

Parameters

QTestState *s QTestState instance to operate on.

Description

Reads a QMP message from QEMU and returns the response. Buffers all the events received meanwhile, until a call to `qtest_qmp_eventwait`

void **qtest_qmp_eventwait** (QTestState *s, const char *event)

Parameters

QTestState *s QTestState instance to operate on.

const char *event event to wait for.

Description

Continuously polls for QMP responses until it receives the desired event.

QDict * **qtest_qmp_eventwait_ref** (QTestState *s, const char *event)

Parameters

QTestState *s QTestState instance to operate on.

const char *event event to wait for.

Description

Continuously polls for QMP responses until it receives the desired event. Returns a copy of the event for further investigation.

QDict * **qtest_qmp_event_ref** (QTestState *s, const char *event)

Parameters

QTestState *s QTestState instance to operate on.

const char *event event to return.

Description

Removes non-matching events from the buffer that was set by `qtest_qmp_receive`, until an event bearing the given name is found, and returns it. If no event matches, clears the buffer and returns NULL.

char * **qtest_hmp** (QTestState *s, const char *fmt, ...)

Parameters

QTestState *s QTestState instance to operate on.

const char *fmt HMP command to send to QEMU, formats arguments like `sprintf()`.

... variable arguments

Description

Send HMP command to QEMU via QMP's human-monitor-command. QMP events are discarded.

Return

the command's output. The caller should `g_free()` it.

char * **qtest_vhmp** (QTestState *s, const char *fmt, va_list ap)

Parameters

QTestState *s QTestState instance to operate on.

const char *fmt HMP command to send to QEMU, formats arguments like `vsprintf()`.

va_list ap HMP command arguments

Description

Send HMP command to QEMU via QMP's human-monitor-command. QMP events are discarded.

Return

the command's output. The caller should `g_free()` it.

bool **qtest_get_irq**(QTestState *s, int num)

Parameters

QTestState *s QTestState instance to operate on.

int num Interrupt to observe.

Return

The level of the **num** interrupt.

void **qtest_irq_intercept_in**(QTestState *s, const char *string)

Parameters

QTestState *s QTestState instance to operate on.

const char *string QOM path of a device.

Description

Associate qtest irqs with the GPIO-in pins of the device whose path is specified by **string**.

void **qtest_irq_intercept_out**(QTestState *s, const char *string)

Parameters

QTestState *s QTestState instance to operate on.

const char *string QOM path of a device.

Description

Associate qtest irqs with the GPIO-out pins of the device whose path is specified by **string**.

void **qtest_set_irq_in**(QTestState *s, const char *string, const char *name, int irq, int level)

Parameters

QTestState *s QTestState instance to operate on.

const char *string QOM path of a device

const char *name IRQ name

int irq IRQ number

int level IRQ level

Description

Force given device/irq GPIO-in pin to the given level.

void **qtest_outb**(QTestState *s, uint16_t addr, uint8_t value)

Parameters

QTestState *s QTestState instance to operate on.

uint16_t addr I/O port to write to.

uint8_t value Value being written.

Description

Write an 8-bit value to an I/O port.

void **qtest_outw** (QTestState *s, uint16_t *addr*, uint16_t *value*)

Parameters

QTestState *s QTestState instance to operate on.

uint16_t addr I/O port to write to.

uint16_t value Value being written.

Description

Write a 16-bit value to an I/O port.

void **qtest_outl** (QTestState *s, uint16_t *addr*, uint32_t *value*)

Parameters

QTestState *s QTestState instance to operate on.

uint16_t addr I/O port to write to.

uint32_t value Value being written.

Description

Write a 32-bit value to an I/O port.

uint8_t **qtest_inb** (QTestState *s, uint16_t *addr*)

Parameters

QTestState *s QTestState instance to operate on.

uint16_t addr I/O port to read from.

Description

Returns an 8-bit value from an I/O port.

uint16_t **qtest_inw** (QTestState *s, uint16_t *addr*)

Parameters

QTestState *s QTestState instance to operate on.

uint16_t addr I/O port to read from.

Description

Returns a 16-bit value from an I/O port.

uint32_t **qtest_inl** (QTestState *s, uint16_t *addr*)

Parameters

QTestState *s QTestState instance to operate on.

uint16_t addr I/O port to read from.

Description

Returns a 32-bit value from an I/O port.

void **qtest_writeb** (QTestState *s, uint64_t *addr*, uint8_t *value*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to write to.

uint8_t value Value being written.

Description

Writes an 8-bit value to memory.

void **qtest_writew** (QTestState *s, uint64_t *addr*, uint16_t *value*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to write to.

uint16_t value Value being written.

Description

Writes a 16-bit value to memory.

void **qtest_writel** (QTestState *s, uint64_t *addr*, uint32_t *value*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to write to.

uint32_t value Value being written.

Description

Writes a 32-bit value to memory.

void **qtest_writeq** (QTestState *s, uint64_t *addr*, uint64_t *value*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to write to.

uint64_t value Value being written.

Description

Writes a 64-bit value to memory.

uint8_t **qtest_readb** (QTestState *s, uint64_t *addr*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to read from.

Description

Reads an 8-bit value from memory.

Return

Value read.

uint16_t **qtest_readw** (QTestState *s, uint64_t *addr*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to read from.

Description

Reads a 16-bit value from memory.

Return

Value read.

uint32_t **qtest_readl** (QTestState *s, uint64_t *addr*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to read from.

Description

Reads a 32-bit value from memory.

Return

Value read.

uint64_t **qtest_readq** (QTestState *s, uint64_t *addr*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to read from.

Description

Reads a 64-bit value from memory.

Return

Value read.

void **qtest_memread** (QTestState *s, uint64_t *addr*, void **data*, size_t *size*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to read from.

void *data Pointer to where memory contents will be stored.

size_t size Number of bytes to read.

Description

Read guest memory into a buffer.

uint64_t **qtest_rtas_call** (QTestState *s, const char **name*, uint32_t *nargs*, uint64_t *args*, uint32_t *nret*,
uint64_t *ret*)

Parameters

QTestState *s QTestState instance to operate on.

const char *name name of the command to call.

uint32_t nargs Number of args.

uint64_t args Guest address to read args from.

uint32_t nret Number of return value.

uint64_t ret Guest address to write return values to.

Description

Call an RTAS function

void **qtest_bufread** (QTestState *s, uint64_t *addr*, void **data*, size_t *size*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to read from.

void *data Pointer to where memory contents will be stored.

size_t size Number of bytes to read.

Description

Read guest memory into a buffer and receive using a base64 encoding.

void **qtest_memwrite** (QTestState *s, uint64_t *addr*, const void **data*, size_t *size*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to write to.

const void *data Pointer to the bytes that will be written to guest memory.

size_t size Number of bytes to write.

Description

Write a buffer to guest memory.

void **qtest_bufwrite** (QTestState *s, uint64_t *addr*, const void **data*, size_t *size*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to write to.

const void *data Pointer to the bytes that will be written to guest memory.

size_t size Number of bytes to write.

Description

Write a buffer to guest memory and transmit using a base64 encoding.

void **qtest_memset** (QTestState *s, uint64_t *addr*, uint8_t *patt*, size_t *size*)

Parameters

QTestState *s QTestState instance to operate on.

uint64_t addr Guest address to write to.

uint8_t patt Byte pattern to fill the guest memory region with.

size_t size Number of bytes to write.

Description

Write a pattern to guest memory.

`int64_t QTestState *s` `QTestState` instance to operate on.

Parameters

QTestState *s `QTestState` instance to operate on.

Description

Advance the QEMU_CLOCK_VIRTUAL to the next deadline.

Return

The current value of the QEMU_CLOCK_VIRTUAL in nanoseconds.

`int64_t QTestState *s, int64_t step`

Parameters

QTestState *s `QTestState` instance to operate on.

int64_t step Number of nanoseconds to advance the clock by.

Description

Advance the QEMU_CLOCK_VIRTUAL by **step** nanoseconds.

Return

The current value of the QEMU_CLOCK_VIRTUAL in nanoseconds.

`int64_t QTestState *s, int64_t val`

Parameters

QTestState *s `QTestState` instance to operate on.

int64_t val Nanoseconds value to advance the clock to.

Description

Advance the QEMU_CLOCK_VIRTUAL to **val** nanoseconds since the VM was launched.

Return

The current value of the QEMU_CLOCK_VIRTUAL in nanoseconds.

`bool QTestState *s`

Parameters

QTestState *s `QTestState` instance to operate on.

Return

True if the architecture under test has a big endian configuration.

`const char * QTestState *s`

Parameters

void no arguments

Return

The architecture for the QEMU executable under test.

`void QTestState *s, const char *str, void (*fn)(void)`

Parameters

const char *str Test case path.

void (*fn) (void) Test case function

Description

Add a GTester testcase with the given name and function. The path is prefixed with the architecture under test, as returned by `qtest_get_arch()`.

void **qtest_add_data_func** (const char *str, const void *data, void (*fn)(const void *))

Parameters

const char *str Test case path.

const void *data Test case data

void (*fn) (const void *) Test case function

Description

Add a GTester testcase with the given name, data and function. The path is prefixed with the architecture under test, as returned by `qtest_get_arch()`.

void **qtest_add_data_func_full** (const char *str, void *data, void (*fn)(const void *), GDestroyNotify data_free_func)

Parameters

const char *str Test case path.

void *data Test case data

void (*fn) (const void *) Test case function

GDestroyNotify data_free_func GDestroyNotify for data

Description

Add a GTester testcase with the given name, data and function. The path is prefixed with the architecture under test, as returned by `qtest_get_arch()`.

data is passed to **data_free_func()** on test completion.

qtest_add (testpath, Fixture, tdata, fsetup, ftest, fteardown)

Parameters

testpath Test case path

Fixture Fixture type

tdata Test case data

fsetup Test case setup function

ftest Test case function

fteardown Test case teardown function

Description

Add a GTester testcase with the given name, data and functions. The path is prefixed with the architecture under test, as returned by `qtest_get_arch()`.

void **qtest_qmp_assert_success** (QTestState *qts, const char *fmt, ...)

Parameters

QTestState *qts QTestState instance to operate on

const char *fmt QMP message to send to qemu, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'.
 ... variable arguments

Description

Sends a QMP message to QEMU and asserts that a 'return' key is present in the response.

void **qtest_cb_for_every_machine** (void (*cb)(const char *machine), bool *skip_old_versioned*)

Parameters

void (*cb) (const char *machine) Pointer to the callback function

bool skip_old_versioned true if versioned old machine types should be skipped

Call a callback function for every name of all available machines.

void **qtest_qmp_device_add_qdict** (QTestState *qts, const char *drv, const QDict *arguments)

Parameters

QTestState *qts QTestState instance to operate on

const char *drv Name of the device that should be added

const QDict *arguments QDict with properties for the device to initialize

Description

Generic hot-plugging test via the device_add QMP command with properties supplied in form of QDict. Use NULL for empty properties list.

void **qtest_qmp_device_add** (QTestState *qts, const char *driver, const char *id, const char *fmt, ...)

Parameters

QTestState *qts QTestState instance to operate on

const char *driver Name of the device that should be added

const char *id Identification string

const char *fmt QMP message to send to qemu, formatted like `qobject_from_jsonf_nofail()`. See `parse_interpolation()` for what's supported after '%'.
 ... variable arguments

Description

Generic hot-plugging test via the device_add QMP command.

void **qtest_qmp_device_del** (QTestState *qts, const char *id)

Parameters

QTestState *qts QTestState instance to operate on

const char *id Identification string

Description

Generic hot-unplugging test via the device_del QMP command.

bool **qmp_rsp_is_err** (QDict *rsp)

Parameters

QDict *rsp QMP response to check for error

Description

Test **rsp** for error and discard **rsp**. Returns 'true' if there is error in **rsp** and 'false' otherwise.

void **qmp_expect_error_and_unref** (QDict **rsp*, const char **class*)

Parameters

QDict *rsp QMP response to check for error

const char *class an error class

Description

Assert the response has the given error class and discard **rsp**.

bool **qtest_probe_child** (QTestState **s*)

Parameters

QTestState *s QTestState instance to operate on.

Return

true if the child is still alive.

void **qtest_set_expected_status** (QTestState **s*, int *status*)

Parameters

QTestState *s QTestState instance to operate on.

int status an expected exit status.

Description

Set expected exit status of the child.

6.12 Decodetree Specification

A *decodetree* is built from instruction *patterns*. A pattern may represent a single architectural instruction or a group of same, depending on what is convenient for further processing.

Each pattern has both *fixedbits* and *fixedmask*, the combination of which describes the condition under which the pattern is matched:

```
(insn & fixedmask) == fixedbits
```

Each pattern may have *fields*, which are extracted from the *insn* and passed along to the translator. Examples of such are registers, immediates, and sub-opcodes.

In support of patterns, one may declare *fields*, *argument sets*, and *formats*, each of which may be re-used to simplify further definitions.

6.12.1 Fields

Syntax:

```
field_def      := '%' identifier ( unnamed_field ) * ( !function=identifier ) ?  
unnamed_field := number ':' ( 's' ) number
```

For *unnamed_field*, the first number is the least-significant bit position of the field and the second number is the length of the field. If the ‘s’ is present, the field is considered signed. If multiple *unnamed_fields* are present, they are concatenated. In this way one can define disjoint fields.

If *!function* is specified, the concatenated result is passed through the named function, taking and returning an integral value.

One may use *!function* with zero *unnamed_fields*. This case is called a *parameter*, and the named function is only passed the *DisasContext* and returns an integral value extracted from there.

A field with no *unnamed_fields* and no *!function* is in error.

FIXME: the fields of the structure into which this result will be stored is restricted to *int*. Which means that we cannot expand 64-bit items.

Field examples:

Input	Generated code
<code>%disp 0:s16</code>	<code>sextract(i, 0, 16)</code>
<code>%imm9 16:6 10:3</code>	<code>extract(i, 16, 6) << 3 extract(i, 10, 3)</code>
<code>%disp12 0:s1 1:1 2:10</code>	<code>sextract(i, 0, 1) << 11 extract(i, 1, 1) << 10 extract(i, 2, 10)</code>
<code>%shimm8 5:s8 13:1 !function=expand_shimm8</code>	<code>expand_shimm8(sextract(i, 5, 8) << 1 extract(i, 13, 1))</code>

6.12.2 Argument Sets

Syntax:

```
args_def      := '&' identifier ( args_elt )+ ( !extern )?
args_elt      := identifier
```

Each *args_elt* defines an argument within the argument set. Each argument set will be rendered as a C structure “arg_\$name” with each of the fields being one of the member arguments.

If *!extern* is specified, the backing structure is assumed to have been already declared, typically via a second decoder.

Argument sets are useful when one wants to define helper functions for the translator functions that can perform operations on a common set of arguments. This can ensure, for instance, that the AND pattern and the OR pattern put their operands into the same named structure, so that a common *gen_logic_insn* may be able to handle the operations common between the two.

Argument set examples:

```
&reg3      ra rb rc
&loadstore reg base offset
```

6.12.3 Formats

Syntax:

```

fmt_def      := '@' identifier ( fmt_elt )+
fmt_elt      := fixedbit_elt | field_elt | field_ref | args_ref
fixedbit_elt := [01.-]+
field_elt    := identifier ':' 's'? number
field_ref    := '%' identifier | identifier '=' '%' identifier
args_ref     := '&' identifier

```

Defining a format is a handy way to avoid replicating groups of fields across many instruction patterns.

A *fixedbit_elt* describes a contiguous sequence of bits that must be 1, 0, or don't care. The difference between '.' and '-' is that '.' means that the bit will be covered with a field or a final 0 or 1 from the pattern, and '-' means that the bit is really ignored by the cpu and will not be specified.

A *field_elt* describes a simple field only given a width; the position of the field is implied by its position with respect to other *fixedbit_elt* and *field_elt*.

If any *fixedbit_elt* or *field_elt* appear, then all bits must be defined. Padding with a *fixedbit_elt* of all '.' is an easy way to accomplish that.

A *field_ref* incorporates a field by reference. This is the only way to add a complex field to a format. A field may be renamed in the process via assignment to another identifier. This is intended to allow the same argument set be used with disjoint named fields.

A single *args_ref* may specify an argument set to use for the format. The set of fields in the format must be a subset of the arguments in the argument set. If an argument set is not specified, one will be inferred from the set of fields.

It is recommended, but not required, that all *field_ref* and *args_ref* appear at the end of the line, not interleaving with *fixedbit_elt* or *field_elt*.

Format examples:

```

@opr      ..... ra:5 rb:5 ... 0 ..... rc:5
@opi      ..... ra:5 lit:8  1 ..... rc:5

```

6.12.4 Patterns

Syntax:

```

pat_def      := identifier ( pat_elt )+
pat_elt      := fixedbit_elt | field_elt | field_ref | args_ref | fmt_ref | const_elt
fmt_ref      := '@' identifier
const_elt    := identifier '=' number

```

The *fixedbit_elt* and *field_elt* specifiers are unchanged from formats. A pattern that does not specify a named format will have one inferred from a referenced argument set (if present) and the set of fields.

A *const_elt* allows a argument to be set to a constant value. This may come in handy when fields overlap between patterns and one has to include the values in the *fixedbit_elt* instead.

The decoder will call a translator function for each pattern matched.

Pattern examples:

```

addl_r    010000 ..... 0000000 ..... @opr
addl_i    010000 ..... 0000000 ..... @opi

```

which will, in part, invoke:

```
trans_addl_r(ctx, &arg_opr, insn)
```

and:

```
trans_addl_i(ctx, &arg_opi, insn)
```

6.12.5 Pattern Groups

Syntax:

```
group          := overlap_group | no_overlap_group
overlap_group  := '{' ( pat_def | group )+ '}'
no_overlap_group := '[' ( pat_def | group )+ ']'
```

A *group* begins with a lone open-brace or open-bracket, with all subsequent lines indented two spaces, and ending with a lone close-brace or close-bracket. Groups may be nested, increasing the required indentation of the lines within the nested group to two spaces per nesting level.

Patterns within overlap groups are allowed to overlap. Conflicts are resolved by selecting the patterns in order. If all of the fixedbits for a pattern match, its translate function will be called. If the translate function returns false, then subsequent patterns within the group will be matched.

Patterns within no-overlap groups are not allowed to overlap, just the same as ungrouped patterns. Thus no-overlap groups are intended to be nested inside overlap groups.

The following example from PA-RISC shows specialization of the *or* instruction:

```
{
{
  nop   000010  -----  0000 001001 0 00000
  copy  000010 00000  r1:5  0000 001001 0 rt:5
}
or      000010 rt2:5 r1:5  cf:4 001001 0 rt:5
}
```

When the *cf* field is zero, the instruction has no side effects, and may be specialized. When the *rt* field is zero, the output is discarded and so the instruction has no effect. When the *rt2* field is zero, the operation is `reg[r1] | 0` and so encodes the canonical register copy operation.

The output from the generator might look like:

```
switch (insn & 0xfc000fe0) {
case 0x08000240:
  /* 000010.. ..... 0010 010..... */
  if ((insn & 0x0000f000) == 0x00000000) {
    /* 000010.. ..... 00000010 010..... */
    if ((insn & 0x0000001f) == 0x00000000) {
      /* 000010.. ..... 00000010 01000000 */
      extract_decode_Fmt_0(&u.f_decode0, insn);
      if (trans_nop(ctx, &u.f_decode0)) return true;
    }
    if ((insn & 0x03e00000) == 0x00000000) {
      /* 00001000 000..... 00000010 010..... */
      extract_decode_Fmt_1(&u.f_decode1, insn);
      if (trans_copy(ctx, &u.f_decode1)) return true;
    }
  }
}
```

(continues on next page)

(continued from previous page)

```
}
extract_decode_Fmt_2(&u.f_decode2, insn);
if (trans_or(ctx, &u.f_decode2)) return true;
return false;
}
```

6.13 Secure Coding Practices

This document covers topics that both developers and security researchers must be aware of so that they can develop safe code and audit existing code properly.

6.13.1 Reporting Security Bugs

For details on how to report security bugs or ask questions about potential security bugs, see the [Security Process wiki page](#).

6.13.2 General Secure C Coding Practices

Most CVEs (security bugs) reported against QEMU are not specific to virtualization or emulation. They are simply C programming bugs. Therefore it's critical to be aware of common classes of security bugs.

There is a wide selection of resources available covering secure C coding. For example, the [CERT C Coding Standard](#) covers the most important classes of security bugs.

Instead of describing them in detail here, only the names of the most important classes of security bugs are mentioned:

- Buffer overflows
- Use-after-free and double-free
- Integer overflows
- Format string vulnerabilities

Some of these classes of bugs can be detected by analyzers. Static analysis is performed regularly by Coverity and the most obvious of these bugs are even reported by compilers. Dynamic analysis is possible with valgrind, tsan, and asan.

6.13.3 Input Validation

Inputs from the guest or external sources (e.g. network, files) cannot be trusted and may be invalid. Inputs must be checked before using them in a way that could crash the program, expose host memory to the guest, or otherwise be exploitable by an attacker.

The most sensitive attack surface is device emulation. All hardware register accesses and data read from guest memory must be validated. A typical example is a device that contains multiple units that are selectable by the guest via an index register:

```
typedef struct {
    ProcessingUnit unit[2];
    ...
} MyDeviceState;
```

(continues on next page)

(continued from previous page)

```
static void mydev_writel(void *opaque, uint32_t addr, uint32_t val)
{
    MyDeviceState *mydev = opaque;
    ProcessingUnit *unit;

    switch (addr) {
    case MYDEV_SELECT_UNIT:
        unit = &mydev->unit[val];    <-- this input wasn't validated!
        ...
    }
}
```

If `val` is not in range `[0, 1]` then an out-of-bounds memory access will take place when `unit` is dereferenced. The code must check that `val` is 0 or 1 and handle the case where it is invalid.

6.13.4 Unexpected Device Accesses

The guest may access device registers in unusual orders or at unexpected moments. Device emulation code must not assume that the guest follows the typical “theory of operation” presented in driver writer manuals. The guest may make nonsense accesses to device registers such as starting operations before the device has been fully initialized.

A related issue is that device emulation code must be prepared for unexpected device register accesses while asynchronous operations are in progress. A well-behaved guest might wait for a completion interrupt before accessing certain device registers. Device emulation code must handle the case where the guest overwrites registers or submits further requests before an ongoing request completes. Unexpected accesses must not cause memory corruption or leaks in QEMU.

Invalid device register accesses can be reported with `qemu_log_mask(LOG_GUEST_ERROR, ...)`. The `-d guest_errors` command-line option enables these log messages.

6.13.5 Live Migration

Device state can be saved to disk image files and shared with other users. Live migration code must validate inputs when loading device state so an attacker cannot gain control by crafting invalid device states. Device state is therefore considered untrusted even though it is typically generated by QEMU itself.

6.13.6 Guest Memory Access Races

Guests with multiple vCPUs may modify guest RAM while device emulation code is running. Device emulation code must copy in descriptors and other guest RAM structures and only process the local copy. This prevents time-of-check-to-time-of-use (TOCTOU) race conditions that could cause QEMU to crash when a vCPU thread modifies guest RAM while device emulation is processing it.

6.14 Translator Internals

QEMU is a dynamic translator. When it first encounters a piece of code, it converts it to the host instruction set. Usually dynamic translators are very complicated and highly CPU dependent. QEMU uses some tricks which make it relatively easily portable and simple while achieving good performances.

QEMU's dynamic translation backend is called TCG, for "Tiny Code Generator". For more information, please take a look at `tcg/README`.

Some notable features of QEMU's dynamic translator are:

6.14.1 CPU state optimisations

The target CPUs have many internal states which change the way it evaluates instructions. In order to achieve a good speed, the translation phase considers that some state information of the virtual CPU cannot change in it. The state is recorded in the Translation Block (TB). If the state changes (e.g. privilege level), a new TB will be generated and the previous TB won't be used anymore until the state matches the state recorded in the previous TB. The same idea can be applied to other aspects of the CPU state. For example, on x86, if the SS, DS and ES segments have a zero base, then the translator does not even generate an addition for the segment base.

6.14.2 Direct block chaining

After each translated basic block is executed, QEMU uses the simulated Program Counter (PC) and other cpu state information (such as the CS segment base value) to find the next basic block.

In order to accelerate the most common cases where the new simulated PC is known, QEMU can patch a basic block so that it jumps directly to the next one.

The most portable code uses an indirect jump. An indirect jump makes it easier to make the jump target modification atomic. On some host architectures (such as x86 or PowerPC), the `JUMP` opcode is directly patched so that the block chaining has no overhead.

6.14.3 Self-modifying code and translated code invalidation

Self-modifying code is a special challenge in x86 emulation because no instruction cache invalidation is signaled by the application when code is modified.

User-mode emulation marks a host page as write-protected (if it is not already read-only) every time translated code is generated for a basic block. Then, if a write access is done to the page, Linux raises a SEGV signal. QEMU then invalidates all the translated code in the page and enables write accesses to the page. For system emulation, write protection is achieved through the software MMU.

Correct translated code invalidation is done efficiently by maintaining a linked list of every translated block contained in a given page. Other linked lists are also maintained to undo direct block chaining.

On RISC targets, correctly written software uses memory barriers and cache flushes, so some of the protection above would not be necessary. However, QEMU still requires that the generated code always matches the target instructions in memory in order to handle exceptions correctly.

6.14.4 Exception support

`longjmp()` is used when an exception such as division by zero is encountered.

The host SIGSEGV and SIGBUS signal handlers are used to get invalid memory accesses. QEMU keeps a map from host program counter to target program counter, and looks up where the exception happened based on the host program counter at the exception point.

On some targets, some bits of the virtual CPU's state are not flushed to the memory until the end of the translation block. This is done for internal emulation state that is rarely accessed directly by the program and/or changes very often throughout the execution of a translation block—this includes condition codes on x86, delay slots on SPARC, conditional execution on Arm, and so on. This state is stored for each target instruction, and looked up on exceptions.

6.14.5 MMU emulation

For system emulation QEMU uses a software MMU. In that mode, the MMU virtual to physical address translation is done at every memory access.

QEMU uses an address translation cache (TLB) to speed up the translation. In order to avoid flushing the translated code each time the MMU mappings change, all caches in QEMU are physically indexed. This means that each basic block is indexed with its physical address.

In order to avoid invalidating the basic block chain when MMU mappings change, chaining is only performed when the destination of the jump shares a page with the basic block that is performing the jump.

The MMU can also distinguish RAM and ROM memory areas from MMIO memory areas. Access is faster for RAM and ROM because the translation cache also hosts the offset between guest address and host memory. Accessing MMIO memory areas instead calls out to C code for device emulation. Finally, the MMU helps tracking dirty pages and pages pointed to by translation blocks.

6.15 TCG Instruction Counting

TCG has long supported a feature known as `icount` which allows for instruction counting during execution. This should not be confused with cycle accurate emulation - QEMU does not attempt to emulate how long an instruction would take on real hardware. That is a job for other more detailed (and slower) tools that simulate the rest of a micro-architecture.

This feature is only available for system emulation and is incompatible with multi-threaded TCG. It can be used to better align execution time with wall-clock time so a “slow” device doesn’t run too fast on modern hardware. It can also provide for a degree of deterministic execution and is an essential part of the record/replay support in QEMU.

6.15.1 Core Concepts

At its heart `icount` is simply a count of executed instructions which is stored in the `TimersState` of QEMU’s timer sub-system. The number of executed instructions can then be used to calculate `QEMU_CLOCK_VIRTUAL` which represents the amount of elapsed time in the system since execution started. Depending on the `icount` mode this may either be a fixed number of ns per instruction or adjusted as execution continues to keep wall clock time and virtual time in sync.

To be able to calculate the number of executed instructions the translator starts by allocating a budget of instructions to be executed. The budget of instructions is limited by how long it will be until the next timer will expire. We store this budget as part of a `vCPU icount_decr` field which shared with the machinery for handling `cpu_exit()`. The whole field is checked at the start of every translated block and will cause a return to the outer loop to deal with whatever caused the exit.

In the case of `icount`, before the flag is checked we subtract the number of instructions the translation block would execute. If this would cause the instruction budget to go negative we exit the main loop and regenerate a new translation block with exactly the right number of instructions to take the budget to 0 meaning whatever timer was due to expire will expire exactly when we exit the main run loop.

Dealing with MMIO

While we can adjust the instruction budget for known events like timer expiry we cannot do the same for MMIO. Every load/store we execute might potentially trigger an I/O event, at which point we will need an up to date and accurate reading of the `icount` number.

To deal with this case, when an I/O access is made we:

- restore un-executed instructions to the icount budget
- re-compile a single¹ instruction block for the current PC
- exit the cpu loop and execute the re-compiled block

The new block is created with the CF_LAST_IO compile flag which ensures the final instruction translation starts with a call to `gen_io_start()` so we don't enter a perpetual loop constantly recompiling a single instruction block. For translators using the common `translator_loop` this is done automatically.

Other I/O operations

MMIO isn't the only type of operation for which we might need a correct and accurate clock. IO port instructions and accesses to system registers are the common examples here. These instructions have to be handled by the individual translators which have the knowledge of which operations are I/O operations.

When the translator is handling an instruction of this kind:

- **it must call `gen_io_start()` if `icount` is enabled, at some point before the generation of the code which actually does the I/O, using a code fragment similar to:**

```
if (tb_cflags(s->base.tb) & CF_USE_ICOUNT) {  
    gen_io_start();  
}
```

- it must end the TB immediately after this instruction

Note that some older front-ends call a "`gen_io_end()`" function: this is obsolete and should not be used.

6.16 Tracing

6.16.1 Introduction

This document describes the tracing infrastructure in QEMU and how to use it for debugging, profiling, and observing execution.

6.16.2 Quickstart

Enable tracing of `memory_region_ops_read` and `memory_region_ops_write` events:

```
$ qemu --trace "memory_region_ops_*" ...  
...  
719585@1608130130.441188:memory_region_ops_read cpu 0 mr 0x562fdffb3820 addr 0x3cc_  
↪value 0x67 size 1  
719585@1608130130.441190:memory_region_ops_write cpu 0 mr 0x562fdfdb2f00 addr 0x3d4_  
↪value 0x70e size 2
```

This output comes from the "log" trace backend that is enabled by default when `./configure --enable-trace-backends=BACKENDS` was not explicitly specified.

Multiple patterns can be specified by repeating the `--trace` option:

```
$ qemu --trace "kvm_*" --trace "virtio_*" ...
```

¹ sometimes two instructions if dealing with delay slots

When patterns are used frequently it is more convenient to store them in a file to avoid long command-line options:

```
$ echo "memory_region_ops_*" >/tmp/events
$ echo "kvm_*" >>/tmp/events
$ qemu --trace events=/tmp/events ...
```

6.16.3 Trace events

Sub-directory setup

Each directory in the source tree can declare a set of trace events in a local “trace-events” file. All directories which contain “trace-events” files must be listed in the “trace_events_subdirs” variable in the top level meson.build file. During build, the “trace-events” file in each listed subdirectory will be processed by the “tracetool” script to generate code for the trace events.

The individual “trace-events” files are merged into a “trace-events-all” file, which is also installed into “/usr/share/qemu” with the name “trace-events”. This merged file is to be used by the “simpletrace.py” script to later analyse traces in the simpletrace data format.

The following files are automatically generated in <builddir>/trace/ during the build:

- trace-<subdir>.c - the trace event state declarations
- trace-<subdir>.h - the trace event enums and probe functions
- trace-dtrace-<subdir>.h - DTrace event probe specification
- trace-dtrace-<subdir>.dtrace - DTrace event probe helper declaration
- trace-dtrace-<subdir>.o - binary DTrace provider (generated by dtrace)
- trace-ust-<subdir>.h - UST event probe helper declarations

Here <subdir> is the sub-directory path with ‘/’ replaced by ‘_’. For example, “accel/kvm” becomes “accel_kvm” and the final filename for “trace-<subdir>.c” becomes “trace-accel_kvm.c”.

Source files in the source tree do not directly include generated files in “<builddir>/trace/”. Instead they #include the local “trace.h” file, without any sub-directory path prefix. eg io/channel-buffer.c would do:

```
#include "trace.h"
```

The “io/trace.h” file must be created manually with an #include of the corresponding “trace/trace-<subdir>.h” file that will be generated in the builddir:

```
$ echo '#include "trace/trace-io.h"' >io/trace.h
```

While it is possible to include a trace.h file from outside a source file’s own sub-directory, this is discouraged in general. It is strongly preferred that all events be declared directly in the sub-directory that uses them. The only exception is where there are some shared trace events defined in the top level directory trace-events file. The top level directory generates trace files with a filename prefix of “trace/trace-root” instead of just “trace”. This is to avoid ambiguity between a trace.h in the current directory, vs the top level directory.

Using trace events

Trace events are invoked directly from source code like this:

```
#include "trace.h" /* needed for trace event prototype */

void *qemu_vmalloc(size_t size)
{
    void *ptr;
    size_t align = QEMU_VMALLOC_ALIGN;

    if (size < align) {
        align = getpagesize();
    }
    ptr = qemu_memalign(align, size);
    trace_qemu_vmalloc(size, ptr);
    return ptr;
}
```

Declaring trace events

The “tracetool” script produces the trace.h header file which is included by every source file that uses trace events. Since many source files include trace.h, it uses a minimum of types and other header files included to keep the namespace clean and compile times and dependencies down.

Trace events should use types as follows:

- Use `stdint.h` types for fixed-size types. Most offsets and guest memory addresses are best represented with `uint32_t` or `uint64_t`. Use fixed-size types over primitive types whose size may change depending on the host (32-bit versus 64-bit) so trace events don’t truncate values or break the build.
- Use `void *` for pointers to structs or for arrays. The trace.h header cannot include all user-defined struct declarations and it is therefore necessary to use `void *` for pointers to structs.
- For everything else, use primitive scalar types (`char`, `int`, `long`) with the appropriate signedness.
- Avoid floating point types (`float` and `double`) because SystemTap does not support them. In most cases it is possible to round to an integer type instead. This may require scaling the value first by multiplying it by 1000 or the like when digits after the decimal point need to be preserved.

Format strings should reflect the types defined in the trace event. Take special care to use `PRId64` and `PRId64` for `int64_t` and `uint64_t` types, respectively. This ensures portability between 32- and 64-bit platforms. Format strings must not end with a newline character. It is the responsibility of backends to adapt line ending for proper logging.

Each event declaration will start with the event name, then its arguments, finally a format string for pretty-printing. For example:

```
qemu_vmalloc(size_t size, void *ptr) "size %zu ptr %p"
qemu_vfree(void *ptr) "ptr %p"
```

Hints for adding new trace events

1. Trace state changes in the code. Interesting points in the code usually involve a state change like starting, stopping, allocating, freeing. State changes are good trace events because they can be used to understand the execution of the system.
2. Trace guest operations. Guest I/O accesses like reading device registers are good trace events because they can be used to understand guest interactions.

3. Use correlator fields so the context of an individual line of trace output can be understood. For example, trace the pointer returned by malloc and used as an argument to free. This way mallocs and frees can be matched up. Trace events with no context are not very useful.
4. Name trace events after their function. If there are multiple trace events in one function, append a unique distinguisher at the end of the name.

6.16.4 Generic interface and monitor commands

You can programmatically query and control the state of trace events through a backend-agnostic interface provided by the header “trace/control.h”.

Note that some of the backends do not provide an implementation for some parts of this interface, in which case QEMU will just print a warning (please refer to header “trace/control.h” to see which routines are backend-dependent).

The state of events can also be queried and modified through monitor commands:

- `info trace-events` View available trace events and their state. State 1 means enabled, state 0 means disabled.
- `trace-event NAME on|off` Enable/disable a given trace event or a group of events (using wildcards).

The “`–trace events=<file>`” command line argument can be used to enable the events listed in <file> from the very beginning of the program. This file must contain one event name per line.

If a line in the “`–trace events=<file>`” file begins with a ‘-’, the trace event will be disabled instead of enabled. This is useful when a wildcard was used to enable an entire family of events but one noisy event needs to be disabled.

Wildcard matching is supported in both the monitor command “`trace-event`” and the events list file. That means you can enable/disable the events having a common prefix in a batch. For example, `virtio-blk` trace events could be enabled using the following monitor command:

```
trace-event virtio_blk_* on
```

6.16.5 Trace backends

The “`tracetool`” script automates tedious trace event code generation and also keeps the trace event declarations independent of the trace backend. The trace events are not tightly coupled to a specific trace backend, such as LTTng or SystemTap. Support for trace backends can be added by extending the “`tracetool`” script.

The trace backends are chosen at configure time:

```
./configure --enable-trace-backends=simple,dtrace
```

For a list of supported trace backends, try `./configure –help` or see below. If multiple backends are enabled, the trace is sent to them all.

If no backends are explicitly selected, configure will default to the “`log`” backend.

The following subsections describe the supported trace backends.

Nop

The “`nop`” backend generates empty trace event functions so that the compiler can optimize out trace events completely. This imposes no performance penalty.

Note that regardless of the selected trace backend, events with the “`disable`” property will be generated with the “`nop`” backend.

Log

The “log” backend sends trace events directly to standard error. This effectively turns trace events into debug printf's.

This is the simplest backend and can be used together with existing code that uses `DPRINTF()`.

The `-msg timestamp=on|off` command-line option controls whether or not to print the tid/timestamp prefix for each trace event.

Simpletrace

The “simple” backend writes binary trace logs to a file from a thread, making it lower overhead than the “log” backend. A Python API is available for writing offline trace file analysis scripts. It may not be as powerful as platform-specific or third-party trace backends but it is portable and has no special library dependencies.

Monitor commands

- `trace-file on|off|flush|set <path>` Enable/disable/flush the trace file or set the trace file name.

Analyzing trace files

The “simple” backend produces binary trace files that can be formatted with the `simpletrace.py` script. The script takes the “trace-events-all” file and the binary trace:

```
./scripts/simpletrace.py trace-events-all trace-12345
```

You must ensure that the same “trace-events-all” file was used to build QEMU, otherwise trace event declarations may have changed and output will not be consistent.

Ftrace

The “ftrace” backend writes trace data to ftrace marker. This effectively sends trace events to ftrace ring buffer, and you can compare qemu trace data and kernel(especially `kvm.ko` when using KVM) trace data.

if you use KVM, enable `kvm` events in ftrace:

```
# echo 1 > /sys/kernel/debug/tracing/events/kvm/enable
```

After running qemu by root user, you can get the trace:

```
# cat /sys/kernel/debug/tracing/trace
```

Restriction: “ftrace” backend is restricted to Linux only.

Syslog

The “syslog” backend sends trace events using the POSIX syslog API. The log is opened specifying the `LOG_DAEMON` facility and `LOG_PID` option (so events are tagged with the pid of the particular QEMU process that generated them). All events are logged at `LOG_INFO` level.

NOTE: syslog may squash duplicate consecutive trace events and apply rate limiting.

Restriction: “syslog” backend is restricted to POSIX compliant OS.

LTTng Userspace Tracer

The “ust” backend uses the LTTng Userspace Tracer library. There are no monitor commands built into QEMU, instead UST utilities should be used to list, enable/disable, and dump traces.

Package lttng-tools is required for userspace tracing. You must ensure that the current user belongs to the “tracing” group, or manually launch the lttng-sessiond daemon for the current user prior to running any instance of QEMU.

While running an instrumented QEMU, LTTng should be able to list all available events:

```
lttng list -u
```

Create tracing session:

```
lttng create mysession
```

Enable events:

```
lttng enable-event qemu:g_malloc -u
```

Where the events can either be a comma-separated list of events, or “-a” to enable all tracepoint events. Start and stop tracing as needed:

```
lttng start
lttng stop
```

View the trace:

```
lttng view
```

Destroy tracing session:

```
lttng destroy
```

Babeltrace can be used at any later time to view the trace:

```
babeltrace $HOME/lttng-traces/mysession-<date>-<time>
```

SystemTap

The “dtrace” backend uses DTrace sdt probes but has only been tested with SystemTap. When SystemTap support is detected a .stp file with wrapper probes is generated to make use in scripts more convenient. This step can also be performed manually after a build in order to change the binary name in the .stp probes:

```
scripts/tracetool.py --backends=dtrace --format=stap \
    --binary path/to/qemu-binary \
    --target-type system \
    --target-name x86_64 \
    --group=all \
    trace-events-all \
    qemu.stp
```

To facilitate simple usage of systemtap where there merely needs to be printf logging of certain probes, a helper script “qemu-trace-stap” is provided. Consult its manual page for guidance on its usage.

6.16.6 Trace event properties

Each event in the “trace-events-all” file can be prefixed with a space-separated list of zero or more of the following event properties.

“disable”

If a specific trace event is going to be invoked a huge number of times, this might have a noticeable performance impact even when the event is programmatically disabled.

In this case you should declare such event with the “disable” property. This will effectively disable the event at compile time (by using the “nop” backend), thus having no performance impact at all on regular builds (i.e., unless you edit the “trace-events-all” file).

In addition, there might be cases where relatively complex computations must be performed to generate values that are only used as arguments for a trace function. In these cases you can use ‘trace_event_get_state_backends()’ to guard such computations, so they are skipped if the event has been either compile-time disabled or run-time disabled. If the event is compile-time disabled, this check will have no performance impact.

```
#include "trace.h" /* needed for trace event prototype */

void *qemu_vmalloc(size_t size)
{
    void *ptr;
    size_t align = QEMU_VMALLOC_ALIGN;

    if (size < align) {
        align = getpagesize();
    }
    ptr = qemu_memalign(align, size);
    if (trace_event_get_state_backends	TRACE_QEMU_VMALLOC) {
        void *complex;
        /* some complex computations to produce the 'complex' value */
        trace_qemu_vmalloc(size, ptr, complex);
    }
    return ptr;
}
```

“tcg”

Guest code generated by TCG can be traced by defining an event with the “tcg” event property. Internally, this property generates two events: “<eventname>_trans” to trace the event at translation time, and “<eventname>_exec” to trace the event at execution time.

Instead of using these two events, you should instead use the function “trace_<eventname>_tcg” during translation (TCG code generation). This function will automatically call “trace_<eventname>_trans”, and will generate the necessary TCG code to call “trace_<eventname>_exec” during guest code execution.

Events with the “tcg” property can be declared in the “trace-events” file with a mix of native and TCG types, and “trace_<eventname>_tcg” will gracefully forward them to the “<eventname>_trans” and “<eventname>_exec” events. Since TCG values are not known at translation time, these are ignored by the “<eventname>_trans” event. Because of this, the entry in the “trace-events” file needs two printing formats (separated by a comma):

```
tcg foo(uint8_t a1, TCGv_i32 a2) "a1=%d", "a1=%d a2=%d"
```

For example:

```
#include "trace-tcg.h"

void some_disassembly_func (...)
{
    uint8_t a1 = ...;
    TCGv_i32 a2 = ...;
    trace_foo_tcg(a1, a2);
}
```

This will immediately call:

```
void trace_foo_trans(uint8_t a1);
```

and will generate the TCG code to call:

```
void trace_foo(uint8_t a1, uint32_t a2);
```

“vcpu”

Identifies events that trace vCPU-specific information. It implicitly adds a “CPUState*” argument, and extends the tracing print format to show the vCPU information. If used together with the “tcg” property, it adds a second “TCGv_env” argument that must point to the per-target global TCG register that points to the vCPU when guest code is executed (usually the “cpu_env” variable).

The “tcg” and “vcpu” properties are currently only honored in the root ./trace-events file.

The following example events:

```
foo(uint32_t a) "a=%x"
vcpu bar(uint32_t a) "a=%x"
tcg vcpu baz(uint32_t a) "a=%x", "a=%x"
```

Can be used as:

```
#include "trace-tcg.h"

CPUArchState *env;
TCGv_ptr cpu_env;

void some_disassembly_func(...)
{
    /* trace emitted at this point */
    trace_foo(0xd1);
    /* trace emitted at this point */
    trace_bar(env_cpu(env), 0xd2);
    /* trace emitted at this point (env) and when guest code is executed (cpu_env) */
    trace_baz_tcg(env_cpu(env), cpu_env, 0xd3);
}
```

If the translating vCPU has address 0xc1 and code is later executed by vCPU 0xc2, this would be an example output:

```
// at guest code translation
foo a=0xd1
bar cpu=0xc1 a=0xd2
baz_trans cpu=0xc1 a=0xd3
// at guest code execution
baz_exec cpu=0xc2 a=0xd3
```

6.17 Introduction

This document outlines the design for multi-threaded TCG (a.k.a MTTCG) system-mode emulation. user-mode emulation has always mirrored the thread structure of the translated executable although some of the changes done for MTTCG system emulation have improved the stability of linux-user emulation.

The original system-mode TCG implementation was single threaded and dealt with multiple CPUs with simple round-robin scheduling. This simplified a lot of things but became increasingly limited as systems being emulated gained additional cores and per-core performance gains for host systems started to level off.

6.18 vCPU Scheduling

We introduce a new running mode where each vCPU will run on its own user-space thread. This is enabled by default for all FE/BE combinations where the host memory model is able to accommodate the guest (TCG_GUEST_DEFAULT_MO & ~TCG_TARGET_DEFAULT_MO is zero) and the guest has had the required work done to support this safely (TARGET_SUPPORTS_MTTTCG).

System emulation will fall back to the original round robin approach if:

- forced by `-accel tcg,thread=single`
- enabling `-icount` mode
- 64 bit guests on 32 bit hosts (TCG_OVERSIZED_GUEST)

In the general case of running translated code there should be no inter-vCPU dependencies and all vCPUs should be able to run at full speed. Synchronisation will only be required while accessing internal shared data structures or when the emulated architecture requires a coherent representation of the emulated machine state.

6.19 Shared Data Structures

6.19.1 Main Run Loop

Even when there is no code being generated there are a number of structures associated with the hot-path through the main run-loop. These are associated with looking up the next translation block to execute. These include:

`tb_jump_cache` (per-vCPU, cache of recent jumps) `tb_ctx.htable` (global hash table, phys address->tb lookup)

As TB linking only occurs when blocks are in the same page this code is critical to performance as looking up the next TB to execute is the most common reason to exit the generated code.

DESIGN REQUIREMENT: Make access to lookup structures safe with multiple reader/writer threads. Minimise any lock contention to do it.

The hot-path avoids using locks where possible. The `tb_jump_cache` is updated with atomic accesses to ensure consistent results. The fall back QHT based hash table is also designed for lockless lookups. Locks are only taken when code generation is required or TranslationBlocks have their block-to-block jumps patched.

6.19.2 Global TCG State

User-mode emulation

We need to protect the entire code generation cycle including any post generation patching of the translated code. This also implies a shared translation buffer which contains code running on all cores. Any execution path that comes to the main run loop will need to hold a mutex for code generation. This also includes times when we need flush code or entries from any shared lookups/caches. Structures held on a per-vCPU basis won't need locking unless other vCPUs will need to modify them.

DESIGN REQUIREMENT: Add locking around all code generation and TB patching.

(Current solution)

Code generation is serialised with `mmap_lock()`.

!User-mode emulation

Each vCPU has its own TCG context and associated TCG region, thereby requiring no locking during translation.

6.19.3 Translation Blocks

Currently the whole system shares a single code generation buffer which when full will force a flush of all translations and start from scratch again. Some operations also force a full flush of translations including:

- debugging operations (breakpoint insertion/removal)
- some CPU helper functions
- linux-user spawning its first thread

This is done with the `async_safe_run_on_cpu()` mechanism to ensure all vCPUs are quiescent when changes are being made to shared global structures.

More granular translation invalidation events are typically due to a change of the state of a physical page:

- code modification (self modify code, patching code)
- page changes (new page mapping in linux-user mode)

While setting the invalid flag in a TranslationBlock will stop it being used when looked up in the hot-path there are a number of other book-keeping structures that need to be safely cleared.

Any TranslationBlocks which have been patched to jump directly to the now invalid blocks need the jump patches reversing so they will return to the C code.

There are a number of look-up caches that need to be properly updated including the:

- jump lookup cache
- the physical-to-tb lookup hash table
- the global page table

The global page table (`tlb_map`) which provides a multi-level look-up for `PageDesc` structures which contain pointers to the start of a linked list of all Translation Blocks in that page (see `page_next`).

Both the jump patching and the page cache involve linked lists that the invalidated TranslationBlock needs to be removed from.

DESIGN REQUIREMENT: Safely handle invalidation of TBs

- safely patch/revert direct jumps
- remove central `PageDesc` lookup entries

- ensure lookup caches/hashes are safely updated

(Current solution)

The direct jump themselves are updated atomically by the TCG `tb_setjmp_target()` code. Modification to the linked lists that allow searching for linked pages are done under the protection of `tb->jmp_lock`, where `tb` is the destination block of a jump. Each origin block keeps a pointer to its destinations so that the appropriate lock can be acquired before iterating over a jump list.

The global page table is a lockless radix tree; `cmpxchg` is used to atomically insert new elements.

The lookup caches are updated atomically and the lookup hash uses QHT which is designed for concurrent safe lookup.

Parallel code generation is supported. QHT is used at insertion time as the synchronization point across threads, thereby ensuring that we only keep track of a single TranslationBlock for each guest code block.

6.19.4 Memory maps and TLBs

The memory handling code is fairly critical to the speed of memory access in the emulated system. The SoftMMU code is designed so the hot-path can be handled entirely within translated code. This is handled with a per-vCPU TLB structure which once populated will allow a series of accesses to the page to occur without exiting the translated code. It is possible to set flags in the TLB address which will ensure the slow-path is taken for each access. This can be done to support:

- Memory regions (dividing up access to PIO, MMIO and RAM)
- Dirty page tracking (for code gen, SMC detection, migration and display)
- Virtual TLB (for translating guest address->real address)

When the TLB tables are updated by a vCPU thread other than their own we need to ensure it is done in a safe way so no inconsistent state is seen by the vCPU thread.

Some operations require updating a number of vCPUs TLBs at the same time in a synchronised manner.

DESIGN REQUIREMENTS:

- TLB Flush All/Page - can be across-vCPUs - cross vCPU TLB flush may need other vCPU brought to halt - change may need to be visible to the calling vCPU immediately
- TLB Flag Update - usually cross-vCPU - want change to be visible as soon as possible
- TLB Update (update a CPUTLBEntry, via `tlb_set_page_with_attrs`) - This is a per-vCPU table - by definition can't race - updated by its own thread when the slow-path is forced

(Current solution)

We have updated `cputlb.c` to defer operations when a cross-vCPU operation with `async_run_on_cpu()` which ensures each vCPU sees a coherent state when it next runs its work (in a few instructions time).

A new set up operations (`tlb_flush_*_all_cpus`) take an additional flag which when set will force synchronisation by setting the source vCPUs work as “safe work” and exiting the `cpu run loop`. This ensure by the time execution restarts all flush operations have completed.

TLB flag updates are all done atomically and are also protected by the corresponding page lock.

(Known limitation)

Not really a limitation but the wait mechanism is overly strict for some architectures which only need flushes completed by a barrier instruction. This could be a future optimisation.

6.19.5 Emulated hardware state

Currently thanks to KVM work any access to IO memory is automatically protected by the global iothread mutex, also known as the BQL (Big Qemu Lock). Any IO region that doesn't use global mutex is expected to do its own locking.

However IO memory isn't the only way emulated hardware state can be modified. Some architectures have model specific registers that trigger hardware emulation features. Generally any translation helper that needs to update more than a single vCPUs of state should take the BQL.

As the BQL, or global iothread mutex is shared across the system we push the use of the lock as far down into the TCG code as possible to minimise contention.

(Current solution)

MMIO access automatically serialises hardware emulation by way of the BQL. Currently Arm targets serialise all ARM_CP_IO register accesses and also defer the reset/startup of vCPUs to the vCPU context by way of `async_run_on_cpu()`.

Updates to interrupt state are also protected by the BQL as they can often be cross vCPU.

6.20 Memory Consistency

Between emulated guests and host systems there are a range of memory consistency models. Even emulating weakly ordered systems on strongly ordered hosts needs to ensure things like store-after-load re-ordering can be prevented when the guest wants to.

6.20.1 Memory Barriers

Barriers (sometimes known as fences) provide a mechanism for software to enforce a particular ordering of memory operations from the point of view of external observers (e.g. another processor core). They can apply to any memory operations as well as just loads or stores.

The Linux kernel has an excellent [write-up](#) on the various forms of memory barrier and the guarantees they can provide.

Barriers are often wrapped around synchronisation primitives to provide explicit memory ordering semantics. However they can be used by themselves to provide safe lockless access by ensuring for example a change to a signal flag will only be visible once the changes to payload are.

DESIGN REQUIREMENT: Add a new `tcg_memory_barrier` op

This would enforce a strong load/store ordering so all loads/stores complete at the memory barrier. On single-core non-SMP strongly ordered backends this could become a NOP.

Aside from explicit standalone memory barrier instructions there are also implicit memory ordering semantics which comes with each guest memory access instruction. For example all x86 load/stores come with fairly strong guarantees of sequential consistency whereas Arm has special variants of load/store instructions that imply acquire/release semantics.

In the case of a strongly ordered guest architecture being emulated on a weakly ordered host the scope for a heavy performance impact is quite high.

DESIGN REQUIREMENTS: Be efficient with use of memory barriers

- host systems with stronger implied guarantees can skip some barriers
- merge consecutive barriers to the strongest one

(Current solution)

The system currently has a `tcg_gen_mb()` which will add memory barrier operations if code generation is being done in a parallel context. The `tcg_optimize()` function attempts to merge barriers up to their strongest form before any load/store operations. The solution was originally developed and tested for linux-user based systems. All backends have been converted to emit fences when required. So far the following front-ends have been updated to emit fences when required:

- target-i386
- target-arm
- target-aarch64
- target-alpha
- target-mips

6.20.2 Memory Control and Maintenance

This includes a class of instructions for controlling system cache behaviour. While QEMU doesn't model cache behaviour these instructions are often seen when code modification has taken place to ensure the changes take effect.

6.20.3 Synchronisation Primitives

There are two broad types of synchronisation primitives found in modern ISAs: atomic instructions and exclusive regions.

The first type offer a simple atomic instruction which will guarantee some sort of test and conditional store will be truly atomic w.r.t. other cores sharing access to the memory. The classic example is the x86 `cmpxchg` instruction.

The second type offer a pair of load/store instructions which offer a guarantee that a region of memory has not been touched between the load and store instructions. An example of this is Arm's `ldrex/strex` pair where the `strex` instruction will return a flag indicating a successful store only if no other CPU has accessed the memory region since the `ldrex`.

Traditionally TCG has generated a series of operations that work because they are within the context of a single translation block so will have completed before another CPU is scheduled. However with the ability to have multiple threads running to emulate multiple CPUs we will need to explicitly expose these semantics.

DESIGN REQUIREMENTS:

- Support classic atomic instructions
- Support load/store exclusive (or load link/store conditional) pairs
- Generic enough infrastructure to support all guest architectures

CURRENT OPEN QUESTIONS:

- How problematic is the ABA problem in general?

(Current solution)

The TCG provides a number of atomic helpers (`tcg_gen_atomic_*`) which can be used directly or combined to emulate other instructions like Arm's `ldrex/strex` instructions. While they are susceptible to the ABA problem so far common guests have not implemented patterns where this may be a problem - typically presenting a locking ABI which assumes `cmpxchg` like semantics.

The code also includes a fall-back for cases where multi-threaded TCG ops can't work (e.g. guest atomic width > host atomic width). In this case an EXCP_ATOMIC exit occurs and the instruction is emulated with an exclusive lock which ensures all emulation is serialised.

While the atomic helpers look good enough for now there may be a need to look at solutions that can more closely model the guest architectures semantics.

6.21 QEMU TCG Plugins

QEMU TCG plugins provide a way for users to run experiments taking advantage of the total system control emulation can have over a guest. It provides a mechanism for plugins to subscribe to events during translation and execution and optionally callback into the plugin during these events. TCG plugins are unable to change the system state only monitor it passively. However they can do this down to an individual instruction granularity including potentially subscribing to all load and store operations.

6.21.1 API Stability

This is a new feature for QEMU and it does allow people to develop out-of-tree plugins that can be dynamically linked into a running QEMU process. However the project reserves the right to change or break the API should it need to do so. The best way to avoid this is to submit your plugin upstream so they can be updated if/when the API changes.

API versioning

All plugins need to declare a symbol which exports the plugin API version they were built against. This can be done simply by:

```
QEMU_PLUGIN_EXPORT int qemu_plugin_version = QEMU_PLUGIN_VERSION;
```

The core code will refuse to load a plugin that doesn't export a *qemu_plugin_version* symbol or if plugin version is outside of QEMU's supported range of API versions.

Additionally the *qemu_info_t* structure which is passed to the *qemu_plugin_install* method of a plugin will detail the minimum and current API versions supported by QEMU. The API version will be incremented if new APIs are added. The minimum API version will be incremented if existing APIs are changed or removed.

Exposure of QEMU internals

The plugin architecture actively avoids leaking implementation details about how QEMU's translation works to the plugins. While there are conceptions such as translation time and translation blocks the details are opaque to plugins. The plugin is able to query select details of instructions and system configuration only through the exported *qemu_plugin* functions.

Query Handle Lifetime

Each callback provides an opaque anonymous information handle which can usually be further queried to find out information about a translation, instruction or operation. The handles themselves are only valid during the lifetime of the callback so it is important that any information that is needed is extracted during the callback and saved by the plugin.

6.21.2 Usage

The QEMU binary needs to be compiled for plugin support:

```
configure --enable-plugins
```

Once built a program can be run with multiple plugins loaded each with their own arguments:

```
$QEMU $OTHER_QEMU_ARGS \  
-plugin tests/plugin/libhowvec.so,arg=inline,arg=hint \  
-plugin tests/plugin/libhotblocks.so
```

Arguments are plugin specific and can be used to modify their behaviour. In this case the howvec plugin is being asked to use inline ops to count and break down the hint instructions by type.

6.21.3 Plugin Life cycle

First the plugin is loaded and the public `qemu_plugin_install` function is called. The plugin will then register callbacks for various plugin events. Generally plugins will register a handler for the *atexit* if they want to dump a summary of collected information once the program/system has finished running.

When a registered event occurs the plugin callback is invoked. The callbacks may provide additional information. In the case of a translation event the plugin has an option to enumerate the instructions in a block of instructions and optionally register callbacks to some or all instructions when they are executed.

There is also a facility to add an inline event where code to increment a counter can be directly inlined with the translation. Currently only a simple increment is supported. This is not atomic so can miss counts. If you want absolute precision you should use a callback which can then ensure atomicity itself.

Finally when QEMU exits all the registered *atexit* callbacks are invoked.

6.21.4 Internals

Locking

We have to ensure we cannot deadlock, particularly under MTTCG. For this we acquire a lock when called from plugin code. We also keep the list of callbacks under RCU so that we do not have to hold the lock when calling the callbacks. This is also for performance, since some callbacks (e.g. memory access callbacks) might be called very frequently.

- A consequence of this is that we keep our own list of CPUs, so that we do not have to worry about locking order wrt `cpu_list_lock`.
- Use a recursive lock, since we can get registration calls from callbacks.

As a result registering/unregistering callbacks is “slow”, since it takes a lock. But this is very infrequent; we want performance when calling (or not calling) callbacks, not when registering them. Using RCU is great for this.

We support the uninstallation of a plugin at any time (e.g. from plugin callbacks). This allows plugins to remove themselves if they no longer want to instrument the code. This operation is asynchronous which means callbacks may still occur after the uninstall operation is requested. The plugin isn’t completely uninstalled until the safe work has executed while all vCPUs are quiescent.

6.21.5 Example Plugins

There are a number of plugins included with QEMU and you are encouraged to contribute your own plugins upstream. There is a *contrib/plugins* directory where they can go.

- tests/plugins

These are some basic plugins that are used to test and exercise the API during the *make check-tcg* target.

- contrib/plugins/hotblocks.c

The hotblocks plugin allows you to examine the where hot paths of execution are in your program. Once the program has finished you will get a sorted list of blocks reporting the starting PC, translation count, number of instructions and execution count. This will work best with linux-user execution as system emulation tends to generate re-translations as blocks from different programs get swapped in and out of system memory.

If your program is single-threaded you can use the *inline* option for slightly faster (but not thread safe) counters.

Example:

```
./aarch64-linux-user/qemu-aarch64 \
  -plugin contrib/plugins/libhotblocks.so -d plugin \
  ./tests/tcg/aarch64-linux-user/sha1
SHA1=15dd99a1991e0b3826fede3deffclfeba42278e6
collected 903 entries in the hash table
pc, tcount, icount, ecount
0x0000000041ed10, 1, 5, 66087
0x000000004002b0, 1, 4, 66087
...
```

- contrib/plugins/hotpages.c

Similar to hotblocks but this time tracks memory accesses:

```
./aarch64-linux-user/qemu-aarch64 \
  -plugin contrib/plugins/libhotpages.so -d plugin \
  ./tests/tcg/aarch64-linux-user/sha1
SHA1=15dd99a1991e0b3826fede3deffclfeba42278e6
Addr, RCPUs, Reads, WCPUs, Writes
0x000055007fe000, 0x0001, 31747952, 0x0001, 8835161
0x000055007ff000, 0x0001, 29001054, 0x0001, 8780625
0x00005500800000, 0x0001, 687465, 0x0001, 335857
0x0000000048b000, 0x0001, 130594, 0x0001, 355
0x0000000048a000, 0x0001, 1826, 0x0001, 11
```

- contrib/plugins/howvec.c

This is an instruction classifier so can be used to count different types of instructions. It has a number of options to refine which get counted. You can give an argument for a class of instructions to break it down fully, so for example to see all the system registers accesses:

```
./aarch64-softmmu/qemu-system-aarch64 $(QEMU_ARGS) \
  -append "root=/dev/sda2 systemd.unit=benchmark.service" \
  -smp 4 -plugin ./contrib/plugins/libhowvec.so,arg=sreg -d plugin
```

which will lead to a sorted list after the class breakdown:

```
Instruction Classes:
Class:  UDEF                not counted
Class:  SVE                  (68 hits)
Class:  PCrel addr           (47789483 hits)
Class:  Add/Sub (imm)         (192817388 hits)
Class:  Logical (imm)         (93852565 hits)
Class:  Move Wide (imm)       (76398116 hits)
Class:  Bitfield              (44706084 hits)
```

(continues on next page)

(continued from previous page)

```

Class: Extract (5499257 hits)
Class: Cond Branch (imm) (147202932 hits)
Class: Exception Gen (193581 hits)
Class: NOP not counted
Class: Hints (6652291 hits)
Class: Barriers (8001661 hits)
Class: PSTATE (1801695 hits)
Class: System Insn (6385349 hits)
Class: System Reg counted individually
Class: Branch (reg) (69497127 hits)
Class: Branch (imm) (84393665 hits)
Class: Cmp & Branch (110929659 hits)
Class: Tst & Branch (44681442 hits)
Class: AdvSimd ldstmult (736 hits)
Class: ldst excl (9098783 hits)
Class: Load Reg (lit) (87189424 hits)
Class: ldst noalloc pair (3264433 hits)
Class: ldst pair (412526434 hits)
Class: ldst reg (imm) (314734576 hits)
Class: Loads & Stores (2117774 hits)
Class: Data Proc Reg (223519077 hits)
Class: Scalar FP (31657954 hits)
Individual Instructions:
Instr: mrs x0, sp_el0 (2682661 hits) (op=0xd5384100/ System Reg)
Instr: mrs x1, tpidr_el2 (1789339 hits) (op=0xd53cd041/ System Reg)
Instr: mrs x2, tpidr_el2 (1513494 hits) (op=0xd53cd042/ System Reg)
Instr: mrs x0, tpidr_el2 (1490823 hits) (op=0xd53cd040/ System Reg)
Instr: mrs x1, sp_el0 (933793 hits) (op=0xd5384101/ System Reg)
Instr: mrs x2, sp_el0 (699516 hits) (op=0xd5384102/ System Reg)
Instr: mrs x4, tpidr_el2 (528437 hits) (op=0xd53cd044/ System Reg)
Instr: mrs x30, ttbr1_el1 (480776 hits) (op=0xd538203e/ System Reg)
Instr: msr ttbr1_el1, x30 (480713 hits) (op=0xd518203e/ System Reg)
Instr: msr vbar_el1, x30 (480671 hits) (op=0xd518c01e/ System Reg)
...

```

To find the argument shorthand for the class you need to examine the source code of the plugin at the moment, specifically the **opt* argument in the *InsnClassExecCount* tables.

- contrib/plugins/lockstep.c

This is a debugging tool for developers who want to find out when and where execution diverges after a subtle change to TCG code generation. It is not an exact science and results are likely to be mixed once asynchronous events are introduced. While the use of *-icount* can introduce determinism to the execution flow it doesn't always follow the translation sequence will be exactly the same. Typically this is caused by a timer firing to service the GUI causing a block to end early. However in some cases it has proved to be useful in pointing people at roughly where execution diverges. The only argument you need for the plugin is a path for the socket the two instances will communicate over:

```

./sparc-softmmu/qemu-system-sparc -monitor none -parallel none \
-net none -M SS-20 -m 256 -kernel day11/zImage.elf \
-plugin ./contrib/plugins/liblockstep.so,arg=lockstep-sparc.sock \
-d plugin,nochain

```

which will eventually report:

```

qemu-system-sparc: warning: nic lance.0 has no peer
@ 0x000000ffd06678 vs 0x000000ffd001e0 (2/1 since last)
@ 0x000000ffd07d9c vs 0x000000ffd06678 (3/1 since last)

```

(continues on next page)

(continued from previous page)

```

Δ insn_count @ 0x000000ffd07d9c (809900609) vs 0x000000ffd06678 (809900612)
  previously @ 0x000000ffd06678/10 (809900609 insns)
  previously @ 0x000000ffd001e0/4 (809900599 insns)
  previously @ 0x000000ffd080ac/2 (809900595 insns)
  previously @ 0x000000ffd08098/5 (809900593 insns)
  previously @ 0x000000ffd080c0/1 (809900588 insns)

```

- contrib/plugins/hwprofile

The hwprofile tool can only be used with system emulation and allows the user to see what hardware is accessed how often. It has a number of options:

- arg=read or arg=write

By default the plugin tracks both reads and writes. You can use one of these options to limit the tracking to just one class of accesses.

- arg=source

Will include a detailed break down of what the guest PC that made the access was. Not compatible with arg=pattern. Example output:

```

cirrus-low-memory @ 0xfffffd00000a0000
pc:fffffc0000005cdc, 1, 256
pc:fffffc0000005ce8, 1, 256
pc:fffffc0000005cec, 1, 256

```

- arg=pattern

Instead break down the accesses based on the offset into the HW region. This can be useful for seeing the most used registers of a device. Example output:

```

pci0-conf @ 0xfffffd01fe000000
off:00000004, 1, 1
off:00000010, 1, 3
off:00000014, 1, 3
off:00000018, 1, 2
off:0000001c, 1, 2
off:00000020, 1, 2
...

```

6.22 Bitwise operations

The header `qemu/bitops.h` provides utility functions for performing bitwise operations.

void **set_bit** (long *nr*, unsigned long **addr*)

Set a bit in memory

Parameters

long **nr** the bit to set

unsigned long ***addr** the address to start counting from

void **set_bit_atomic** (long *nr*, unsigned long **addr*)

Set a bit in memory atomically

Parameters

long nr the bit to set

unsigned long *addr the address to start counting from

void **clear_bit** (long *nr*, unsigned long **addr*)
Clears a bit in memory

Parameters

long nr Bit to clear

unsigned long *addr Address to start counting from

void **change_bit** (long *nr*, unsigned long **addr*)
Toggle a bit in memory

Parameters

long nr Bit to change

unsigned long *addr Address to start counting from

int **test_and_set_bit** (long *nr*, unsigned long **addr*)
Set a bit and return its old value

Parameters

long nr Bit to set

unsigned long *addr Address to count from

int **test_and_clear_bit** (long *nr*, unsigned long **addr*)
Clear a bit and return its old value

Parameters

long nr Bit to clear

unsigned long *addr Address to count from

int **test_and_change_bit** (long *nr*, unsigned long **addr*)
Change a bit and return its old value

Parameters

long nr Bit to change

unsigned long *addr Address to count from

int **test_bit** (long *nr*, const unsigned long **addr*)
Determine whether a bit is set

Parameters

long nr bit number to test

const unsigned long *addr Address to start counting from

unsigned long **find_last_bit** (const unsigned long **addr*, unsigned long *size*)
find the last set bit in a memory region

Parameters

const unsigned long *addr The address to start the search at

unsigned long size The maximum size to search

Description

Returns the bit number of the first set bit, or size.

unsigned long **find_next_bit** (const unsigned long **addr*, unsigned long *size*, unsigned long *offset*)
 find the next set bit in a memory region

Parameters

const unsigned long *addr The address to base the search on

unsigned long size The bitmap size in bits

unsigned long offset The bitnumber to start searching at

unsigned long **find_next_zero_bit** (const unsigned long **addr*, unsigned long *size*, unsigned long *offset*)
 find the next cleared bit in a memory region

Parameters

const unsigned long *addr The address to base the search on

unsigned long size The bitmap size in bits

unsigned long offset The bitnumber to start searching at

unsigned long **find_first_bit** (const unsigned long **addr*, unsigned long *size*)
 find the first set bit in a memory region

Parameters

const unsigned long *addr The address to start the search at

unsigned long size The maximum size to search

Description

Returns the bit number of the first set bit.

unsigned long **find_first_zero_bit** (const unsigned long **addr*, unsigned long *size*)
 find the first cleared bit in a memory region

Parameters

const unsigned long *addr The address to start the search at

unsigned long size The maximum size to search

Description

Returns the bit number of the first cleared bit.

uint8_t **rol8** (uint8_t *word*, unsigned int *shift*)
 rotate an 8-bit value left

Parameters

uint8_t word value to rotate

unsigned int shift bits to roll

uint8_t **ror8** (uint8_t *word*, unsigned int *shift*)
 rotate an 8-bit value right

Parameters

uint8_t word value to rotate

unsigned int shift bits to roll

`uint16_t rol16 (uint16_t word, unsigned int shift)`
rotate a 16-bit value left

Parameters

`uint16_t word` value to rotate

`unsigned int shift` bits to roll

`uint16_t ror16 (uint16_t word, unsigned int shift)`
rotate a 16-bit value right

Parameters

`uint16_t word` value to rotate

`unsigned int shift` bits to roll

`uint32_t rol32 (uint32_t word, unsigned int shift)`
rotate a 32-bit value left

Parameters

`uint32_t word` value to rotate

`unsigned int shift` bits to roll

`uint32_t ror32 (uint32_t word, unsigned int shift)`
rotate a 32-bit value right

Parameters

`uint32_t word` value to rotate

`unsigned int shift` bits to roll

`uint64_t rol64 (uint64_t word, unsigned int shift)`
rotate a 64-bit value left

Parameters

`uint64_t word` value to rotate

`unsigned int shift` bits to roll

`uint64_t ror64 (uint64_t word, unsigned int shift)`
rotate a 64-bit value right

Parameters

`uint64_t word` value to rotate

`unsigned int shift` bits to roll

`uint32_t extract32 (uint32_t value, int start, int length)`

Parameters

`uint32_t value` the value to extract the bit field from

`int start` the lowest bit in the bit field (numbered from 0)

`int length` the length of the bit field

Description

Extract from the 32 bit input **value** the bit field specified by the **start** and **length** parameters, and return it. The bit field must lie entirely within the 32 bit word. It is valid to request that all 32 bits are returned (ie **length** 32 and **start** 0).

Return

the value of the bit field extracted from the input value.

`uint8_t extract8 (uint8_t value, int start, int length)`

Parameters

`uint8_t value` the value to extract the bit field from

`int start` the lowest bit in the bit field (numbered from 0)

`int length` the length of the bit field

Description

Extract from the 8 bit input **value** the bit field specified by the **start** and **length** parameters, and return it. The bit field must lie entirely within the 8 bit word. It is valid to request that all 8 bits are returned (ie **length** 8 and **start** 0).

Return

the value of the bit field extracted from the input value.

`uint16_t extract16 (uint16_t value, int start, int length)`

Parameters

`uint16_t value` the value to extract the bit field from

`int start` the lowest bit in the bit field (numbered from 0)

`int length` the length of the bit field

Description

Extract from the 16 bit input **value** the bit field specified by the **start** and **length** parameters, and return it. The bit field must lie entirely within the 16 bit word. It is valid to request that all 16 bits are returned (ie **length** 16 and **start** 0).

Return

the value of the bit field extracted from the input value.

`uint64_t extract64 (uint64_t value, int start, int length)`

Parameters

`uint64_t value` the value to extract the bit field from

`int start` the lowest bit in the bit field (numbered from 0)

`int length` the length of the bit field

Description

Extract from the 64 bit input **value** the bit field specified by the **start** and **length** parameters, and return it. The bit field must lie entirely within the 64 bit word. It is valid to request that all 64 bits are returned (ie **length** 64 and **start** 0).

Return

the value of the bit field extracted from the input value.

`int32_t sextract32 (uint32_t value, int start, int length)`

Parameters

`uint32_t value` the value to extract the bit field from

`int start` the lowest bit in the bit field (numbered from 0)

int length the length of the bit field

Description

Extract from the 32 bit input **value** the bit field specified by the **start** and **length** parameters, and return it, sign extended to an `int32_t` (ie with the most significant bit of the field propagated to all the upper bits of the return value). The bit field must lie entirely within the 32 bit word. It is valid to request that all 32 bits are returned (ie **length** 32 and **start** 0).

Return

the sign extended value of the bit field extracted from the input value.

`int64_t sextract64 (uint64_t value, int start, int length)`

Parameters

uint64_t value the value to extract the bit field from

int start the lowest bit in the bit field (numbered from 0)

int length the length of the bit field

Description

Extract from the 64 bit input **value** the bit field specified by the **start** and **length** parameters, and return it, sign extended to an `int64_t` (ie with the most significant bit of the field propagated to all the upper bits of the return value). The bit field must lie entirely within the 64 bit word. It is valid to request that all 64 bits are returned (ie **length** 64 and **start** 0).

Return

the sign extended value of the bit field extracted from the input value.

`uint32_t deposit32 (uint32_t value, int start, int length, uint32_t fieldval)`

Parameters

uint32_t value initial value to insert bit field into

int start the lowest bit in the bit field (numbered from 0)

int length the length of the bit field

uint32_t fieldval the value to insert into the bit field

Description

Deposit **fieldval** into the 32 bit **value** at the bit field specified by the **start** and **length** parameters, and return the modified **value**. Bits of **value** outside the bit field are not modified. Bits of **fieldval** above the least significant **length** bits are ignored. The bit field must lie entirely within the 32 bit word. It is valid to request that all 32 bits are modified (ie **length** 32 and **start** 0).

Return

the modified **value**.

`uint64_t deposit64 (uint64_t value, int start, int length, uint64_t fieldval)`

Parameters

uint64_t value initial value to insert bit field into

int start the lowest bit in the bit field (numbered from 0)

int length the length of the bit field

uint64_t fieldval the value to insert into the bit field

Description

Deposit **fieldval** into the 64 bit **value** at the bit field specified by the **start** and **length** parameters, and return the modified **value**. Bits of **value** outside the bit field are not modified. Bits of **fieldval** above the least significant **length** bits are ignored. The bit field must lie entirely within the 64 bit word. It is valid to request that all 64 bits are modified (ie **length** 64 and **start** 0).

Return

the modified **value**.

uint32_t **half_shuffle32** (uint32_t *x*)

Parameters

uint32_t x 32-bit value (of which only the bottom 16 bits are of interest)

Description

Given an input value:

```
xxxx xxxx xxxx xxxx ABCD EFGH IJKL MNOP
```

return the value where the bottom 16 bits are spread out into the odd bits in the word, and the even bits are zeroed:

```
0A0B 0C0D 0E0F 0G0H 0I0J 0K0L 0M0N 0O0P
```

Any bits set in the top half of the input are ignored.

Return

the shuffled bits.

uint64_t **half_shuffle64** (uint64_t *x*)

Parameters

uint64_t x 64-bit value (of which only the bottom 32 bits are of interest)

Description

Given an input value:

```
xxxx xxxx xxxx .... xxxx xxxx ABCD EFGH IJKL MNOP QRST UVWX YZab cdef
```

return the value where the bottom 32 bits are spread out into the odd bits in the word, and the even bits are zeroed:

```
0A0B 0C0D 0E0F 0G0H 0I0J 0K0L 0M0N .... 0U0V 0W0X 0Y0Z 0a0b 0c0d 0e0f
```

Any bits set in the top half of the input are ignored.

Return

the shuffled bits.

uint32_t **half_unshuffle32** (uint32_t *x*)

Parameters

uint32_t x 32-bit value (of which only the odd bits are of interest)

Description

Given an input value:

```
xAxB xCxD xExF xGxH xIxJ xKxL xMxN xOxP
```

return the value where all the odd bits are compressed down into the low half of the word, and the high half is zeroed:

```
0000 0000 0000 0000 ABCD EFGH IJKL MNOP
```

Any even bits set in the input are ignored.

Return

the unshuffled bits.

`uint64_t half_unshuffle64 (uint64_t x)`

Parameters

uint64_t x 64-bit value (of which only the odd bits are of interest)

Description

Given an input value:

```
xAxB xCxD xExF xGxH xIxJ xKxL xMxN . . . . xUxV xWxX xYxZ xaxb xcxd xexf
```

return the value where all the odd bits are compressed down into the low half of the word, and the high half is zeroed:

```
0000 0000 0000 . . . . 0000 0000 ABCD EFGH IJKL MNOP QRST UVWX YZab cdef
```

Any even bits set in the input are ignored.

Return

the unshuffled bits.

6.23 Reset in QEMU: the Resettable interface

The reset of qemu objects is handled using the resettable interface declared in `include/hw/resettable.h`.

This interface allows objects to be grouped (on a tree basis); so that the whole group can be reset consistently. Each individual member object does not have to care about others; in particular, problems of order (which object is reset first) are addressed.

As of now DeviceClass and BusClass implement this interface.

6.23.1 Triggering reset

This section documents the APIs which “users” of a resettable object should use to control it. All resettable control functions must be called while holding the `iothread` lock.

You can apply a reset to an object using `resettable_assert_reset()`. You need to call `resettable_release_reset()` to release the object from reset. To instantly reset an object, without keeping it in reset state, just call `resettable_reset()`. These functions take two parameters: a pointer to the object to reset and a reset type.

Several types of reset will be supported. For now only cold reset is defined; others may be added later. The Resettable interface handles reset types with an enum:

RESET_TYPE_COLD Cold reset is supported by every resettable object. In QEMU, it means we reset to the initial state corresponding to the start of QEMU; this might differ from what is a real hardware cold reset. It differs from other resets (like warm or bus resets) which may keep certain parts untouched.

Calling `resettable_reset()` is equivalent to calling `resettable_assert_reset()` then `resettable_release_reset()`. It is possible to interleave multiple calls to these three functions. There may be several reset sources/controllers of a given object. The interface handles everything and the different reset controllers do not need to know anything about each others. The object will leave reset state only when each other controllers end their reset operation. This point is handled internally by maintaining a count of in-progress resets; it is crucial to call `resettable_release_reset()` one time and only one time per `resettable_assert_reset()` call.

For now migration of a device or bus in reset is not supported. Care must be taken not to delay `resettable_release_reset()` after its `resettable_assert_reset()` counterpart.

Note that, since resettable is an interface, the API takes a simple Object as parameter. Still, it is a programming error to call a resettable function on a non-resettable object and it will trigger a run time assert error. Since most calls to resettable interface are done through base class functions, such an error is not likely to happen.

For Devices and Buses, the following helper functions exist:

- `device_cold_reset()`
- `bus_cold_reset()`

These are simple wrappers around `resettable_reset()` function; they only cast the Device or Bus into an Object and pass the cold reset type. When possible prefer to use these functions instead of `resettable_reset()`.

Device and bus functions co-exist because there can be semantic differences between resetting a bus and resetting the controller bridge which owns it. For example, consider a SCSI controller. Resetting the controller puts all its registers back to what reset state was as well as reset everything on the SCSI bus, whereas resetting just the SCSI bus only resets everything that's on it but not the controller.

6.23.2 Multi-phase mechanism

This section documents the internals of the resettable interface.

The resettable interface uses a multi-phase system to relieve objects and machines from reset ordering problems. To address this, the reset operation of an object is split into three well defined phases.

When resetting several objects (for example the whole machine at simulation startup), all first phases of all objects are executed, then all second phases and then all third phases.

The three phases are:

1. The **enter** phase is executed when the object enters reset. It resets only local state of the object; it must not do anything that has a side-effect on other objects, such as raising or lowering a `qemu_irq` line or reading or writing guest memory.
2. The **hold** phase is executed for entry into reset, once every object in the group which is being reset has had its *enter* phase executed. At this point devices can do actions that affect other objects.
3. The **exit** phase is executed when the object leaves the reset state. Actions affecting other objects are permitted.

As said in previous section, the interface maintains a count of reset. This count is used to ensure phases are executed only when required. *enter* and *hold* phases are executed only when asserting reset for the first time (if an object is already in reset state when calling `resettable_assert_reset()` or `resettable_reset()`, they are not executed). The *exit* phase is executed only when the last reset operation ends. Therefore the object does not need to care how many of reset controllers it has and how many of them have started a reset.

6.23.3 Handling reset in a resettable object

This section documents the APIs that an implementation of a resettable object must provide and what functions it has access to. It is intended for people who want to implement or convert a class which has the resettable interface; for example when specializing an existing device or bus.

Methods to implement

Three methods should be defined or left empty. Each method corresponds to a phase of the reset; they are `name.phases.enter()`, `name.phases.hold()` and `name.phases.exit()`. They all take the object as parameter. The `enter` method also take the reset type as second parameter.

When extending an existing class, these methods may need to be extended too. The `resettable_class_set_parent_phases()` class function may be used to backup parent class methods.

Here follows an example to implement reset for a Device which sets an IO while in reset.

```
static void mydev_reset_enter(Object *obj, ResetType type)
{
    MyDevClass *myclass = MYDEV_GET_CLASS(obj);
    MyDevState *mydev = MYDEV(obj);
    /* call parent class enter phase */
    if (myclass->parent_phases.enter) {
        myclass->parent_phases.enter(obj, type);
    }
    /* initialize local state only */
    mydev->var = 0;
}

static void mydev_reset_hold(Object *obj)
{
    MyDevClass *myclass = MYDEV_GET_CLASS(obj);
    MyDevState *mydev = MYDEV(obj);
    /* call parent class hold phase */
    if (myclass->parent_phases.hold) {
        myclass->parent_phases.hold(obj);
    }
    /* set an IO */
    qemu_set_irq(mydev->irq, 1);
}

static void mydev_reset_exit(Object *obj)
{
    MyDevClass *myclass = MYDEV_GET_CLASS(obj);
    MyDevState *mydev = MYDEV(obj);
    /* call parent class exit phase */
    if (myclass->parent_phases.exit) {
        myclass->parent_phases.exit(obj);
    }
    /* clear an IO */
    qemu_set_irq(mydev->irq, 0);
}

typedef struct MyDevClass {
    MyParentClass parent_class;
    /* to store eventual parent reset methods */

```

(continues on next page)

(continued from previous page)

```

    ResettablePhases parent_phases;
} MyDevClass;

static void mydev_class_init(ObjectClass *class, void *data)
{
    MyDevClass *myclass = MYDEV_CLASS(class);
    ResettableClass *rc = RESETTABLE_CLASS(class);
    resettable_class_set_parent_reset_phases(rc,
                                             mydev_reset_enter,
                                             mydev_reset_hold,
                                             mydev_reset_exit,
                                             &myclass->parent_phases);
}

```

In the above example, we override all three phases. It is possible to override only some of them by passing `NULL` instead of a function pointer to `resettable_class_set_parent_reset_phases()`. For example, the following will only override the *enter* phase and leave *hold* and *exit* untouched:

```

resettable_class_set_parent_reset_phases(rc, mydev_reset_enter,
                                         NULL, NULL,
                                         &myclass->parent_phases);

```

This is equivalent to providing a trivial implementation of the hold and exit phases which does nothing but call the parent class's implementation of the phase.

Polling the reset state

Resettable interface provides the `resettable_is_in_reset()` function. This function returns true if the object parameter is currently under reset.

An object is under reset from the beginning of the *init* phase to the end of the *exit* phase. During all three phases, the function will return that the object is in reset.

This function may be used if the object behavior has to be adapted while in reset state. For example if a device has an irq input, it will probably need to ignore it while in reset; then it can for example check the reset state at the beginning of the irq callback.

Note that until migration of the reset state is supported, an object should not be left in reset. So apart from being currently executing one of the reset phases, the only cases when this function will return true is if an external interaction (like changing an io) is made during *hold* or *exit* phase of another object in the same reset group.

Helpers `device_is_in_reset()` and `bus_is_in_reset()` are also provided for devices and buses and should be preferred.

6.23.4 Base class handling of reset

This section documents parts of the reset mechanism that you only need to know about if you are extending it to work with a new base class other than `DeviceClass` or `BusClass`, or maintaining the existing code in those classes. Most people can ignore it.

Methods to implement

There are two other methods that need to exist in a class implementing the interface: `get_state()` and `child_foreach()`.

`get_state()` is simple. *resettable* is an interface and, as a consequence, does not have any class state structure. But in order to factorize the code, we need one. This method must return a pointer to `ResettableState` structure. The structure must be allocated by the base class; preferably it should be located inside the object instance structure.

`child_foreach()` is more complex. It should execute the given callback on every reset child of the given resettable object. All children must be resettable too. Additional parameters (a reset type and an opaque pointer) must be passed to the callback too.

In `DeviceClass` and `BusClass` the `ResettableState` is located `DeviceState` and `BusState` structure. `child_foreach()` is implemented to follow the bus hierarchy; for a bus, it calls the function on every child device; for a device, it calls the function on every bus child. When we reset the main system bus, we reset the whole machine bus tree.

Changing a resettable parent

One thing which should be taken care of by the base class is handling reset hierarchy changes.

The reset hierarchy is supposed to be static and built during machine creation. But there are actually some exceptions. To cope with this, the resettable API provides `resettable_change_parent()`. This function allows to set, update or remove the parent of a resettable object after machine creation is done. As parameters, it takes the object being moved, the old parent if any and the new parent if any.

This function can be used at any time when not in a reset operation. During a reset operation it must be used only in *hold* phase. Using it in *enter* or *exit* phase is an error. Also it should not be used during machine creation, although it is harmless to do so: the function is a no-op as long as old and new parent are NULL or not in reset.

There is currently 2 cases where this function is used:

1. *device hotplug*; it means a new device is introduced on a live bus.
2. *hot bus change*; it means an existing live device is added, moved or removed in the bus hierarchy. At the moment, it occurs only in the raspi machines for changing the sdbus used by sd card.

6.24 Booting from real channel-attached devices on s390x

6.24.1 s390 hardware IPL

The s390 hardware IPL process consists of the following steps.

1. A READ IPL ccw is constructed in memory location `0x0`. This ccw, by definition, reads the IPL1 record which is located on the disk at cylinder 0 track 0 record 1. Note that the chain flag is on in this ccw so when it is complete another ccw will be fetched and executed from memory location `0x08`.
2. Execute the Read IPL ccw at `0x00`, thereby reading IPL1 data into `0x00`. IPL1 data is 24 bytes in length and consists of the following pieces of information: `[psw][read ccw][tic ccw]`. When the machine executes the Read IPL ccw it read the 24-bytes of IPL1 to be read into memory starting at location `0x0`. Then the ccw program at `0x08` which consists of a read ccw and a tic ccw is automatically executed because of the chain flag from the original READ IPL ccw. The read ccw will read the IPL2 data into memory and the TIC (Transfer In Channel) will transfer control to the channel program contained in the IPL2 data. The TIC channel command is the equivalent of a branch/jump/goto instruction for channel programs.

NOTE: The ccws in IPL1 are defined by the architecture to be format 0.

3. Execute IPL2. The TIC ccw instruction at the end of the IPL1 channel program will begin the execution of the IPL2 channel program. IPL2 is stage-2 of the boot process and will contain a larger channel program than IPL1. The point of IPL2 is to find and load either the operating system or a small program that loads the operating system from disk. At the end of this step all or some of the real operating system is loaded into memory and we

are ready to hand control over to the guest operating system. At this point the guest operating system is entirely responsible for loading any more data it might need to function.

NOTE: The IPL2 channel program might read data into memory location 0x0 thereby overwriting the IPL1 psw and channel program. This is ok as long as the data placed in location 0x0 contains a psw whose instruction address points to the guest operating system code to execute at the end of the IPL/boot process.

NOTE: The ccws in IPL2 are defined by the architecture to be format 0.

4. Start executing the guest operating system. The psw that was loaded into memory location 0x0 as part of the ipl process should contain the needed flags for the operating system we have loaded. The psw's instruction address will point to the location in memory where we want to start executing the operating system. This psw is loaded (via LPSW instruction) causing control to be passed to the operating system code.

In a non-virtualized environment this process, handled entirely by the hardware, is kicked off by the user initiating a "Load" procedure from the hardware management console. This "Load" procedure crafts a special "Read IPL" ccw in memory location 0x0 that reads IPL1. It then executes this ccw thereby kicking off the reading of IPL1 data. Since the channel program from IPL1 will be written immediately after the special "Read IPL" ccw, the IPL1 channel program will be executed immediately (the special read ccw has the chaining bit turned on). The TIC at the end of the IPL1 channel program will cause the IPL2 channel program to be executed automatically. After this sequence completes the "Load" procedure then loads the psw from 0x0.

6.24.2 How this all pertains to QEMU (and the kernel)

In theory we should merely have to do the following to IPL/boot a guest operating system from a DASD device:

1. Place a "Read IPL" ccw into memory location 0x0 with chaining bit on.
2. Execute channel program at 0x0.
3. LPSW 0x0.

However, our emulation of the machine's channel program logic within the kernel is missing one key feature that is required for this process to work: non-prefetch of ccw data.

When we start a channel program we pass the channel subsystem parameters via an ORB (Operation Request Block). One of those parameters is a prefetch bit. If the bit is on then the vfio-ccw kernel driver is allowed to read the entire channel program from guest memory before it starts executing it. This means that any channel commands that read additional channel commands will not work as expected because the newly read commands will only exist in guest memory and NOT within the kernel's channel subsystem memory. The kernel vfio-ccw driver currently requires this bit to be on for all channel programs. This is a problem because the IPL process consists of transferring control from the "Read IPL" ccw immediately to the IPL1 channel program that was read by "Read IPL".

Not being able to turn off prefetch will also prevent the TIC at the end of the IPL1 channel program from transferring control to the IPL2 channel program.

Lastly, in some cases (the ziplt bootloader for example) the IPL2 program also transfers control to another channel program segment immediately after reading it from the disk. So we need to be able to handle this case.

6.24.3 What QEMU does

Since we are forced to live with prefetch we cannot use the very simple IPL procedure we defined in the preceding section. So we compensate by doing the following.

1. Place "Read IPL" ccw into memory location 0x0, but turn off chaining bit.
2. Execute "Read IPL" at 0x0.

So now IPL1's psw is at 0x0 and IPL1's channel program is at 0x08.

3. Write a custom channel program that will seek to the IPL2 record and then execute the READ and TIC ccws from IPL1. Normally the seek is not required because after reading the IPL1 record the disk is automatically positioned to read the very next record which will be IPL2. But since we are not reading both IPL1 and IPL2 as part of the same channel program we must manually set the position.
4. Grab the target address of the TIC instruction from the IPL1 channel program. This address is where the IPL2 channel program starts.

Now IPL2 is loaded into memory somewhere, and we know the address.

5. Execute the IPL2 channel program at the address obtained in step #4.

Because this channel program can be dynamic, we must use a special algorithm that detects a READ immediately followed by a TIC and breaks the ccw chain by turning off the chain bit in the READ ccw. When control is returned from the kernel/hardware to the QEMU bios code we immediately issue another start subchannel to execute the remaining TIC instruction. This causes the entire channel program (starting from the TIC) and all needed data to be refetched thereby stepping around the limitation that would otherwise prevent this channel program from executing properly.

Now the operating system code is loaded somewhere in guest memory and the psw in memory location 0x0 will point to entry code for the guest operating system.

6. LPSW 0x0

LPSW transfers control to the guest operating system and we're done.

6.25 Modelling a clock tree in QEMU

6.25.1 What are clocks?

Clocks are QOM objects developed for the purpose of modelling the distribution of clocks in QEMU.

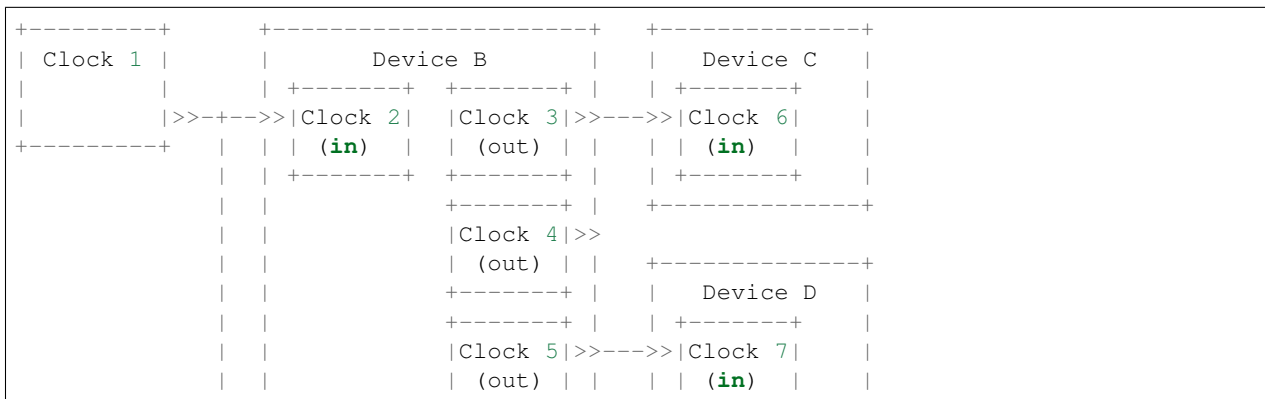
They allow us to model the clock distribution of a platform and detect configuration errors in the clock tree such as badly configured PLL, clock source selection or disabled clock.

The object is *Clock* and its QOM name is `clock` (in C code, the macro `TYPE_CLOCK`).

Clocks are typically used with devices where they are used to model inputs and outputs. They are created in a similar way to GPIOs. Inputs and outputs of different devices can be connected together.

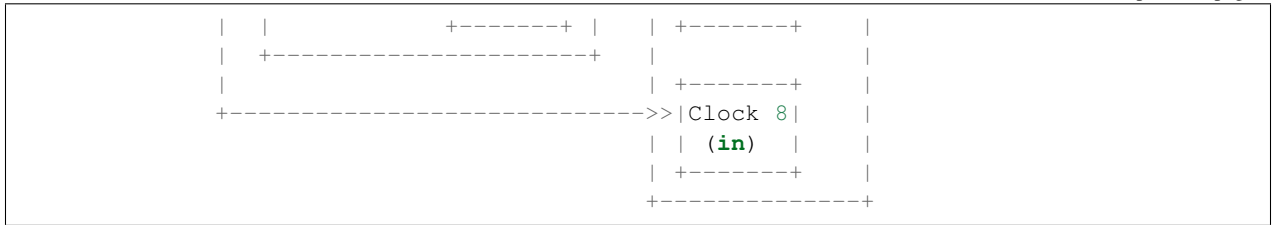
In these cases a Clock object is a child of a Device object, but this is not a requirement. Clocks can be independent of devices. For example it is possible to create a clock outside of any device to model the main clock source of a machine.

Here is an example of clocks:



(continues on next page)

(continued from previous page)



Clocks are defined in the `include/hw/clock.h` header and device related functions are defined in the `include/hw/qdev-clock.h` header.

6.25.2 The clock state

The state of a clock is its period; it is stored as an integer representing it in units of 2^{-32} ns. The special value of 0 is used to represent the clock being inactive or gated. The clocks do not model the signal itself (pin toggling) or other properties such as the duty cycle.

All clocks contain this state: outputs as well as inputs. This allows the current period of a clock to be fetched at any time. When a clock is updated, the value is immediately propagated to all connected clocks in the tree.

To ease interaction with clocks, helpers with a unit suffix are defined for every clock state setter or getter. The suffixes are:

- `_ns` for handling periods in nanoseconds
- `_hz` for handling frequencies in hertz

The 0 period value is converted to 0 in hertz and vice versa. 0 always means that the clock is disabled.

6.25.3 Adding a new clock

Adding clocks to a device must be done during the init method of the Device instance.

To add an input clock to a device, the function `qdev_init_clock_in()` must be used. It takes the name, a callback and an opaque parameter for the callback (this will be explained in a following section). Output is simpler; only the name is required. Typically:

```
qdev_init_clock_in(DEVICE(dev), "clk_in", clk_in_callback, dev);
qdev_init_clock_out(DEVICE(dev), "clk_out");
```

Both functions return the created Clock pointer, which should be saved in the device's state structure for further use.

These objects will be automatically deleted by the QOM reference mechanism.

Note that it is possible to create a static array describing clock inputs and outputs. The function `qdev_init_clocks()` must be called with the array as parameter to initialize the clocks: it has the same behaviour as calling the `qdev_init_clock_in/out()` for each clock in the array. To ease the array construction, some macros are defined in `include/hw/qdev-clock.h`. As an example, the following creates 2 clocks to a device: one input and one output.

```
/* device structure containing pointers to the clock objects */
typedef struct MyDeviceState {
    DeviceState parent_obj;
    Clock *clk_in;
    Clock *clk_out;
} MyDeviceState;
```

(continues on next page)

(continued from previous page)

```

/*
 * callback for the input clock (see "Callback on input clock
 * change" section below for more information).
 */
static void clk_in_callback(void *opaque);

/*
 * static array describing clocks:
 * + a clock input named "clk_in", whose pointer is stored in
 *   the clk_in field of a MyDeviceState structure with callback
 *   clk_in_callback.
 * + a clock output named "clk_out" whose pointer is stored in
 *   the clk_out field of a MyDeviceState structure.
 */
static const ClockPortInitArray mydev_clocks = {
    QDEV_CLOCK_IN(MyDeviceState, clk_in, clk_in_callback),
    QDEV_CLOCK_OUT(MyDeviceState, clk_out),
    QDEV_CLOCK_END
};

/* device initialization function */
static void mydev_init(Object *obj)
{
    /* cast to MyDeviceState */
    MyDeviceState *mydev = MYDEVICE(obj);
    /* create and fill the pointer fields in the MyDeviceState */
    qdev_init_clocks(mydev, mydev_clocks);
    [...]
}

```

An alternative way to create a clock is to simply call `object_new(TYPE_CLOCK)`. In that case the clock will neither be an input nor an output of a device. After the whole QOM hierarchy of the clock has been set `clock_setup_canonical_path()` should be called.

At creation, the period of the clock is 0: the clock is disabled. You can change it using `clock_set_ns()` or `clock_set_hz()`.

Note that if you are creating a clock with a fixed period which will never change (for example the main clock source of a board), then you'll have nothing else to do. This value will be propagated to other clocks when connecting the clocks together and devices will fetch the right value during the first reset.

6.25.4 Retrieving clocks from a device

`qdev_get_clock_in()` and `dev_get_clock_out()` are available to get the clock inputs or outputs of a device. For example:

```
Clock *clk = qdev_get_clock_in(DEVICE(mydev), "clk_in");
```

or:

```
Clock *clk = dev_get_clock_out(DEVICE(mydev), "clk_out");
```


(continued from previous page)

```

return;
}

```

Note that this only checks that the clock has been wired up; it is still possible that the output clock connected to it is disabled or has not yet been configured, in which case the period will be zero. You should use the clock callback to find out when the clock period changes.

6.25.7 Fetching clock frequency/period

To get the current state of a clock, use the functions `clock_get()` or `clock_get_hz()`.

`clock_get()` returns the period of the clock in its fully precise internal representation, as an unsigned 64-bit integer in units of 2^{32} nanoseconds. (For many purposes `clock_ticks_to_ns()` will be more convenient; see the section below on expiry deadlines.)

`clock_get_hz()` returns the frequency of the clock, rounded to the next lowest integer. This implies some inaccuracy due to the rounding, so be cautious about using it in calculations.

It is also possible to register a callback on clock frequency changes. Here is an example:

```

void clock_callback(void *opaque) {
    MyDeviceState *s = (MyDeviceState *) opaque;
    /*
     * 'opaque' is the argument passed to qdev_init_clock_in();
     * usually this will be the device state pointer.
     */

    /* do something with the new period */
    fprintf(stdout, "device new period is %" PRIu64 " * 2^32 ns\n",
            clock_get(dev->my_clk_input));
}

```

If you are only interested in the frequency for displaying it to humans (for instance in debugging), use `clock_display_freq()`, which returns a prettified string-representation, e.g. “33.3 MHz”. The caller must free the string with `g_free()` after use.

6.25.8 Calculating expiry deadlines

A commonly required operation for a clock is to calculate how long it will take for the clock to tick N times; this can then be used to set a timer expiry deadline. Use the function `clock_ticks_to_ns()`, which takes an unsigned 64-bit count of ticks and returns the length of time in nanoseconds required for the clock to tick that many times.

It is important not to try to calculate expiry deadlines using a shortcut like multiplying a “period of clock in nanoseconds” value by the tick count, because clocks can have periods which are not a whole number of nanoseconds, and the accumulated error in the multiplication can be significant.

For a clock with a very long period and a large number of ticks, the result of this function could in theory be too large to fit in a 64-bit value. To avoid overflow in this case, `clock_ticks_to_ns()` saturates the result to `INT64_MAX` (because this is the largest valid input to the QEMUTimer APIs). Since `INT64_MAX` nanoseconds is almost 300 years, anything with an expiry later than that is in the “will never happen” category. Callers of `clock_ticks_to_ns()` should therefore generally not special-case the possibility of a saturated result but just allow the timer to be set to that far-future value. (If you are performing further calculations on the returned value rather than simply passing it to a QEMUTimer function like `timer_mod_ns()` then you should be careful to avoid overflow in those calculations, of course.)

6.25.9 Changing a clock period

A device can change its outputs using the `clock_update()`, `clock_update_ns()` or `clock_update_hz()` function. It will trigger updates on every connected input.

For example, let's say that we have an output clock `clkout` and we have a pointer to it in the device state because we did the following in init phase:

```
dev->clkout = qdev_init_clock_out(DEVICE(dev), "clkout");
```

Then at any time (apart from the cases listed below), it is possible to change the clock value by doing:

```
clock_update_hz(dev->clkout, 1000 * 1000 * 1000); /* 1GHz */
```

Because updating a clock may trigger any side effects through connected clocks and their callbacks, this operation must be done while holding the qemu io lock.

For the same reason, one can update clocks only when it is allowed to have side effects on other objects. In consequence, it is forbidden:

- during migration,
- and in the enter phase of reset.

Note that calling `clock_update[_ns|_hz]()` is equivalent to calling `clock_set[_ns|_hz]()` (with the same arguments) then `clock_propagate()` on the clock. Thus, setting the clock value can be separated from triggering the side-effects. This is often required to factorize code to handle reset and migration in devices.

6.25.10 Aliasing clocks

Sometimes, one needs to forward, or inherit, a clock from another device. Typically, when doing device composition, a device might expose a sub-device's clock without interfering with it. The function `qdev_alias_clock()` can be used to achieve this behaviour. Note that it is possible to expose the clock under a different name. `qdev_alias_clock()` works for both input and output clocks.

For example, if device B is a child of device A, `device_a_instance_init()` may do something like this:

```
void device_a_instance_init(Object *obj)
{
    AState *A = DEVICE_A(obj);
    BState *B;
    /* create object B as child of A */
    [...]
    qdev_alias_clock(B, "clk", A, "b_clk");
    /*
     * Now A has a clock "b_clk" which is an alias to
     * the clock "clk" of its child B.
     */
}
```

This function does not return any clock object. The new clock has the same direction (input or output) as the original one. This function only adds a link to the existing clock. In the above example, object B remains the only object allowed to use the clock and device A must not try to change the clock period or set a callback to the clock. This diagram describes the example with an input clock:

```
+-----+
|           |
|   Device A   |
|           |
+-----+
```

(continues on next page)

(continued from previous page)

```

|          +-----+ |
|          | Device B | |
|          | +-----+ | |
>>"b_clk">>>| "clk" | | |
| (in) | | (in) | | |
|          | +-----+ | |
|          +-----+ |
+-----+

```

6.25.11 Migration

Clock state is not migrated automatically. Every device must handle its clock migration. Alias clocks must not be migrated.

To ensure clock states are restored correctly during migration, there are two solutions.

Clock states can be migrated by adding an entry into the device vmstate description. You should use the `VMSTATE_CLOCK` macro for this. This is typically used to migrate an input clock state. For example:

```

MyDeviceState {
    DeviceState parent_obj;
    [...] /* some fields */
    Clock *clk;
};

VMStateDescription my_device_vmstate = {
    .name = "my_device",
    .fields = (VMStateField[]) {
        [...], /* other migrated fields */
        VMSTATE_CLOCK(clk, MyDeviceState),
        VMSTATE_END_OF_LIST()
    }
};

```

The second solution is to restore the clock state using information already at our disposal. This can be used to restore output clock states using the device state. The functions `clock_set[_ns|_hz]()` can be used during the `post_load()` migration callback.

When adding clock support to an existing device, if you care about migration compatibility you will need to be careful, as simply adding a `VMSTATE_CLOCK()` line will break compatibility. Instead, you can put the `VMSTATE_CLOCK()` line into a vmstate subsection with a suitable `needed` function, and use `clock_set()` in a `pre_load()` function to set the default value that will be used if the source virtual machine in the migration does not send the clock state.

Care should be taken not to use `clock_update[_ns|_hz]()` or `clock_propagate()` during the whole migration procedure because it will trigger side effects to other devices in an unknown state.

6.26 The QEMU Object Model (QOM)

The QEMU Object Model provides a framework for registering user creatable types and instantiating objects from those types. QOM provides the following features:

- System for dynamically registering types
- Support for single-inheritance of types

- Multiple inheritance of stateless interfaces

Listing 1: Creating a minimal type

```
#include "qdev.h"

#define TYPE_MY_DEVICE "my-device"

// No new virtual functions: we can reuse the typedef for the
// superclass.
typedef DeviceClass MyDeviceClass;
typedef struct MyDevice
{
    DeviceState parent;

    int reg0, reg1, reg2;
} MyDevice;

static const TypeInfo my_device_info = {
    .name = TYPE_MY_DEVICE,
    .parent = TYPE_DEVICE,
    .instance_size = sizeof(MyDevice),
};

static void my_device_register_types(void)
{
    type_register_static(&my_device_info);
}

type_init(my_device_register_types)
```

In the above example, we create a simple type that is described by #TypeInfo. #TypeInfo describes information about the type including what it inherits from, the instance and class size, and constructor/destructor hooks.

Alternatively several static types could be registered using helper macro DEFINE_TYPES()

```
static const TypeInfo device_types_info[] = {
    {
        .name = TYPE_MY_DEVICE_A,
        .parent = TYPE_DEVICE,
        .instance_size = sizeof(MyDeviceA),
    },
    {
        .name = TYPE_MY_DEVICE_B,
        .parent = TYPE_DEVICE,
        .instance_size = sizeof(MyDeviceB),
    },
};

DEFINE_TYPES(device_types_info)
```

Every type has an #ObjectClass associated with it. #ObjectClass derivatives are instantiated dynamically but there is only ever one instance for any given type. The #ObjectClass typically holds a table of function pointers for the virtual methods implemented by this type.

Using object_new(), a new #Object derivative will be instantiated. You can cast an #Object to a subclass (or base-class) type using object_dynamic_cast(). You typically want to define macro wrappers around OBJECT_CHECK() and OBJECT_CLASS_CHECK() to make it easier to convert to a specific type:

Listing 2: Typecasting macros

```
#define MY_DEVICE_GET_CLASS(obj) \
    OBJECT_GET_CLASS(MyDeviceClass, obj, TYPE_MY_DEVICE)
#define MY_DEVICE_CLASS(klass) \
    OBJECT_CLASS_CHECK(MyDeviceClass, klass, TYPE_MY_DEVICE)
#define MY_DEVICE(obj) \
    OBJECT_CHECK(MyDevice, obj, TYPE_MY_DEVICE)
```

6.26.1 Class Initialization

Before an object is initialized, the class for the object must be initialized. There is only one class object for all instance objects that is created lazily.

Classes are initialized by first initializing any parent classes (if necessary). After the parent class object has initialized, it will be copied into the current class object and any additional storage in the class object is zero filled.

The effect of this is that classes automatically inherit any virtual function pointers that the parent class has already initialized. All other fields will be zero filled.

Once all of the parent classes have been initialized, #TypeInfo::class_init is called to let the class being instantiated provide default initialize for its virtual functions. Here is how the above example might be modified to introduce an overridden virtual function:

Listing 3: Overriding a virtual function

```
#include "qdev.h"

void my_device_class_init(ObjectClass *klass, void *class_data)
{
    DeviceClass *dc = DEVICE_CLASS(klass);
    dc->reset = my_device_reset;
}

static const TypeInfo my_device_info = {
    .name = TYPE_MY_DEVICE,
    .parent = TYPE_DEVICE,
    .instance_size = sizeof(MyDevice),
    .class_init = my_device_class_init,
};
```

Introducing new virtual methods requires a class to define its own struct and to add a .class_size member to the #TypeInfo. Each method will also have a wrapper function to call it easily:

Listing 4: Defining an abstract class

```
#include "qdev.h"

typedef struct MyDeviceClass
{
    DeviceClass parent;

    void (*froblicate) (MyDevice *obj);
} MyDeviceClass;

static const TypeInfo my_device_info = {
```

(continues on next page)

(continued from previous page)

```

.name = TYPE_MY_DEVICE,
.parent = TYPE_DEVICE,
.instance_size = sizeof(MyDevice),
.abstract = true, // or set a default in my_device_class_init
.class_size = sizeof(MyDeviceClass),
};

void my_device_frobnicate(MyDevice *obj)
{
    MyDeviceClass *klass = MY_DEVICE_GET_CLASS(obj);

    klass->frobnicate(obj);
}

```

6.26.2 Interfaces

Interfaces allow a limited form of multiple inheritance. Instances are similar to normal types except for the fact that they are only defined by their classes and never carry any state. As a consequence, a pointer to an interface instance should always be of incomplete type in order to be sure it cannot be dereferenced. That is, you should define the ‘typedef struct SomethingIf SomethingIf’ so that you can pass around `SomethingIf *si` arguments, but not define a `struct SomethingIf { ... }`. The only things you can validly do with a `SomethingIf *` are to pass it as an argument to a method on its corresponding `SomethingIfClass`, or to dynamically cast it to an object that implements the interface.

6.26.3 Methods

A *method* is a function within the namespace scope of a class. It usually operates on the object instance by passing it as a strongly-typed first argument. If it does not operate on an object instance, it is dubbed *class method*.

Methods cannot be overloaded. That is, the #ObjectClass and method name uniquely identify the function to be called; the signature does not vary except for trailing varargs.

Methods are always *virtual*. Overriding a method in #TypeInfo.class_init of a subclass leads to any user of the class obtained via OBJECT_GET_CLASS() accessing the overridden function. The original function is not automatically invoked. It is the responsibility of the overriding class to determine whether and when to invoke the method being overridden.

To invoke the method being overridden, the preferred solution is to store the original value in the overriding class before overriding the method. This corresponds to `{super,base}.method(...)` in Java and C# respectively; this frees the overriding class from hardcoding its parent class, which someone might choose to change at some point.

Listing 5: Overriding a virtual method

```

typedef struct MyState MyState;

typedef void (*MyDoSomething) (MyState *obj);

typedef struct MyClass {
    ObjectClass parent_class;

    MyDoSomething do_something;
} MyClass;

static void my_do_something(MyState *obj)

```

(continues on next page)

(continued from previous page)

```

{
    // do something
}

static void my_class_init(ObjectClass *oc, void *data)
{
    MyClass *mc = MY_CLASS(oc);

    mc->do_something = my_do_something;
}

static const TypeInfo my_type_info = {
    .name = TYPE_MY,
    .parent = TYPE_OBJECT,
    .instance_size = sizeof(MyState),
    .class_size = sizeof(MyClass),
    .class_init = my_class_init,
};

typedef struct DerivedClass {
    MyClass parent_class;

    MyDoSomething parent_do_something;
} DerivedClass;

static void derived_do_something(MyState *obj)
{
    DerivedClass *dc = DERIVED_GET_CLASS(obj);

    // do something here
    dc->parent_do_something(obj);
    // do something else here
}

static void derived_class_init(ObjectClass *oc, void *data)
{
    MyClass *mc = MY_CLASS(oc);
    DerivedClass *dc = DERIVED_CLASS(oc);

    dc->parent_do_something = mc->do_something;
    mc->do_something = derived_do_something;
}

static const TypeInfo derived_type_info = {
    .name = TYPE_DERIVED,
    .parent = TYPE_MY,
    .class_size = sizeof(DerivedClass),
    .class_init = derived_class_init,
};

```

Alternatively, `object_class_by_name()` can be used to obtain the class and its non-overridden methods for a specific type. This would correspond to `MyClass::method(...)` in C++.

The first example of such a QOM method was `#CPUClass.reset`, another example is `#DeviceClass.realize`.

6.26.4 Standard type declaration and definition macros

A lot of the code outlined above follows a standard pattern and naming convention. To reduce the amount of boilerplate code that needs to be written for a new type there are two sets of macros to generate the common parts in a standard format.

A type is declared using the `OBJECT_DECLARE` macro family. In types which do not require any virtual functions in the class, the `OBJECT_DECLARE_SIMPLE_TYPE` macro is suitable, and is commonly placed in the header file:

Listing 6: Declaring a simple type

```
OBJECT_DECLARE_SIMPLE_TYPE(MyDevice, my_device,
                           MY_DEVICE, DEVICE)
```

This is equivalent to the following:

Listing 7: Expansion from declaring a simple type

```
typedef struct MyDevice MyDevice;
typedef struct MyDeviceClass MyDeviceClass;

G_DEFINE_AUTOPTR_CLEANUP_FUNC(MyDeviceClass, object_unref)

#define MY_DEVICE_GET_CLASS(void *obj) \
    OBJECT_GET_CLASS(MyDeviceClass, obj, TYPE_MY_DEVICE)
#define MY_DEVICE_CLASS(void *klass) \
    OBJECT_CLASS_CHECK(MyDeviceClass, klass, TYPE_MY_DEVICE)
#define MY_DEVICE(void *obj) \
    OBJECT_CHECK(MyDevice, obj, TYPE_MY_DEVICE)

struct MyDeviceClass {
    DeviceClass parent_class;
};
```

The ‘struct MyDevice’ needs to be declared separately. If the type requires virtual functions to be declared in the class struct, then the alternative `OBJECT_DECLARE_TYPE()` macro can be used. This does the same as `OBJECT_DECLARE_SIMPLE_TYPE()`, but without the ‘struct MyDeviceClass’ definition.

To implement the type, the `OBJECT_DEFINE` macro family is available. In the simple case the `OBJECT_DEFINE_TYPE` macro is suitable:

Listing 8: Defining a simple type

```
OBJECT_DEFINE_TYPE(MyDevice, my_device, MY_DEVICE, DEVICE)
```

This is equivalent to the following:

Listing 9: Expansion from defining a simple type

```
static void my_device_finalize(Object *obj);
static void my_device_class_init(ObjectClass *oc, void *data);
static void my_device_init(Object *obj);

static const TypeInfo my_device_info = {
    .parent = TYPE_DEVICE,
    .name = TYPE_MY_DEVICE,
    .instance_size = sizeof(MyDevice),
    .instance_init = my_device_init,
```

(continues on next page)

(continued from previous page)

```
.instance_finalize = my_device_finalize,
.class_size = sizeof(MyDeviceClass),
.class_init = my_device_class_init,
};

static void
my_device_register_types(void)
{
    type_register_static(&my_device_info);
}
type_init(my_device_register_types);
```

This is sufficient to get the type registered with the type system, and the three standard methods now need to be implemented along with any other logic required for the type.

If the type needs to implement one or more interfaces, then the `OBJECT_DEFINE_TYPE_WITH_INTERFACES()` macro can be used instead. This accepts an array of interface type names.

Listing 10: Defining a simple type implementing interfaces

```
OBJECT_DEFINE_TYPE_WITH_INTERFACES(MyDevice, my_device,
                                   MY_DEVICE, DEVICE,
                                   { TYPE_USER_CREATABLE },
                                   { NULL })
```

If the type is not intended to be instantiated, then then the `OBJECT_DEFINE_ABSTRACT_TYPE()` macro can be used instead:

Listing 11: Defining a simple abstract type

```
OBJECT_DEFINE_ABSTRACT_TYPE(MyDevice, my_device,
                             MY_DEVICE, DEVICE)
```

API Reference

ObjectPropertyAccessor

Typedef:

Syntax

```
void ObjectPropertyAccessor (Object *obj, Visitor *v, const char
                             *name, void *opaque, Error **errp)
```

Parameters

Object *obj the object that owns the property

Visitor *v the visitor that contains the property data

const char *name the name of the property

void *opaque the object property opaque

Error **errp a pointer to an Error that is filled if getting/setting fails.

Description

Called when trying to get/set a property.

ObjectPropertyResolve**Typedef:****Syntax**

```
Object * ObjectPropertyResolve (Object *obj, void *opaque, const char
*part)
```

Parameters**Object *obj** the object that owns the property**void *opaque** the opaque registered with the property**const char *part** the name of the property**Description**Resolves the *Object* corresponding to property **part**.The returned object can also be used as a starting point to resolve a relative path starting with “**part**”.**Return**If **path** is the path that led to **obj**, the function returns the *Object* corresponding to “**path/part**”. If “**path/part**” is not a valid object path, it returns NULL.**ObjectPropertyRelease****Typedef:****Syntax**

```
void ObjectPropertyRelease (Object *obj, const char *name, void
*opaque)
```

Parameters**Object *obj** the object that owns the property**const char *name** the name of the property**void *opaque** the opaque registered with the property**Description**

Called when a property is removed from a object.

ObjectPropertyInit**Typedef:****Syntax**

```
void ObjectPropertyInit (Object *obj, ObjectProperty *prop)
```

Parameters**Object *obj** the object that owns the property**ObjectProperty *prop** the property to set**Description**

Called when a property is initialized.

ObjectUnparent**Typedef:****Syntax**

```
void ObjectUnparent (Object *obj)
```

Parameters

Object *obj the object that is being removed from the composition tree

Description

Called when an object is being removed from the QOM composition tree. The function should remove any backlinks from children objects to **obj**.

ObjectFree

Typedef:

Syntax

```
void ObjectFree (void *obj)
```

Parameters

void *obj the object being freed

Description

Called when an object's last reference is removed.

struct **ObjectClass**

Definition

```
struct ObjectClass {  
};
```

Members**Description**

The base for all classes. The only thing that *ObjectClass* contains is an integer type handle.

struct **Object**

Definition

```
struct Object {  
};
```

Members**Description**

The base for all objects. The first member of this object is a pointer to a *ObjectClass*. Since C guarantees that the first member of a structure always begins at byte 0 of that structure, as long as any sub-object places its parent as the first member, we can cast directly to a *Object*.

As a result, *Object* contains a reference to the objects type as its first member. This allows identification of the real type of the object at run time.

DECLARE_INSTANCE_CHECKER (InstanceType, OBJ_NAME, TYPENAME)

Parameters

InstanceType instance struct name

OBJ_NAME the object name in uppercase with underscore separators

TYPENAME type name

Description

Direct usage of this macro should be avoided, and the complete `OBJECT_DECLARE_TYPE` macro is recommended instead.

This macro will provide the instance type cast functions for a QOM type.

DECLARE_CLASS_CHECKERS (ClassType, OBJ_NAME, TYPENAME)

Parameters

ClassType class struct name

OBJ_NAME the object name in uppercase with underscore separators

TYPENAME type name

Description

Direct usage of this macro should be avoided, and the complete `OBJECT_DECLARE_TYPE` macro is recommended instead.

This macro will provide the class type cast functions for a QOM type.

DECLARE_OBJ_CHECKERS (InstanceType, ClassType, OBJ_NAME, TYPENAME)

Parameters

InstanceType instance struct name

ClassType class struct name

OBJ_NAME the object name in uppercase with underscore separators

TYPENAME type name

Description

Direct usage of this macro should be avoided, and the complete `OBJECT_DECLARE_TYPE` macro is recommended instead.

This macro will provide the three standard type cast functions for a QOM type.

OBJECT_DECLARE_TYPE (InstanceType, ClassType, MODULE_OBJ_NAME)

Parameters

InstanceType instance struct name

ClassType class struct name

MODULE_OBJ_NAME the object name in uppercase with underscore separators

Description

This macro is typically used in a header file, and will:

- create the typedefs for the object and class structs
- register the type for use with `g_autoptr`
- provide three standard type cast functions

The object struct and class struct need to be declared manually.

OBJECT_DECLARE_SIMPLE_TYPE (InstanceType, MODULE_OBJ_NAME)

Parameters

InstanceType instance struct name

MODULE_OBJ_NAME the object name in uppercase with underscore separators

Description

This does the same as `OBJECT_DECLARE_TYPE()`, but with no class struct declared.

This macro should be used unless the class struct needs to have virtual methods declared.

OBJECT_DEFINE_TYPE_EXTENDED (ModuleObjName, module_obj_name, MODULE_OBJ_NAME, PARENT_MODULE_OBJ_NAME, ABSTRACT, ...)

Parameters

ModuleObjName the object name with initial caps

module_obj_name the object name in lowercase with underscore separators

MODULE_OBJ_NAME the object name in uppercase with underscore separators

PARENT_MODULE_OBJ_NAME the parent object name in uppercase with underscore separators

ABSTRACT boolean flag to indicate whether the object can be instantiated

... list of initializers for “InterfaceInfo” to declare implemented interfaces

Description

This macro is typically used in a source file, and will:

- declare prototypes for `_finalize`, `_class_init` and `_init` methods
- declare the `TypeInfo` struct instance
- provide the constructor to register the type

After using this macro, implementations of the `_finalize`, `_class_init`, and `_init` methods need to be written. Any of these can be zero-line no-op impls if no special logic is required for a given type.

This macro should rarely be used, instead one of the more specialized macros is usually a better choice.

OBJECT_DEFINE_TYPE (ModuleObjName, module_obj_name, MODULE_OBJ_NAME, PARENT_MODULE_OBJ_NAME)

Parameters

ModuleObjName the object name with initial caps

module_obj_name the object name in lowercase with underscore separators

MODULE_OBJ_NAME the object name in uppercase with underscore separators

PARENT_MODULE_OBJ_NAME the parent object name in uppercase with underscore separators

Description

This is a specialization of `OBJECT_DEFINE_TYPE_EXTENDED`, which is suitable for the common case of a non-abstract type, without any interfaces.

OBJECT_DEFINE_TYPE_WITH_INTERFACES (ModuleObjName, module_obj_name, MODULE_OBJ_NAME, PARENT_MODULE_OBJ_NAME, ...)

Parameters

ModuleObjName the object name with initial caps

module_obj_name the object name in lowercase with underscore separators

MODULE_OBJ_NAME the object name in uppercase with underscore separators

PARENT_MODULE_OBJ_NAME the parent object name in uppercase with underscore separators

... list of initializers for “InterfaceInfo” to declare implemented interfaces

Description

This is a specialization of `OBJECT_DEFINE_TYPE_EXTENDED`, which is suitable for the common case of a non-abstract type, with one or more implemented interfaces.

Note when passing the list of interfaces, be sure to include the final NULL entry, e.g. { `TYPE_USER_CREATABLE` }, { NULL }

OBJECT_DEFINE_ABSTRACT_TYPE (ModuleObjName, module_obj_name, MODULE_OBJ_NAME, PARENT_MODULE_OBJ_NAME)

Parameters

ModuleObjName the object name with initial caps

module_obj_name the object name in lowercase with underscore separators

MODULE_OBJ_NAME the object name in uppercase with underscore separators

PARENT_MODULE_OBJ_NAME the parent object name in uppercase with underscore separators

Description

This is a specialization of `OBJECT_DEFINE_TYPE_EXTENDED`, which is suitable for defining an abstract type, without any interfaces.

struct **TypeInfo**

Definition

```
struct TypeInfo {
    const char *name;
    const char *parent;
    size_t instance_size;
    size_t instance_align;
    void (*instance_init)(Object *obj);
    void (*instance_post_init)(Object *obj);
    void (*instance_finalize)(Object *obj);
    bool abstract;
    size_t class_size;
    void (*class_init)(ObjectClass *klass, void *data);
    void (*class_base_init)(ObjectClass *klass, void *data);
    void *class_data;
    InterfaceInfo *interfaces;
};
```

Members

name The name of the type.

parent The name of the parent type.

instance_size The size of the object (derivative of *Object*). If **instance_size** is 0, then the size of the object will be the size of the parent object.

instance_align The required alignment of the object. If **instance_align** is 0, then normal malloc alignment is sufficient; if non-zero, then we must use `qemu_memalign` for allocation.

instance_init This function is called to initialize an object. The parent class will have already been initialized so the type is only responsible for initializing its own members.

instance_post_init This function is called to finish initialization of an object, after all **instance_init** functions were called.

instance_finalize This function is called during object destruction. This is called before the parent **instance_finalize** function has been called. An object should only free the members that are unique to its type in this function.

abstract If this field is true, then the class is considered abstract and cannot be directly instantiated.

class_size The size of the class object (derivative of *ObjectClass*) for this object. If **class_size** is 0, then the size of the class will be assumed to be the size of the parent class. This allows a type to avoid implementing an explicit class type if they are not adding additional virtual functions.

class_init This function is called after all parent class initialization has occurred to allow a class to set its default virtual method pointers. This is also the function to use to override virtual methods from a parent class.

class_base_init This function is called for all base classes after all parent class initialization has occurred, but before the class itself is initialized. This is the function to use to undo the effects of memcopy from the parent class to the descendants.

class_data Data to pass to the **class_init**, **class_base_init**. This can be useful when building dynamic classes.

interfaces The list of interfaces associated with this type. This should point to a static array that's terminated with a zero filled element.

OBJECT (obj)

Parameters

obj A derivative of *Object*

Description

Converts an object to a *Object*. Since all objects are *Objects*, this function will always succeed.

OBJECT_CLASS (class)

Parameters

class A derivative of *ObjectClass*.

Description

Converts a class to an *ObjectClass*. Since all objects are *Objects*, this function will always succeed.

OBJECT_CHECK (type, obj, name)

Parameters

type The C type to use for the return value.

obj A derivative of **type** to cast.

name The QOM typename of **type**

Description

A type safe version of **object_dynamic_cast_assert**. Typically each class will define a macro based on this type to perform type safe dynamic_casts to this object type.

If an invalid object is passed to this function, a run time assert will be generated.

OBJECT_CLASS_CHECK (class_type, class, name)

Parameters

class_type The C type to use for the return value.

class A derivative class of **class_type** to cast.

name the QOM typename of **class_type**.

Description

A type safe version of `object_class_dynamic_cast_assert`. This macro is typically wrapped by each type to perform type safe casts of a class to a specific class type.

OBJECT_GET_CLASS (class, obj, name)

Parameters

class The C type to use for the return value.

obj The object to obtain the class for.

name The QOM typename of **obj**.

Description

This function will return a specific class for a given object. Its generally used by each type to provide a type safe macro to get a specific class type from an object.

struct **InterfaceInfo**

Definition

```
struct InterfaceInfo {
    const char *type;
};
```

Members

type The name of the interface.

Description

The information associated with an interface.

struct **InterfaceClass**

Definition

```
struct InterfaceClass {
    ObjectClass parent_class;
};
```

Members

parent_class the base class

Description

The class for all interfaces. Subclasses of this class should only add virtual methods.

INTERFACE_CLASS (klass)

Parameters

klass class to cast from

Return

An *InterfaceClass* or raise an error if cast is invalid

INTERFACE_CHECK (interface, obj, name)

Parameters

interface the type to return

obj the object to convert to an interface

name the interface type name

Return

obj casted to **interface** if cast is valid, otherwise raise error.

Object * **object_new_with_class** (*ObjectClass* **klass*)

Parameters

ObjectClass ***klass** The class to instantiate.

Description

This function will initialize a new object using heap allocated memory. The returned object has a reference count of 1, and will be freed when the last reference is dropped.

Return

The newly allocated and instantiated object.

Object * **object_new** (const char **typename*)

Parameters

const char ***typename** The name of the type of the object to instantiate.

Description

This function will initialize a new object using heap allocated memory. The returned object has a reference count of 1, and will be freed when the last reference is dropped.

Return

The newly allocated and instantiated object.

Object * **object_new_with_props** (const char **typename*, *Object* **parent*, const char **id*, Error ***errp*, ...)

Parameters

const char ***typename** The name of the type of the object to instantiate.

Object ***parent** the parent object

const char ***id** The unique ID of the object

Error ****errp** pointer to error object

... list of property names and values

Description

This function will initialize a new object using heap allocated memory. The returned object has a reference count of 1, and will be freed when the last reference is dropped.

The **id** parameter will be used when registering the object as a child of **parent** in the composition tree.

The variadic parameters are a list of pairs of (propname, propvalue) strings. The propname of NULL indicates the end of the property list. If the object implements the user creatable interface, the object will be marked complete once all the properties have been processed.

Listing 12: Creating an object with properties

```

Error *err = NULL;
Object *obj;

obj = object_new_with_props(TYPE_MEMORY_BACKEND_FILE,
                           object_get_objects_root(),
                           "hostmem0",
                           &err,
                           "share", "yes",
                           "mem-path", "/dev/shm/somefile",
                           "prealloc", "yes",
                           "size", "1048576",
                           NULL);

if (!obj) {
    error_reportf_err(err, "Cannot create memory backend: ");
}

```

The returned object will have one stable reference maintained for as long as it is present in the object hierarchy.

Return

The newly allocated, instantiated & initialized object.

Object * **object_new_with_props** (const char **typename*, *Object* **parent*, const char **id*, Error ***errp*, va_list *vargs*)

Parameters

const char *typename The name of the type of the object to instantiate.

Object *parent the parent object

const char *id The unique ID of the object

Error **errp pointer to error object

va_list vargs list of property names and values

Description

See `object_new_with_props()` for documentation.

bool **object_set_props** (*Object* **obj*, Error ***errp*, ...)

Parameters

Object *obj the object instance to set properties on

Error **errp pointer to error object

... list of property names and values

Description

This function will set a list of properties on an existing object instance.

The variadic parameters are a list of pairs of (proppname, propvalue) strings. The proppname of NULL indicates the end of the property list.

Listing 13: Update an object's properties

```
Error *err = NULL;
Object *obj = ...get / create object...;

if (!object_set_props(obj,
                      &err,
                      "share", "yes",
                      "mem-path", "/dev/shm/somefile",
                      "prealloc", "yes",
                      "size", "1048576",
                      NULL)) {
    error_reportf_err(err, "Cannot set properties: ");
}
```

The returned object will have one stable reference maintained for as long as it is present in the object hierarchy.

Return

true on success, false on error.

bool **object_set_propv** (*Object* *obj, Error **errp, va_list args)

Parameters

Object *obj the object instance to set properties on

Error **errp pointer to error object

va_list args list of property names and values

Description

See object_set_props() for documentation.

Return

true on success, false on error.

void **object_initialize** (void *obj, size_t size, const char *typename)

Parameters

void *obj A pointer to the memory to be used for the object.

size_t size The maximum size available at obj for the object.

const char *typename The name of the type of the object to instantiate.

Description

This function will initialize an object. The memory for the object should have already been allocated. The returned object has a reference count of 1, and will be finalized when the last reference is dropped.

bool **object_initialize_child_with_props** (*Object* *parentobj, const char *propname, void *childobj, size_t size, const char *type, Error **errp, ...)

Parameters

Object *parentobj The parent object to add a property to

const char *propname The name of the property

void *childobj A pointer to the memory to be used for the object.

size_t size The maximum size available at childobj for the object.

const char *type The name of the type of the object to instantiate.

Error **errp If an error occurs, a pointer to an area to store the error

... list of property names and values

Description

This function will initialize an object. The memory for the object should have already been allocated. The object will then be added as child property to a parent with `object_property_add_child()` function. The returned object has a reference count of 1 (for the “child<...>” property from the parent), so the object will be finalized automatically when the parent gets removed.

The variadic parameters are a list of pairs of (propname, propvalue) strings. The propname of `NULL` indicates the end of the property list. If the object implements the user creatable interface, the object will be marked complete once all the properties have been processed.

Return

true on success, false on failure.

```
bool object_initialize_child_with_propsv (Object *parentobj, const char *propname,
                                         void *childobj, size_t size, const char *type,
                                         Error **errp, va_list vargs)
```

Parameters

Object *parentobj The parent object to add a property to

const char *propname The name of the property

void *childobj A pointer to the memory to be used for the object.

size_t size The maximum size available at **childobj** for the object.

const char *type The name of the type of the object to instantiate.

Error **errp If an error occurs, a pointer to an area to store the error

va_list vargs list of property names and values

Description

See `object_initialize_child()` for documentation.

Return

true on success, false on failure.

```
object_initialize_child (parent, propname, child, type)
```

Parameters

parent The parent object to add a property to

propname The name of the property

child A precisely typed pointer to the memory to be used for the object.

type The name of the type of the object to instantiate.

Description

This is like:

```
object_initialize_child_with_props (parent, propname,
                                   child, sizeof(*child), type,
                                   &error_abort, NULL)
```

Object * **object_dynamic_cast** (*Object* **obj*, const char **typename*)

Parameters

Object **obj* The object to cast.

const char **typename* The *typename* to cast to.

Description

This function will determine if **obj** is-a **typename**. **obj** can refer to an object or an interface associated with an object.

Return

This function returns **obj** on success or NULL on failure.

Object * **object_dynamic_cast_assert** (*Object* **obj*, const char **typename*, const char **file*, int *line*,
const char **func*)

Parameters

Object **obj* The object to cast.

const char **typename* The *typename* to cast to.

const char **file* Source code file where function was called

int *line* Source code line where function was called

const char **func* Name of function where this function was called

Description

See `object_dynamic_cast()` for a description of the parameters of this function. The only difference in behavior is that this function asserts instead of returning NULL on failure if QOM cast debugging is enabled. This function is not meant to be called directly, but only through the wrapper macro `OBJECT_CHECK`.

ObjectClass * **object_get_class** (*Object* **obj*)

Parameters

Object **obj* A derivative of *Object*

Return

The *ObjectClass* of the type associated with **obj**.

const char * **object_get_typename** (const *Object* **obj*)

Parameters

const Object **obj* A derivative of *Object*.

Return

The QOM typename of **obj**.

Type **type_register_static** (const *TypeInfo* **info*)

Parameters

const TypeInfo **info* The *TypeInfo* of the new type.

Description

info and all of the strings it points to should exist for the life time that the type is registered.

Return

the new Type.

Type **type_register** (const *TypeInfo* *info)

Parameters

const TypeInfo *info The *TypeInfo* of the new type

Description

Unlike `type_register_static()`, this call does not require **info** or its string members to continue to exist after the call returns.

Return

the new *Type*.

void **type_register_static_array** (const *TypeInfo* *infos, int nr_infos)

Parameters

const TypeInfo *infos The array of the new type *TypeInfo* structures.

int nr_infos number of entries in **infos**

Description

infos and all of the strings it points to should exist for the life time that the type is registered.

DEFINE_TYPES (type_array)

Parameters

type_array The array containing *TypeInfo* structures to register

Description

type_array should be static constant that exists for the life time that the type is registered.

ObjectClass * **object_class_dynamic_cast_assert** (*ObjectClass* *klass, const char *typename,
const char *file, int line, const char *func)

Parameters

ObjectClass ***klass** The *ObjectClass* to attempt to cast.

const char *typename The QOM typename of the class to cast to.

const char *file Source code file where function was called

int line Source code line where function was called

const char *func Name of function where this function was called

Description

See `object_class_dynamic_cast()` for a description of the parameters of this function. The only difference in behavior is that this function asserts instead of returning NULL on failure if QOM cast debugging is enabled. This function is not meant to be called directly, but only through the wrapper macro `OBJECT_CLASS_CHECK`.

ObjectClass * **object_class_dynamic_cast** (*ObjectClass* *klass, const char *typename)

Parameters

ObjectClass ***klass** The *ObjectClass* to attempt to cast.

const char *typename The QOM typename of the class to cast to.

Return

If **typename** is a class, this function returns **klass** if **typename** is a subtype of **klass**, else returns NULL.

Description

If **typename** is an interface, this function returns the interface definition for **klass** if **klass** implements it unambiguously; `NULL` is returned if **klass** does not implement the interface or if multiple classes or interfaces on the hierarchy leading to **klass** implement it. (FIXME: perhaps this can be detected at type definition time?)

ObjectClass * **object_class_get_parent** (*ObjectClass* **klass*)

Parameters

ObjectClass ***klass** The class to obtain the parent for.

Return

The parent for **klass** or `NULL` if none.

const char * **object_class_get_name** (*ObjectClass* **klass*)

Parameters

ObjectClass ***klass** The class to obtain the QOM typename for.

Return

The QOM typename for **klass**.

bool **object_class_is_abstract** (*ObjectClass* **klass*)

Parameters

ObjectClass ***klass** The class to obtain the abstractness for.

Return

true if **klass** is abstract, false otherwise.

ObjectClass * **object_class_by_name** (const char **typename*)

Parameters

const char ***typename** The QOM typename to obtain the class for.

Return

The class for **typename** or `NULL` if not found.

ObjectClass * **module_object_class_by_name** (const char **typename*)

Parameters

const char ***typename** The QOM typename to obtain the class for.

Description

For objects which might be provided by a module. Behaves like `object_class_by_name`, but additionally tries to load the module needed in case the class is not available.

Return

The class for **typename** or `NULL` if not found.

GSList * **object_class_get_list** (const char **implements_type*, bool *include_abstract*)

Parameters

const char ***implements_type** The type to filter for, including its derivatives.

bool **include_abstract** Whether to include abstract classes.

Return

A singly-linked list of the classes in reverse hashtable order.

GSList * **object_class_get_list_sorted** (const char **implements_type*, bool *include_abstract*)

Parameters

const char *implements_type The type to filter for, including its derivatives.

bool include_abstract Whether to include abstract classes.

Return

A singly-linked list of the classes in alphabetical case-insensitive order.

Object * **object_ref** (void **obj*)

Parameters

void *obj the object

Description

Increase the reference count of a object. A object cannot be freed as long as its reference count is greater than zero.

Return

obj

void **object_unref** (void **obj*)

Parameters

void *obj the object

Description

Decrease the reference count of a object. A object cannot be freed as long as its reference count is greater than zero.

ObjectProperty * **object_property_try_add** (*Object* **obj*, const char **name*, const char **type*, *ObjectPropertyAccessor* **get*, *ObjectPropertyAccessor* **set*, *ObjectPropertyRelease* **release*, void **opaque*, Error ***errp*)

Parameters

Object *obj the object to add a property to

const char *name the name of the property. This can contain any character except for a forward slash. In general, you should use hyphens '-' instead of underscores '_' when naming properties.

const char *type the type name of the property. This namespace is pretty loosely defined. Sub namespaces are constructed by using a prefix and then to angle brackets. For instance, the type 'virtio-net-pci' in the 'link' namespace would be 'link<virtio-net-pci>'.

ObjectPropertyAccessor *get The getter to be called to read a property. If this is NULL, then the property cannot be read.

ObjectPropertyAccessor *set the setter to be called to write a property. If this is NULL, then the property cannot be written.

ObjectPropertyRelease *release called when the property is removed from the object. This is meant to allow a property to free its opaque upon object destruction. This may be NULL.

void *opaque an opaque pointer to pass to the callbacks for the property

Error **errp pointer to error object

Return

The *ObjectProperty*; this can be used to set the **resolve** callback for child and link properties.

```
ObjectProperty * object_property_add(Object *obj, const char *name, const char *type, ObjectPropertyAccessor *get, ObjectPropertyAccessor *set, ObjectPropertyRelease *release, void *opaque)
```

Same as `object_property_try_add()` with **errp** hardcoded to `&error_abort`.

Parameters

Object *obj the object to add a property to

const char *name the name of the property. This can contain any character except for a forward slash. In general, you should use hyphens '-' instead of underscores '_' when naming properties.

const char *type the type name of the property. This namespace is pretty loosely defined. Sub namespaces are constructed by using a prefix and then to angle brackets. For instance, the type 'virtio-net-pci' in the 'link' namespace would be 'link<virtio-net-pci>'.

ObjectPropertyAccessor *get The getter to be called to read a property. If this is NULL, then the property cannot be read.

ObjectPropertyAccessor *set the setter to be called to write a property. If this is NULL, then the property cannot be written.

ObjectPropertyRelease *release called when the property is removed from the object. This is meant to allow a property to free its opaque upon object destruction. This may be NULL.

void *opaque an opaque pointer to pass to the callbacks for the property

`void object_property_set_default_bool (ObjectProperty *prop, bool value)`

Parameters

ObjectProperty *prop the property to set

bool value the value to be written to the property

Description

Set the property default value.

`void object_property_set_default_str (ObjectProperty *prop, const char *value)`

Parameters

ObjectProperty *prop the property to set

const char *value the value to be written to the property

Description

Set the property default value.

`void object_property_set_default_int (ObjectProperty *prop, int64_t value)`

Parameters

ObjectProperty *prop the property to set

int64_t value the value to be written to the property

Description

Set the property default value.

`void object_property_set_default_uint (ObjectProperty *prop, uint64_t value)`

Parameters

ObjectProperty *prop the property to set

uint64_t value the value to be written to the property

Description

Set the property default value.

ObjectProperty * **object_property_find**(*Object* *obj, const char *name)

Parameters

Object *obj the object

const char *name the name of the property

Description

Look up a property for an object.

Return its ObjectProperty if found, or NULL.

ObjectProperty * **object_property_find_err**(*Object* *obj, const char *name, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

Error **errp returns an error if this function fails

Description

Look up a property for an object.

Return its ObjectProperty if found, or NULL.

ObjectProperty * **object_class_property_find**(*ObjectClass* *klass, const char *name)

Parameters

ObjectClass *klass the object class

const char *name the name of the property

Description

Look up a property for an object class.

Return its ObjectProperty if found, or NULL.

ObjectProperty * **object_class_property_find_err**(*ObjectClass* *klass, const char *name, Error **errp)

Parameters

ObjectClass *klass the object class

const char *name the name of the property

Error **errp returns an error if this function fails

Description

Look up a property for an object class.

Return its ObjectProperty if found, or NULL.

void **object_property_iter_init**(ObjectPropertyIterator *iter, *Object* *obj)

Parameters

ObjectPropertyIterator *iter the iterator instance

Object *obj the object

Description

Initializes an iterator for traversing all properties registered against an object instance, its class and all parent classes.

It is forbidden to modify the property list while iterating, whether removing or adding properties.

Typical usage pattern would be

Listing 14: Using object property iterators

```
ObjectProperty *prop;
ObjectPropertyIterator iter;

object_property_iter_init(&iter, obj);
while ((prop = object_property_iter_next(&iter))) {
    ... do something with prop ...
}
```

void **object_class_property_iter_init** (ObjectPropertyIterator *iter, *ObjectClass* *klass)

Parameters

ObjectPropertyIterator *iter the iterator instance

ObjectClass *klass the class

Description

Initializes an iterator for traversing all properties registered against an object class and all parent classes.

It is forbidden to modify the property list while iterating, whether removing or adding properties.

This can be used on abstract classes as it does not create a temporary instance.

ObjectProperty * **object_property_iter_next** (ObjectPropertyIterator *iter)

Parameters

ObjectPropertyIterator *iter the iterator instance

Description

Return the next available property. If no further properties are available, a NULL value will be returned and the **iter** pointer should not be used again after this point without re-initializing it.

Return

the next property, or NULL when all properties have been traversed.

bool **object_property_get** (*Object* *obj, const char *name, Visitor *v, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

Visitor *v the visitor that will receive the property value. This should be an Output visitor and the data will be written with **name** as the name.

Error **errp returns an error if this function fails

Description

Reads a property from a object.

Return

true on success, false on failure.

bool **object_property_set_str** (*Object* *obj, const char *name, const char *value, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

const char *value the value to be written to the property

Error **errp returns an error if this function fails

Description

Writes a string value to a property.

Return

true on success, false on failure.

char * **object_property_get_str** (*Object* *obj, const char *name, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

Error **errp returns an error if this function fails

Return

the value of the property, converted to a C string, or NULL if an error occurs (including when the property value is not a string). The caller should free the string.

bool **object_property_set_link** (*Object* *obj, const char *name, *Object* *value, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

Object *value the value to be written to the property

Error **errp returns an error if this function fails

Description

Writes an object's canonical path to a property.

If the link property was created with OBJ_PROP_LINK_STRONG bit, the old target object is unreferenced, and a reference is added to the new target object.

Return

true on success, false on failure.

Object * **object_property_get_link** (*Object* *obj, const char *name, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

Error **errp returns an error if this function fails

Return

the value of the property, resolved from a path to an Object, or NULL if an error occurs (including when the property value is not a string or not a valid object path).

`bool object_property_set_bool (Object *obj, const char *name, bool value, Error **errp)`

Parameters

Object *obj the object

const char *name the name of the property

bool value the value to be written to the property

Error **errp returns an error if this function fails

Description

Writes a bool value to a property.

Return

true on success, false on failure.

`bool object_property_get_bool (Object *obj, const char *name, Error **errp)`

Parameters

Object *obj the object

const char *name the name of the property

Error **errp returns an error if this function fails

Return

the value of the property, converted to a boolean, or false if an error occurs (including when the property value is not a bool).

`bool object_property_set_int (Object *obj, const char *name, int64_t value, Error **errp)`

Parameters

Object *obj the object

const char *name the name of the property

int64_t value the value to be written to the property

Error **errp returns an error if this function fails

Description

Writes an integer value to a property.

Return

true on success, false on failure.

`int64_t object_property_get_int (Object *obj, const char *name, Error **errp)`

Parameters

Object *obj the object

const char *name the name of the property

Error **errp returns an error if this function fails

Return

the value of the property, converted to an integer, or -1 if an error occurs (including when the property value is not an integer).

```
bool object_property_set_uint (Object *obj, const char *name, uint64_t value, Error **errp)
```

Parameters

Object *obj the object

const char *name the name of the property

uint64_t value the value to be written to the property

Error **errp returns an error if this function fails

Description

Writes an unsigned integer value to a property.

Return

true on success, false on failure.

```
uint64_t object_property_get_uint (Object *obj, const char *name, Error **errp)
```

Parameters

Object *obj the object

const char *name the name of the property

Error **errp returns an error if this function fails

Return

the value of the property, converted to an unsigned integer, or 0 an error occurs (including when the property value is not an integer).

```
int object_property_get_enum (Object *obj, const char *name, const char *typename, Error **errp)
```

Parameters

Object *obj the object

const char *name the name of the property

const char *typename the name of the enum data type

Error **errp returns an error if this function fails

Return

the value of the property, converted to an integer (which can't be negative), or -1 on error (including when the property value is not an enum).

```
bool object_property_set (Object *obj, const char *name, Visitor *v, Error **errp)
```

Parameters

Object *obj the object

const char *name the name of the property

Visitor *v the visitor that will be used to write the property value. This should be an Input visitor and the data will be first read with **name** as the name and then written as the property value.

Error **errp returns an error if this function fails

Description

Writes a property to a object.

Return

true on success, false on failure.

bool **object_property_parse** (*Object* *obj, const char *name, const char *string, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

const char *string the string that will be used to parse the property value.

Error **errp returns an error if this function fails

Description

Parses a string and writes the result into a property of an object.

Return

true on success, false on failure.

char * **object_property_print** (*Object* *obj, const char *name, bool human, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

bool human if true, print for human consumption

Error **errp returns an error if this function fails

Description

Returns a string representation of the value of the property. The caller shall free the string.

const char * **object_property_get_type** (*Object* *obj, const char *name, Error **errp)

Parameters

Object *obj the object

const char *name the name of the property

Error **errp returns an error if this function fails

Return

The type name of the property.

Object * **object_get_root** (void)

Parameters

void no arguments

Return

the root object of the composition tree

Object * **object_get_objects_root** (void)

Parameters

void no arguments

Description

Get the container object that holds user created object instances. This is the object at path “/objects”

Return

the user object container

Object * **object_get_internal_root** (void)

Parameters

void no arguments

Description

Get the container object that holds internally used object instances. Any object which is put into this container must not be user visible, and it will not be exposed in the QOM tree.

Return

the internal object container

const char * **object_get_canonical_path_component** (const *Object* *obj)

Parameters

const *Object* *obj the object

Return

The final component in the object’s canonical path. The canonical path is the path within the composition tree starting from the root. NULL if the object doesn’t have a parent (and thus a canonical path).

char * **object_get_canonical_path** (const *Object* *obj)

Parameters

const *Object* *obj the object

Return

The canonical path for a object, newly allocated. This is the path within the composition tree starting from the root. Use g_free() to free it.

Object * **object_resolve_path** (const char *path, bool *ambiguous)

Parameters

const char *path the path to resolve

bool *ambiguous returns true if the path resolution failed because of an ambiguous match

Description

There are two types of supported paths—absolute paths and partial paths.

Absolute paths are derived from the root object and can follow child<> or link<> properties. Since they can follow link<> properties, they can be arbitrarily long. Absolute paths look like absolute filenames and are prefixed with a leading slash.

Partial paths look like relative filenames. They do not begin with a prefix. The matching rules for partial paths are subtle but designed to make specifying objects easy. At each level of the composition tree, the partial path is matched as an absolute path. The first match is not returned. At least two matches are searched for. A successful result is only returned if only one match is found. If more than one match is found, a flag is returned to indicate that the match was ambiguous.

Return

The matched object or NULL on path lookup failure.

Object * **object_resolve_path_type** (const char **path*, const char **typename*, bool **ambiguous*)

Parameters

const char *path the path to resolve

const char *typename the type to look for.

bool *ambiguous returns true if the path resolution failed because of an ambiguous match

Description

This is similar to `object_resolve_path`. However, when looking for a partial path only matches that implement the given type are considered. This restricts the search and avoids spuriously flagging matches as ambiguous.

For both partial and absolute paths, the return value goes through a dynamic cast to **typename**. This is important if either the link, or the typename itself are of interface types.

Return

The matched object or NULL on path lookup failure.

Object * **object_resolve_path_component** (*Object* **parent*, const char **part*)

Parameters

Object *parent the object in which to resolve the path

const char *part the component to resolve.

Description

This is similar to `object_resolve_path` with an absolute path, but it only resolves one element (**part**) and takes the others from **parent**.

Return

The resolved object or NULL on path lookup failure.

ObjectProperty * **object_property_try_add_child** (*Object* **obj*, const char **name*, *Object* **child*,
Error ***errp*)

Parameters

Object *obj the object to add a property to

const char *name the name of the property

Object *child the child object

Error **errp pointer to error object

Description

Child properties form the composition tree. All objects need to be a child of another object. Objects can only be a child of one object.

There is no way for a child to determine what its parent is. It is not a bidirectional relationship. This is by design.

The value of a child property as a C string will be the child object's canonical path. It can be retrieved using `object_property_get_str()`. The child object itself can be retrieved using `object_property_get_link()`.

Return

The newly added property on success, or NULL on failure.

ObjectProperty * **object_property_add_child**(*Object* *obj, const char *name, *Object* *child)

Parameters

Object *obj the object to add a property to
const char *name the name of the property
Object *child the child object

Description

Same as object_property_try_add_child() with **errp** hardcoded to &error_abort

void **object_property_allow_set_link**(const *Object* *obj, const char *name, *Object* *child, Error **errp)

Parameters

const Object *obj the object to add a property to
const char *name the name of the property
Object *child the child object
Error **errp pointer to error object

Description

The default implementation of the object_property_add_link() check() callback function. It allows the link property to be set and never returns an error.

ObjectProperty * **object_property_add_link**(*Object* *obj, const char *name, const char *type, *Object* **targetp, void (*check)(const *Object* *obj, const char *name, *Object* *val, Error **errp), ObjectPropertyLinkFlags flags)

Parameters

Object *obj the object to add a property to
const char *name the name of the property
const char *type the qobj type of the link
Object **targetp a pointer to where the link object reference is stored
void (*check)(const Object *obj, const char *name, Object *val, Error **errp)
 callback to veto setting or NULL if the property is read-only
ObjectPropertyLinkFlags flags additional options for the link

Description

Links establish relationships between objects. Links are unidirectional although two links can be combined to form a bidirectional relationship between objects.

Links form the graph in the object model.

The **check()** callback is invoked when object_property_set_link() is called and can raise an error to prevent the link being set. If **check** is NULL, the property is read-only and cannot be set.

Ownership of the pointer that **child** points to is transferred to the link property. The reference count for ***child** is managed by the property from after the function returns till the property is deleted with object_property_del(). If the **flags** OBJ_PROP_LINK_STRONG bit is set, the reference count is decremented when the property is deleted or modified.

Return

The newly added property on success, or NULL on failure.

ObjectProperty * **object_property_add_str** (*Object* *obj, const char *name, char *(*get)(*Object* *, Error **), void (*set)(*Object* *, const char *, Error **))

Parameters

Object *obj the object to add a property to

const char *name the name of the property

char *(*get) (**Object** *, **Error** **) the getter or NULL if the property is write-only. This function must return a string to be freed by g_free().

void (*set) (**Object** *, **const char** *, **Error** **) the setter or NULL if the property is read-only

Description

Add a string property using getters/setters. This function will add a property of type 'string'.

Return

The newly added property on success, or NULL on failure.

ObjectProperty * **object_property_add_bool** (*Object* *obj, const char *name, bool (*get)(*Object* *, Error **), void (*set)(*Object* *, bool, Error **))

Parameters

Object *obj the object to add a property to

const char *name the name of the property

bool (*get) (**Object** *, **Error** **) the getter or NULL if the property is write-only.

void (*set) (**Object** *, **bool**, **Error** **) the setter or NULL if the property is read-only

Description

Add a bool property using getters/setters. This function will add a property of type 'bool'.

Return

The newly added property on success, or NULL on failure.

ObjectProperty * **object_property_add_enum** (*Object* *obj, const char *name, const char *typename, const QEnumLookup *lookup, int (*get)(*Object* *, Error **), void (*set)(*Object* *, int, Error **))

Parameters

Object *obj the object to add a property to

const char *name the name of the property

const char *typename the name of the enum data type

const QEnumLookup *lookup enum value namelookup table

int (*get) (**Object** *, **Error** **) the getter or NULL if the property is write-only.

void (*set) (**Object** *, **int**, **Error** **) the setter or NULL if the property is read-only

Description

Add an enum property using getters/setters. This function will add a property of type 'typename'.

Return

The newly added property on success, or NULL on failure.


```
ObjectProperty * object_property_add_tm(Object *obj, const char *name, void (*get)(Object *, struct
tm *, Error **))
```

Parameters

Object *obj the object to add a property to

const char *name the name of the property

void (*get)(Object *, struct tm *, Error **) the getter or NULL if the property is write-only.

Description

Add a read-only struct tm valued property using a getter function. This function will add a property of type 'struct tm'.

Return

The newly added property on success, or NULL on failure.

```
ObjectProperty * object_property_add_uint8_ptr(Object *obj, const char *name, const uint8_t *v,
ObjectPropertyFlags flags)
```

Parameters

Object *obj the object to add a property to

const char *name the name of the property

const uint8_t *v pointer to value

ObjectPropertyFlags flags bitwise-or'd ObjectPropertyFlags

Description

Add an integer property in memory. This function will add a property of type 'uint8'.

Return

The newly added property on success, or NULL on failure.

```
ObjectProperty * object_property_add_uint16_ptr(Object *obj, const char *name, const
uint16_t *v, ObjectPropertyFlags flags)
```

Parameters

Object *obj the object to add a property to

const char *name the name of the property

const uint16_t *v pointer to value

ObjectPropertyFlags flags bitwise-or'd ObjectPropertyFlags

Description

Add an integer property in memory. This function will add a property of type 'uint16'.

Return

The newly added property on success, or NULL on failure.

```
ObjectProperty * object_property_add_uint32_ptr(Object *obj, const char *name, const
uint32_t *v, ObjectPropertyFlags flags)
```

Parameters

Object *obj the object to add a property to

const char *name the name of the property

const uint32_t *v pointer to value

ObjectPropertyFlags flags bitwise-or'd ObjectPropertyFlags

Description

Add an integer property in memory. This function will add a property of type 'uint32'.

Return

The newly added property on success, or NULL on failure.

ObjectProperty * **object_property_add_uint64_ptr**(*Object* *obj, const char *name, const uint64_t *v, ObjectPropertyFlags flags)

Parameters

Object *obj the object to add a property to

const char *name the name of the property

const uint64_t *v pointer to value

ObjectPropertyFlags flags bitwise-or'd ObjectPropertyFlags

Description

Add an integer property in memory. This function will add a property of type 'uint64'.

Return

The newly added property on success, or NULL on failure.

ObjectProperty * **object_property_add_alias**(*Object* *obj, const char *name, *Object* *target_obj, const char *target_name)

Parameters

Object *obj the object to add a property to

const char *name the name of the property

Object *target_obj the object to forward property access to

const char *target_name the name of the property on the forwarded object

Description

Add an alias for a property on an object. This function will add a property of the same type as the forwarded property.

The caller must ensure that **target_obj** stays alive as long as this property exists. In the case of a child object or an alias on the same object this will be the case. For aliases to other objects the caller is responsible for taking a reference.

Return

The newly added property on success, or NULL on failure.

ObjectProperty * **object_property_add_const_link**(*Object* *obj, const char *name, *Object* *target)

Parameters

Object *obj the object to add a property to

const char *name the name of the property

Object *target the object to be referred by the link

Description

Add an unmodifiable link for a property on an object. This function will add a property of type `link<TYPE>` where `TYPE` is the type of **target**.

The caller must ensure that **target** stays alive as long as this property exists. In the case **target** is a child of **obj**, this will be the case. Otherwise, the caller is responsible for taking a reference.

Return

The newly added property on success, or `NULL` on failure.

```
void object_property_set_description (Object *obj, const char *name, const char *description)
```

Parameters

Object *obj the object owning the property

const char *name the name of the property

const char *description the description of the property on the object

Description

Set an object property's description.

Return

`true` on success, `false` on failure.

```
int object_child_foreach (Object *obj, int (*fn)(Object *child, void *opaque), void *opaque)
```

Parameters

Object *obj the object whose children will be navigated

int (*fn) (**Object** *child, **void** *opaque) the iterator function to be called

void *opaque an opaque value that will be passed to the iterator

Description

Call **fn** passing each child of **obj** and **opaque** to it, until **fn** returns non-zero.

It is forbidden to add or remove children from **obj** from the **fn** callback.

Return

The last value returned by **fn**, or 0 if there is no child.

```
int object_child_foreach_recursive (Object *obj, int (*fn)(Object *child, void *opaque),
                                   void *opaque)
```

Parameters

Object *obj the object whose children will be navigated

int (*fn) (**Object** *child, **void** *opaque) the iterator function to be called

void *opaque an opaque value that will be passed to the iterator

Description

Call **fn** passing each child of **obj** and **opaque** to it, until **fn** returns non-zero. Calls recursively, all child nodes of **obj** will also be passed all the way down to the leaf nodes of the tree. Depth first ordering.

It is forbidden to add or remove children from **obj** (or its child nodes) from the **fn** callback.

Return

The last value returned by **fn**, or 0 if there is no child.

Object * **container_get** (*Object* **root*, const char **path*)

Parameters

Object **root* root of the #path, e.g., object_get_root()

const char **path* path to the container

Description

Return a container object whose path is **path**. Create more containers along the path if necessary.

Return

the container object.

size_t **object_type_get_instance_size** (const char **typename*)

Parameters

const char **typename* Name of the Type whose instance_size is required

Description

Returns the instance_size of the given **typename**.

char * **object_property_help** (const char **name*, const char **type*, QObject **defval*, const char **description*)

Parameters

const char **name* the name of the property

const char **type* the type of the property

QObject **defval* the default value

const char **description* description of the property

Return

a user-friendly formatted string describing the property for help purposes.

6.27 block-coroutine-wrapper

A lot of functions in QEMU block layer (see `block/*`) can only be called in coroutine context. Such functions are normally marked by the `coroutine_fn` specifier. Still, sometimes we need to call them from non-coroutine context; for this we need to start a coroutine, run the needed function from it and wait for the coroutine to finish in a `BDRV_POLL_WHILE()` loop. To run a coroutine we need a function with one `void*` argument. So for each `coroutine_fn` function which needs a non-coroutine interface, we should define a structure to pack the parameters, define a separate function to unpack the parameters and call the original function and finally define a new interface function with same list of arguments as original one, which will pack the parameters into a struct, create a coroutine, run it and wait in `BDRV_POLL_WHILE()` loop. It's boring to create such wrappers by hand, so we have a script to generate them.

6.27.1 Usage

Assume we have defined the `coroutine_fn` function `bdrv_co_foo(<some args>)` and need a non-coroutine interface for it, called `bdrv_foo(<same args>)`. In this case the script can help. To trigger the generation:

1. You need `bdrv_foo` declaration somewhere (for example, in `block/coroutines.h`) with the `generated_co_wrapper` mark, like this:

```
int generated_co_wrapper bdrv_foo(<some args>);
```

2. You need to feed this declaration to block-coroutine-wrapper script. For this, add the .h (or .c) file with the declaration to the input: `files(...)` list of `block_gen_c` target declaration in `block/meson.build`

You are done. During the build, coroutine wrappers will be generated in `<BUILD_DIR>/block/block-gen.c`.

6.27.2 Links

1. The script location is `scripts/block-coroutine-wrapper.py`.
2. Generic place for private `generated_co_wrapper` declarations is `block/coroutines.h`, for public declarations: `include/block/block.h`
3. The core API of generated coroutine wrappers is placed in (not generated) `block/block-gen.h`

This is the design document for multi-process QEMU. It does not necessarily reflect the status of the current implementation, which may lack features or be considerably different from what is described in this document. This document is still useful as a description of the goals and general direction of this feature.

Please refer to the following wiki for latest details: <https://wiki.qemu.org/Features/MultiProcessQEMU>

6.28 Multi-process QEMU

QEMU is often used as the hypervisor for virtual machines running in the Oracle cloud. Since one of the advantages of cloud computing is the ability to run many VMs from different tenants in the same cloud infrastructure, a guest that compromised its hypervisor could potentially use the hypervisor's access privileges to access data it is not authorized for.

QEMU can be susceptible to security attacks because it is a large, monolithic program that provides many features to the VMs it services. Many of these features can be configured out of QEMU, but even a reduced configuration QEMU has a large amount of code a guest can potentially attack. Separating QEMU reduces the attack surface by aiding to limit each component in the system to only access the resources that it needs to perform its job.

6.28.1 QEMU services

QEMU can be broadly described as providing three main services. One is a VM control point, where VMs can be created, migrated, re-configured, and destroyed. A second is to emulate the CPU instructions within the VM, often accelerated by HW virtualization features such as Intel's VT extensions. Finally, it provides IO services to the VM by emulating HW IO devices, such as disk and network devices.

A multi-process QEMU

A multi-process QEMU involves separating QEMU services into separate host processes. Each of these processes can be given only the privileges it needs to provide its service, e.g., a disk service could be given access only to the disk images it provides, and not be allowed to access other files, or any network devices. An attacker who compromised this service would not be able to use this exploit to access files or devices beyond what the disk service was given access to.

A QEMU control process would remain, but in multi-process mode, will have no direct interfaces to the VM. During VM execution, it would still provide the user interface to hot-plug devices or live migrate the VM.

A first step in creating a multi-process QEMU is to separate IO services from the main QEMU program, which would continue to provide CPU emulation. i.e., the control process would also be the CPU emulation process. In a later phase, CPU emulation could be separated from the control process.

6.28.2 Separating IO services

Separating IO services into individual host processes is a good place to begin for a couple of reasons. One is the sheer number of IO devices QEMU can emulate provides a large surface of interfaces which could potentially be exploited, and, indeed, have been a source of exploits in the past. Another is the modular nature of QEMU device emulation code provides interface points where the QEMU functions that perform device emulation can be separated from the QEMU functions that manage the emulation of guest CPU instructions. The devices emulated in the separate process are referred to as remote devices.

QEMU device emulation

QEMU uses an object oriented SW architecture for device emulation code. Configured objects are all compiled into the QEMU binary, then objects are instantiated by name when used by the guest VM. For example, the code to emulate a device named “foo” is always present in QEMU, but its instantiation code is only run when the device is included in the target VM. (e.g., via the QEMU command line as *-device foo*)

The object model is hierarchical, so device emulation code names its parent object (such as “pci-device” for a PCI device) and QEMU will instantiate a parent object before calling the device’s instantiation code.

Current separation models

In order to separate the device emulation code from the CPU emulation code, the device object code must run in a different process. There are a couple of existing QEMU features that can run emulation code separately from the main QEMU process. These are examined below.

vhost user model

Virtio guest device drivers can be connected to vhost user applications in order to perform their IO operations. This model uses special virtio device drivers in the guest and vhost user device objects in QEMU, but once the QEMU vhost user code has configured the vhost user application, mission-mode IO is performed by the application. The vhost user application is a daemon process that can be contacted via a known UNIX domain socket.

vhost socket

As mentioned above, one of the tasks of the vhost device object within QEMU is to contact the vhost application and send it configuration information about this device instance. As part of the configuration process, the application can also be sent other file descriptors over the socket, which then can be used by the vhost user application in various ways, some of which are described below.

vhost MMIO store acceleration

VMs are often run using HW virtualization features via the KVM kernel driver. This driver allows QEMU to accelerate the emulation of guest CPU instructions by running the guest in a virtual HW mode. When the guest executes instructions that cannot be executed by virtual HW mode, execution returns to the KVM driver so it can inform QEMU to emulate the instructions in SW.

One of the events that can cause a return to QEMU is when a guest device driver accesses an IO location. QEMU then dispatches the memory operation to the corresponding QEMU device object. In the case of a vhost user device, the memory operation would need to be sent over a socket to the vhost application. This path is accelerated by the QEMU virtio code by setting up an eventfd file descriptor that the vhost application can directly receive MMIO store notifications from the KVM driver, instead of needing them to be sent to the QEMU process first.

vhost interrupt acceleration

Another optimization used by the vhost application is the ability to directly inject interrupts into the VM via the KVM driver, again, bypassing the need to send the interrupt back to the QEMU process first. The QEMU virtio setup code configures the KVM driver with an eventfd that triggers the device interrupt in the guest when the eventfd is written. This irqfd file descriptor is then passed to the vhost user application program.

vhost access to guest memory

The vhost application is also allowed to directly access guest memory, instead of needing to send the data as messages to QEMU. This is also done with file descriptors sent to the vhost user application by QEMU. These descriptors can be passed to `mmap()` by the vhost application to map the guest address space into the vhost application.

IOMMUs introduce another level of complexity, since the address given to the guest virtio device to DMA to or from is not a guest physical address. This case is handled by having vhost code within QEMU register as a listener for IOMMU mapping changes. The vhost application maintains a cache of IOMMU translations: sending translation requests back to QEMU on cache misses, and in turn receiving flush requests from QEMU when mappings are purged.

applicability to device separation

Much of the vhost model can be re-used by separated device emulation. In particular, the ideas of using a socket between QEMU and the device emulation application, using a file descriptor to inject interrupts into the VM via KVM, and allowing the application to `mmap()` the guest should be re used.

There are, however, some notable differences between how a vhost application works and the needs of separated device emulation. The most basic is that vhost uses custom virtio device drivers which always trigger IO with MMIO stores. A separated device emulation model must work with existing IO device models and guest device drivers. MMIO loads break vhost store acceleration since they are synchronous - guest progress cannot continue until the load has been emulated. By contrast, stores are asynchronous, the guest can continue after the store event has been sent to the vhost application.

Another difference is that in the vhost user model, a single daemon can support multiple QEMU instances. This is contrary to the security regime desired, in which the emulation application should only be allowed to access the files or devices the VM it's running on behalf of can access. ##### qemu-io model

Qemu-io is a test harness used to test changes to the QEMU block backend object code. (e.g., the code that implements disk images for disk driver emulation) Qemu-io is not a device emulation application per se, but it does compile the QEMU block objects into a separate binary from the main QEMU one. This could be useful for disk device emulation, since its emulation applications will need to include the QEMU block objects.

6.28.3 New separation model based on proxy objects

A different model based on proxy objects in the QEMU program communicating with remote emulation programs could provide separation while minimizing the changes needed to the device emulation code. The rest of this section is a discussion of how a proxy object model would work.

Remote emulation processes

The remote emulation process will run the QEMU object hierarchy without modification. The device emulation objects will be also be based on the QEMU code, because for anything but the simplest device, it would not be a tractable to re-implement both the object model and the many device backends that QEMU has.

The processes will communicate with the QEMU process over UNIX domain sockets. The processes can be executed either as standalone processes, or be executed by QEMU. In both cases, the host backends the emulation processes will provide are specified on its command line, as they would be for QEMU. For example:

```
disk-proc -blockdev driver=file,node-name=file0,filename=disk-file0 \
-blockdev driver=qcow2,node-name=drive0,file=file0
```

would indicate process *disk-proc* uses a qcow2 emulated disk named *file0* as its backend.

Emulation processes may emulate more than one guest controller. A common configuration might be to put all controllers of the same device class (e.g., disk, network, etc.) in a single process, so that all backends of the same type can be managed by a single QMP monitor.

communication with QEMU

The first argument to the remote emulation process will be a Unix domain socket that connects with the Proxy object. This is a required argument.

```
disk-proc <socket number> <backend list>
```

remote process QMP monitor

Remote emulation processes can be monitored via QMP, similar to QEMU itself. The QMP monitor socket is specified the same as for a QEMU process:

```
disk-proc -qmp unix:/tmp/disk-mon,server
```

can be monitored over the UNIX socket path */tmp/disk-mon*.

QEMU command line

Each remote device emulated in a remote process on the host is represented as a *-device* of type *pci-proxy-dev*. A socket sub-option to this option specifies the Unix socket that connects to the remote process. An *id* sub-option is required, and it should be the same id as used in the remote process.

```
qemu-system-x86_64 ... -device pci-proxy-dev,id=lsi0,socket=3
```

can be used to add a device emulated in a remote process

QEMU management of remote processes

QEMU is not aware of the type of type of the remote PCI device. It is a pass through device as far as QEMU is concerned.

communication with emulation process

primary channel

The primary channel (referred to as `com` in the code) is used to bootstrap the remote process. It is also used to pass on device-agnostic commands like `reset`.

per-device channels

Each remote device communicates with QEMU using a dedicated communication channel. The proxy object sets up this channel using the primary channel during its initialization.

QEMU device proxy objects

QEMU has an object model based on sub-classes inherited from the “object” super-class. The sub-classes that are of interest here are the “device” and “bus” sub-classes whose child sub-classes make up the device tree of a QEMU emulated system.

The proxy object model will use device proxy objects to replace the device emulation code within the QEMU process. These objects will live in the same place in the object and bus hierarchies as the objects they replace. i.e., the proxy object for an LSI SCSI controller will be a sub-class of the “pci-device” class, and will have the same PCI bus parent and the same SCSI bus child objects as the LSI controller object it replaces.

It is worth noting that the same proxy object is used to mediate with all types of remote PCI devices.

object initialization

The Proxy device objects are initialized in the exact same manner in which any other QEMU device would be initialized.

In addition, the Proxy objects perform the following two tasks: - Parses the “socket” sub option and connects to the remote process using this channel - Uses the “id” sub-option to connect to the emulated device on the separate process

class_init

The `class_init()` method of a proxy object will, in general behave similarly to the object it replaces, including setting any static properties and methods needed by the proxy.

instance_init / realize

The `instance_init()` and `realize()` functions would only need to perform tasks related to being a proxy, such as registering its own MMIO handlers, or creating a child bus that other proxy devices can be attached to later.

Other tasks will be device-specific. For example, PCI device objects will initialize the PCI config space in order to make a valid PCI device tree within the QEMU process.

address space registration

Most devices are driven by guest device driver accesses to IO addresses or ports. The QEMU device emulation code uses QEMU's memory region function calls (such as `memory_region_init_io()`) to add callback functions that QEMU will invoke when the guest accesses the device's areas of the IO address space. When a guest driver does access the device, the VM will exit HW virtualization mode and return to QEMU, which will then lookup and execute the corresponding callback function.

A proxy object would need to mirror the memory region calls the actual device emulator would perform in its initialization code, but with its own callbacks. When invoked by QEMU as a result of a guest IO operation, they will forward the operation to the device emulation process.

PCI config space

PCI devices also have a configuration space that can be accessed by the guest driver. Guest accesses to this space is not handled by the device emulation object, but by its PCI parent object. Much of this space is read-only, but certain registers (especially BAR and MSI-related ones) need to be propagated to the emulation process.

PCI parent proxy

One way to propagate guest PCI config accesses is to create a "pci-device-proxy" class that can serve as the parent of a PCI device proxy object. This class's parent would be "pci-device" and it would override the PCI parent's `config_read()` and `config_write()` methods with ones that forward these operations to the emulation program.

interrupt receipt

A proxy for a device that generates interrupts will need to create a socket to receive interrupt indications from the emulation process. An incoming interrupt indication would then be sent up to its bus parent to be injected into the guest. For example, a PCI device object may use `pci_set_irq()`.

live migration

The proxy will register to save and restore any *vmstate* it needs over a live migration event. The device proxy does not need to manage the remote device's *vmstate*; that will be handled by the remote process proxy (see below).

QEMU remote device operation

Generic device operations, such as DMA, will be performed by the remote process proxy by sending messages to the remote process.

DMA operations

DMA operations would be handled much like vhost applications do. One of the initial messages sent to the emulation process is a guest memory table. Each entry in this table consists of a file descriptor and size that the emulation process can `mmap()` to directly access guest memory, similar to `vhost_user_set_mem_table()`. Note guest memory must be backed by file descriptors, such as when QEMU is given the *-mem-path* command line option.

IOMMU operations

When the emulated system includes an IOMMU, the remote process proxy in QEMU will need to create a socket for IOMMU requests from the emulation process. It will handle those requests with an `address_space_get_iotlb_entry()` call. In order to handle IOMMU unmaps, the remote process proxy will also register as a listener on the device's DMA address space. When an IOMMU memory region is created within the DMA address space, an IOMMU notifier for unmaps will be added to the memory region that will forward unmaps to the emulation process over the IOMMU socket.

device hot-plug via QMP

An QMP “`device_add`” command can add a device emulated by a remote process. It will also have “`rid`” option to the command, just as the `-device` command line option does. The remote process may either be one started at QEMU startup, or be one added by the “`add-process`” QMP command described above. In either case, the remote process proxy will forward the new device's JSON description to the corresponding emulation process.

live migration

The remote process proxy will also register for live migration notifications with `vmstate_register()`. When called to save state, the proxy will send the remote process a secondary socket file descriptor to save the remote process's device *vmstate* over. The incoming byte stream length and data will be saved as the proxy's *vmstate*. When the proxy is resumed on its new host, this *vmstate* will be extracted, and a secondary socket file descriptor will be sent to the new remote process through which it receives the *vmstate* in order to restore the devices there.

device emulation in remote process

The parts of QEMU that the emulation program will need include the object model; the memory emulation objects; the device emulation objects of the targeted device, and any dependent devices; and, the device's backends. It will also need code to setup the machine environment, handle requests from the QEMU process, and route machine-level requests (such as interrupts or IOMMU mappings) back to the QEMU process.

initialization

The process initialization sequence will follow the same sequence followed by QEMU. It will first initialize the back-end objects, then device emulation objects. The JSON descriptions sent by the QEMU process will drive which objects need to be created.

- address spaces

Before the device objects are created, the initial address spaces and memory regions must be configured with `memory_map_init()`. This creates a RAM memory region object (*system_memory*) and an IO memory region object (*system_io*).

- RAM

RAM memory region creation will follow how `pc_memory_init()` creates them, but must use `memory_region_init_ram_from_fd()` instead of `memory_region_allocate_system_memory()`. The file descriptors needed will be supplied by the guest memory table from above. Those RAM regions would then be added to the *system_memory* memory region with `memory_region_add_subregion()`.

- PCI

IO initialization will be driven by the JSON descriptions sent from the QEMU process. For a PCI device, a PCI bus will need to be created with `pci_root_bus_new()`, and a PCI memory region will need to be created and added to the `system_memory` memory region with `memory_region_add_subregion_overlap()`. The overlap version is required for architectures where PCI memory overlaps with RAM memory.

MMIO handling

The device emulation objects will use `memory_region_init_io()` to install their MMIO handlers, and `pci_register_bar()` to associate those handlers with a PCI BAR, as they do within QEMU currently.

In order to use `address_space_rw()` in the emulation process to handle MMIO requests from QEMU, the PCI physical addresses must be the same in the QEMU process and the device emulation process. In order to accomplish that, guest BAR programming must also be forwarded from QEMU to the emulation process.

interrupt injection

When device emulation wants to inject an interrupt into the VM, the request climbs the device's bus object hierarchy until the point where a bus object knows how to signal the interrupt to the guest. The details depend on the type of interrupt being raised.

- PCI pin interrupts

On x86 systems, there is an emulated IOAPIC object attached to the root PCI bus object, and the root PCI object forwards interrupt requests to it. The IOAPIC object, in turn, calls the KVM driver to inject the corresponding interrupt into the VM. The simplest way to handle this in an emulation process would be to setup the root PCI bus driver (via `pci_bus_irqs()`) to send a interrupt request back to the QEMU process, and have the device proxy object reflect it up the PCI tree there.

- PCI MSI/X interrupts

PCI MSI/X interrupts are implemented in HW as DMA writes to a CPU-specific PCI address. In QEMU on x86, a KVM APIC object receives these DMA writes, then calls into the KVM driver to inject the interrupt into the VM. A simple emulation process implementation would be to send the MSI DMA address from QEMU as a message at initialization, then install an address space handler at that address which forwards the MSI message back to QEMU.

DMA operations

When a emulation object wants to DMA into or out of guest memory, it first must use `dma_memory_map()` to convert the DMA address to a local virtual address. The emulation process memory region objects setup above will be used to translate the DMA address to a local virtual address the device emulation code can access.

IOMMU

When an IOMMU is in use in QEMU, DMA translation uses IOMMU memory regions to translate the DMA address to a guest physical address before that physical address can be translated to a local virtual address. The emulation process will need similar functionality.

- IOTLB cache

The emulation process will maintain a cache of recent IOMMU translations (the IOTLB). When the `translate()` callback of an IOMMU memory region is invoked, the IOTLB cache will be searched for an entry that will map the DMA address to a guest PA. On a cache miss, a message will be sent back to QEMU requesting the corresponding translation entry, which be both be used to return a guest address and be added to the cache.

- IOTLB purge

The IOMMU emulation will also need to act on unmap requests from QEMU. These happen when the guest IOMMU driver purges an entry from the guest's translation table.

live migration

When a remote process receives a live migration indication from QEMU, it will set up a channel using the received file descriptor with `qio_channel_socket_new_fd()`. This channel will be used to create a *QEMUfile* that can be passed to `qemu_save_device_state()` to send the process's device state back to QEMU. This method will be reversed on restore - the channel will be passed to `qemu_loadvm_state()` to restore the device state.

Accelerating device emulation

The messages that are required to be sent between QEMU and the emulation process can add considerable latency to IO operations. The optimizations described below attempt to ameliorate this effect by allowing the emulation process to communicate directly with the kernel KVM driver. The KVM file descriptors created would be passed to the emulation process via initialization messages, much like the guest memory table is done. ##### MMIO acceleration

Vhost user applications can receive guest virtio driver stores directly from KVM. The issue with the eventfd mechanism used by vhost user is that it does not pass any data with the event indication, so it cannot handle guest loads or guest stores that carry store data. This concept could, however, be expanded to cover more cases.

The expanded idea would require a new type of KVM device: *KVM_DEV_TYPE_USER*. This device has two file descriptors: a master descriptor that QEMU can use for configuration, and a slave descriptor that the emulation process can use to receive MMIO notifications. QEMU would create both descriptors using the KVM driver, and pass the slave descriptor to the emulation process via an initialization message.

data structures

- guest physical range

The guest physical range structure describes the address range that a device will respond to. It includes the base and length of the range, as well as which bus the range resides on (e.g., on an x86machine, it can specify whether the range refers to memory or IO addresses).

A device can have multiple physical address ranges it responds to (e.g., a PCI device can have multiple BARs), so the structure will also include an enumerated identifier to specify which of the device's ranges is being referred to.

Name	Description
addr	range base address
len	range length
bus	addr type (memory or IO)
id	range ID (e.g., PCI BAR)

- MMIO request structure

This structure describes an MMIO operation. It includes which guest physical range the MMIO was within, the offset within that range, the MMIO type (e.g., load or store), and its length and data. It also includes a sequence number that can be used to reply to the MMIO, and the CPU that issued the MMIO.

Name	Description
rid	range MMIO is within
offset	offset withing <i>rid</i>
type	e.g., load or store
len	MMIO length
data	store data
seq	sequence ID

- MMIO request queues

MMIO request queues are FIFO arrays of MMIO request structures. There are two queues: pending queue is for MMIOs that haven't been read by the emulation program, and the sent queue is for MMIOs that haven't been acknowledged. The main use of the second queue is to validate MMIO replies from the emulation program.

- scoreboard

Each CPU in the VM is emulated in QEMU by a separate thread, so multiple MMIOs may be waiting to be consumed by an emulation program and multiple threads may be waiting for MMIO replies. The scoreboard would contain a wait queue and sequence number for the per-CPU threads, allowing them to be individually woken when the MMIO reply is received from the emulation program. It also tracks the number of posted MMIO stores to the device that haven't been replied to, in order to satisfy the PCI constraint that a load to a device will not complete until all previous stores to that device have been completed.

- device shadow memory

Some MMIO loads do not have device side-effects. These MMIOs can be completed without sending a MMIO request to the emulation program if the emulation program shares a shadow image of the device's memory image with the KVM driver.

The emulation program will ask the KVM driver to allocate memory for the shadow image, and will then use `mmap()` to directly access it. The emulation program can control KVM access to the shadow image by sending KVM an access map telling it which areas of the image have no side-effects (and can be completed immediately), and which require a MMIO request to the emulation program. The access map can also inform the KVM drive which size accesses are allowed to the image.

master descriptor

The master descriptor is used by QEMU to configure the new KVM device. The descriptor would be returned by the KVM driver when QEMU issues a `KVM_CREATE_DEVICE ioctl()` with a `KVM_DEV_TYPE_USER` type.

KVM_DEV_TYPE_USER device ops

The `KVM_DEV_TYPE_USER` operations vector will be registered by a `kvm_register_device_ops()` call when the KVM system is initialized by `kvm_init()`. These device ops are called by the KVM driver when QEMU executes certain `ioctl()` operations on its KVM file descriptor. They include:

- create

This routine is called when QEMU issues a `KVM_CREATE_DEVICE ioctl()` on its per-VM file descriptor. It will allocate and initialize a KVM user device specific data structure, and assign the `kvm_device` private field to it.

- ioctl

This routine is invoked when QEMU issues an `ioctl()` on the master descriptor. The `ioctl()` commands supported are defined by the KVM device type. `KVM_DEV_TYPE_USER` ones will need several commands:

`KVM_DEV_USER_SLAVE_FD` creates the slave file descriptor that will be passed to the device emulation program. Only one slave can be created by each master descriptor. The file operations performed by this descriptor are described below.

The *KVM_DEV_USER_PA_RANGE* command configures a guest physical address range that the slave descriptor will receive MMIO notifications for. The range is specified by a guest physical range structure argument. For buses that assign addresses to devices dynamically, this command can be executed while the guest is running, such as the case when a guest changes a device's PCI BAR registers.

KVM_DEV_USER_PA_RANGE will use `kvm_io_bus_register_dev()` to register *kvm_io_device_ops* callbacks to be invoked when the guest performs a MMIO operation within the range. When a range is changed, `kvm_io_bus_unregister_dev()` is used to remove the previous instantiation.

KVM_DEV_USER_TIMEOUT will configure a timeout value that specifies how long KVM will wait for the emulation process to respond to a MMIO indication.

- destroy

This routine is called when the VM instance is destroyed. It will need to destroy the slave descriptor; and free any memory allocated by the driver, as well as the *kvm_device* structure itself.

slave descriptor

The slave descriptor will have its own file operations vector, which responds to system calls on the descriptor performed by the device emulation program.

- read

A read returns any pending MMIO requests from the KVM driver as MMIO request structures. Multiple structures can be returned if there are multiple MMIO operations pending. The MMIO requests are moved from the pending queue to the sent queue, and if there are threads waiting for space in the pending to add new MMIO operations, they will be woken here.

- write

A write also consists of a set of MMIO requests. They are compared to the MMIO requests in the sent queue. Matches are removed from the sent queue, and any threads waiting for the reply are woken. If a store is removed, then the number of posted stores in the per-CPU scoreboard is decremented. When the number is zero, and a non side-effect load was waiting for posted stores to complete, the load is continued.

- ioctl

There are several `ioctl()`s that can be performed on the slave descriptor.

A *KVM_DEV_USER_SHADOW_SIZE* `ioctl()` causes the KVM driver to allocate memory for the shadow image. This memory can later be `mmap()`ed by the emulation process to share the emulation's view of device memory with the KVM driver.

A *KVM_DEV_USER_SHADOW_CTRL* `ioctl()` controls access to the shadow image. It will send the KVM driver a shadow control map, which specifies which areas of the image can complete guest loads without sending the load request to the emulation program. It will also specify the size of load operations that are allowed.

- poll

An emulation program will use the `poll()` call with a *POLLIN* flag to determine if there are MMIO requests waiting to be read. It will return if the pending MMIO request queue is not empty.

- mmap

This call allows the emulation program to directly access the shadow image allocated by the KVM driver. As device emulation updates device memory, changes with no side-effects will be reflected in the shadow, and the KVM driver can satisfy guest loads from the shadow image without needing to wait for the emulation program.

kvm_io_device ops

Each KVM per-CPU thread can handle MMIO operation on behalf of the guest VM. KVM will use the MMIO's guest physical address to search for a matching *kvm_io_device* to see if the MMIO can be handled by the KVM driver instead of exiting back to QEMU. If a match is found, the corresponding callback will be invoked.

- read

This callback is invoked when the guest performs a load to the device. Loads with side-effects must be handled synchronously, with the KVM driver putting the QEMU thread to sleep waiting for the emulation process reply before re-starting the guest. Loads that do not have side-effects may be optimized by satisfying them from the shadow image, if there are no outstanding stores to the device by this CPU. PCI memory ordering demands that a load cannot complete before all older stores to the same device have been completed.

- write

Stores can be handled asynchronously unless the pending MMIO request queue is full. In this case, the QEMU thread must sleep waiting for space in the queue. Stores will increment the number of posted stores in the per-CPU scoreboard, in order to implement the PCI ordering constraint above.

interrupt acceleration

This performance optimization would work much like a vhost user application does, where the QEMU process sets up *eventfds* that cause the device's corresponding interrupt to be triggered by the KVM driver. These irq file descriptors are sent to the emulation process at initialization, and are used when the emulation code raises a device interrupt.

intx acceleration

Traditional PCI pin interrupts are level based, so, in addition to an irq file descriptor, a re-sampling file descriptor needs to be sent to the emulation program. This second file descriptor allows multiple devices sharing an irq to be notified when the interrupt has been acknowledged by the guest, so they can re-trigger the interrupt if their device has not de-asserted its interrupt.

intx irq descriptor

The irq descriptors are created by the proxy object using `event_notifier_init()` to create the irq and re-sampling *eventfds*, and `kvm_vm_ioctl(KVM_IRQFD)` to bind them to an interrupt. The interrupt route can be found with `pci_device_route_intx_to_irq()`.

intx routing changes

Intx routing can be changed when the guest programs the APIC the device pin is connected to. The proxy object in QEMU will use `pci_device_set_intx_routing_notifier()` to be informed of any guest changes to the route. This handler will broadly follow the VFIO interrupt logic to change the route: de-assigning the existing irq descriptor from its route, then assigning it the new route. (see `vfio_intx_update()`)

MSI/X acceleration

MSI/X interrupts are sent as DMA transactions to the host. The interrupt data contains a vector that is programmed by the guest. A device may have multiple MSI interrupts associated with it, so multiple irq descriptors may need to be sent to the emulation program.

MSI/X irq descriptor

This case will also follow the VFIO example. For each MSI/X interrupt, an *eventfd* is created, a virtual interrupt is allocated by `kvm_irqchip_add_msi_route()`, and the virtual interrupt is bound to the eventfd with `kvm_irqchip_add_irqfd_notifier()`.

MSI/X config space changes

The guest may dynamically update several MSI-related tables in the device's PCI config space. These include per-MSI interrupt enables and vector data. Additionally, MSIX tables exist in device memory space, not config space. Much like the BAR case above, the proxy object must look at guest config space programming to keep the MSI interrupt state consistent between QEMU and the emulation program.

6.28.4 Disaggregated CPU emulation

After IO services have been disaggregated, a second phase would be to separate a process to handle CPU instruction emulation from the main QEMU control function. There are no object separation points for this code, so the first task would be to create one.

6.28.5 Host access controls

Separating QEMU relies on the host OS's access restriction mechanisms to enforce that the differing processes can only access the objects they are entitled to. There are a couple types of mechanisms usually provided by general purpose OSs.

Discretionary access control

Discretionary access control allows each user to control who can access their files. In Linux, this type of control is usually too coarse for QEMU separation, since it only provides three separate access controls: one for the same user ID, the second for users IDs with the same group ID, and the third for all other user IDs. Each device instance would need a separate user ID to provide access control, which is likely to be unwieldy for dynamically created VMs.

Mandatory access control

Mandatory access control allows the OS to add an additional set of controls on top of discretionary access for the OS to control. It also adds other attributes to processes and files such as types, roles, and categories, and can establish rules for how processes and files can interact.

Type enforcement

Type enforcement assigns a *type* attribute to processes and files, and allows rules to be written on what operations a process with a given type can perform on a file with a given type. QEMU separation could take advantage of type enforcement by running the emulation processes with different types, both from the main QEMU process, and from the emulation processes of different classes of devices.

For example, guest disk images and disk emulation processes could have types separate from the main QEMU process and non-disk emulation processes, and the type rules could prevent processes other than disk emulation ones from accessing guest disk images. Similarly, network emulation processes can have a type separate from the main QEMU process and non-network emulation process, and only that type can access the host tun/tap device used to provide guest networking.

Category enforcement

Category enforcement assigns a set of numbers within a given range to the process or file. The process is granted access to the file if the process's set is a superset of the file's set. This enforcement can be used to separate multiple instances of devices in the same class.

For example, if there are multiple disk devices provided to a guest, each device emulation process could be provisioned with a separate category. The different device emulation processes would not be able to access each other's backing disk images.

Alternatively, categories could be used in lieu of the type enforcement scheme described above. In this scenario, different categories would be used to prevent device emulation processes in different classes from accessing resources assigned to other classes.

Symbols

- aio=AIO
 - qemu-nbd command line option, [194](#)
- backing-chain
 - qemu-img-common-opts command line option, [181](#)
- bitmaps
 - qemu-img-convert command line option, [181](#)
- blockdev BLOCKDEVDEF
 - qemu-storage-daemon command line option, [191](#)
- cache=CACHE
 - qemu-nbd command line option, [194](#)
- cache=none|auto|always
 - virtiofsd command line option, [201](#)
- chardev CHARDEVDEF
 - qemu-storage-daemon command line option, [191](#)
- detect-zeroes=DETECT_ZEROES
 - qemu-nbd command line option, [194](#)
- discard=DISCARD
 - qemu-nbd command line option, [194](#)
- export [type=]nbd,id=<id>,node-name=<node-name>[,name=<export-name>][,writable=on|off][,block-size=SIZE]
 - qemu-storage-daemon command line option, [191](#)
- fd=FDNUM
 - virtiofsd command line option, [201](#)
- force-share (-U)
 - qemu-img-common-opts command line option, [181](#)
- fork
 - qemu-nbd command line option, [194](#)
- image-opts
 - qemu-img-common-opts command line option, [181](#)
 - qemu-nbd command line option, [193](#)
- monitor MONITORDEF
 - qemu-storage-daemon command line option, [191](#)
- nbd-server addr.type=inet,addr.host=<host>,addr.port=PORT
 - qemu-storage-daemon command line option, [192](#)
- object OBJECTDEF
 - qemu-img-common-opts command line option, [180](#)
- object help
 - qemu-storage-daemon command line option, [192](#)
- object type,id=ID,...props...
 - qemu-nbd command line option, [193](#)
- pid=FILE=PATH
 - qemu-nbd command line option, [194](#)
- pid=PID, -p PID
 - qemu-trace-stap-run command line option, [198](#)
- salvage
 - qemu-img-convert command line option, [182](#)
- socket-group=GROUP
 - virtiofsd command line option, [201](#)
- socket-path=PATH
 - virtiofsd command line option, [201](#)
- syslog
 - virtiofsd command line option, [200](#)
- target-image-opts
 - qemu-img-common-opts command line option, [181](#)
- target-is-zero
 - qemu-img-convert command line option, [182](#)
- thread-pool-size=NUM
 - virtiofsd command line option, [201](#)
- tls-authz=ID
 - qemu-nbd command line option, [194](#)
- tls-creds=ID
 - qemu-nbd command line option, [194](#)
- verbose, -v

```

    qemu-trace-stap command line
        option, 198
-A, -allocation-depth
    qemu-nbd command line option, 193
-B, -bitmap=NAME
    qemu-nbd command line option, 193
-C
    qemu-img-convert command line
        option, 182
-D, -description=DESCRIPTION
    qemu-nbd command line option, 194
-D, -dump-conf
    qemu-ga command line option, 250
-F
    qemu-img-compare command line
        option, 181
-F, -fsfreeze-hook=PATH
    qemu-ga command line option, 250
-L, -list
    qemu-nbd command line option, 194
-S SIZE
    qemu-img-common-opts command line
        option, 181
-T SRC_CACHE
    qemu-img-common-opts command line
        option, 181
-T, -trace [[enable=]PATTERN][,events=FILE][,hide=FILE][,filter=PATH]
    qemu-img command line option, 179
    qemu-nbd command line option, 195
    qemu-pr-helper command line option,
        197
    qemu-storage-daemon command line
        option, 191
-V, -version
    qemu-ga command line option, 250
    qemu-img command line option, 179
    qemu-nbd command line option, 195
    qemu-pr-helper command line option,
        197
    qemu-storage-daemon command line
        option, 191
    virtiofsd command line option, 200
-W
    qemu-img-convert command line
        option, 182
-a
    qemu-img-snapshot command line
        option, 182
-b, -bind=IFACE
    qemu-nbd command line option, 193
-b, -blacklist=LIST
    qemu-ga command line option, 250
-c
    qemu-img-common-opts command line
        option, 181
    qemu-img-snapshot command line
        option, 182
-c, -connect=DEV
    qemu-nbd command line option, 194
-d
    qemu-img-snapshot command line
        option, 182
    virtiofsd command line option, 200
-d, -daemon
    qemu-ga command line option, 250
    qemu-pr-helper command line option,
        196
-d, -disconnect
    qemu-nbd command line option, 194
-e, -shared=NUM
    qemu-nbd command line option, 194
-f
    qemu-img-compare command line
        option, 181
-f, -fd SOCKET_ID
    virtfs-proxy-helper command line
        option, 199
-f, -format=FMT
    qemu-nbd command line option, 193
-f, -filter=PATH
    qemu-ga command line option, 250
    qemu-pr-helper command line option,
        197
-g, -gid GID
    virtfs-proxy-helper command line
        option, 200
-g, -group=GROUP
    qemu-pr-helper command line option,
        197
-h
    qemu-img-common-opts command line
        option, 181
    virtfs-proxy-helper command line
        option, 199
-h, -help
    qemu-ga command line option, 250
    qemu-img command line option, 179
    qemu-nbd command line option, 195
    qemu-pr-helper command line option,
        197
    qemu-storage-daemon command line
        option, 191
    virtiofsd command line option, 200
-k, -socket=PATH
    qemu-nbd command line option, 193
    qemu-pr-helper command line option,
        197

```

-l
 qemu-img-snapshot command line option, 182
 -l, -load-snapshot=SNAPSHOT_PARAM
 qemu-nbd command line option, 193
 -l, -logfile=PATH
 qemu-ga command line option, 250
 -m
 qemu-img-convert command line option, 182
 -m, -method=METHOD
 qemu-ga command line option, 250
 -n
 qemu-img-convert command line option, 182
 -n, -nocache
 qemu-nbd command line option, 194
 -n, -nodaemon
 virtfs-proxy-helper command line option, 200
 -o OPTION
 virtiofsd command line option, 200
 -o, -offset=OFFSET
 qemu-nbd command line option, 193
 -p
 qemu-img-common-opts command line option, 181
 -p, -path PATH
 virtfs-proxy-helper command line option, 199
 -p, -path=PATH
 qemu-ga command line option, 250
 -p, -port=PORT
 qemu-nbd command line option, 193
 -q
 qemu-img-common-opts command line option, 181
 -q, -quiet
 qemu-pr-helper command line option, 197
 -r
 qemu-img-convert command line option, 182
 -r, -read-only
 qemu-nbd command line option, 193
 -s
 qemu-img-compare command line option, 181
 -s, -snapshot
 qemu-nbd command line option, 193
 -s, -socket SOCKET_FILE
 virtfs-proxy-helper command line option, 199
 -t CACHE

 qemu-img-common-opts command line option, 181
 -t, -persistent
 qemu-nbd command line option, 194
 -t, -statedir=PATH
 qemu-ga command line option, 250
 -u, -uid UID
 virtfs-proxy-helper command line option, 200
 -u, -user=USER
 qemu-pr-helper command line option, 197
 -v, -verbose
 qemu-ga command line option, 250
 qemu-nbd command line option, 194
 qemu-pr-helper command line option, 197
 -x, -export-name=NAME
 qemu-nbd command line option, 194

A

address_space_cache_destroy (*C function*), 905
 address_space_cache_invalidate (*C function*), 905
 address_space_destroy (*C function*), 904
 address_space_init (*C function*), 903
 address_space_read (*C function*), 905
 address_space_read_cached (*C function*), 906
 address_space_remove_listeners (*C function*), 904
 address_space_rw (*C function*), 904
 address_space_write (*C function*), 904
 address_space_write_cached (*C function*), 906
 address_space_write_rom (*C function*), 905
 AddressSpace (*C type*), 884
 amend [-object OBJECTDEF] [-image-opts] [-p] [-q] [-f FMT] [-t CACHE] [-force] -o OPTIONS FILENAME
 qemu-img command line option, 180
 qemu-img-commands command line option, 182

B

backing_file
 image-formats command line option, 65
 qcow command line option, 64
 qcow2 command line option, 62
 qed command line option, 64
 backing_fmt
 qcow2 command line option, 62
 qed command line option, 64

bench [-c COUNT] [-d DEPTH] [-f FMT]
 [-flush-interval=FLUSH_INTERVAL]
 [-i AIO] [-n] [-no-drain] [-o
 OFFSET] [-pattern=PATTERN] [-q]
 [-s BUFFER_SIZE] [-S STEP_SIZE]
 [-t CACHE] [-w] [-U] FILENAME
 qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [183](#)

bitmap (-merge SOURCE | -add
 | -remove | -clear |
 -enable | -disable)... [-b
 SOURCE_FILE [-F SOURCE_FMT]]
 [-g GRANULARITY] [-object
 OBJECTDEF] [-image-opts | -f
 FMT] FILENAME BITMAP
 qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [183](#)

block_size
 VHDX command line option, [66](#)

block_state_zero
 VHDX command line option, [66](#)

bochs
 image-formats command line option,
[66](#)

bs=BLOCK_SIZE
 qemu-img-dd command line option, [182](#)

C

change_bit (*C function*), [968](#)

check [-object OBJECTDEF]
 [-image-opts] [-q] [-f FMT]
 [-output=OFMT] [-r [leaks
 | all]] [-T SRC_CACHE] [-U]
 FILENAME
 qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [183](#)

cipher-alg
 luks command line option, [65](#)

cipher-mode
 luks command line option, [65](#)

clear_bit (*C function*), [968](#)

cloop
 image-formats command line option,
[66](#)

cluster_size
 qcow2 command line option, [63](#)
 qed command line option, [64](#)

commit [-object OBJECTDEF]
 [-image-opts] [-q] [-f FMT]
 [-t CACHE] [-b BASE] [-r
 RATE_LIMIT] [-d] [-p] FILENAME

qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [184](#)

compare [-object OBJECTDEF]
 [-image-opts] [-f FMT] [-F FMT]
 [-T SRC_CACHE] [-p] [-q] [-s]
 [-U] FILENAME1 FILENAME2
 qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [184](#)

compat
 qcow2 command line option, [62](#)

compat6
 image-formats command line option,
[65](#)

container_get (*C function*), [1021](#)

convert [-object OBJECTDEF]
 [-image-opts]
 [-target-image-opts]
 [-target-is-zero] [-bitmaps]
 [-U] [-C] [-c] [-p] [-q]
 [-n] [-f FMT] [-t CACHE] [-T
 SRC_CACHE] [-O OUTPUT_FMT] [-B
 BACKING_FILE] [-o OPTIONS]
 [-l SNAPSHOT_PARAM] [-S
 SPARSE_SIZE] [-r RATE_LIMIT]
 [-m NUM_COROUTINES] [-W]
 FILENAME [FILENAME2 [...]]
 OUTPUT_FILENAME
 qemu-img-commands command line
 option, [184](#)

convert [-object OBJECTDEF]
 [-image-opts]
 [-target-image-opts]
 [-target-is-zero] [-bitmaps]
 [-U] [-C] [-c] [-p] [-q]
 [-n] [-f FMT] [-t CACHE] [-T
 SRC_CACHE] [-O OUTPUT_FMT] [-B
 BACKING_FILE] [-o OPTIONS]
 [-l SNAPSHOT_PARAM] [-S
 SPARSE_SIZE] [-r RATE_LIMIT]
 [-m NUM_COROUTINES] [-W]
 [-salvage] FILENAME [FILENAME2
 [...]] OUTPUT_FILENAME
 qemu-img command line option, [180](#)

count=BLOCKS
 qemu-img-dd command line option, [182](#)

create [-object OBJECTDEF] [-q] [-f
 FMT] [-b BACKING_FILE] [-F
 BACKING_FMT] [-u] [-o OPTIONS]
 FILENAME [SIZE]
 qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [185](#)

D

dd [-image-opts] [-U] [-f FMT] [-O
 OUTPUT_FMT] [bs=BLOCK_SIZE]
 [count=BLOCKS] [skip=BLOCKS]
 if=INPUT of=OUTPUT
 qemu-img command line option, 180
 qemu-img-commands command line
 option, 185
 DECLARE_CLASS_CHECKERS (*C function*), 995
 DECLARE_INSTANCE_CHECKER (*C function*), 994
 DECLARE_OBJ_CHECKERS (*C function*), 995
 DEFINE_TYPES (*C function*), 1005
 deposit32 (*C function*), 972
 deposit64 (*C function*), 972
 dmz
 image-formats command line option,
 66

E

encrypt.cipher-alg
 qcow2 command line option, 63
 encrypt.cipher-mode
 qcow2 command line option, 63
 encrypt.format
 qcow command line option, 64
 qcow2 command line option, 62
 encrypt.hash-alg
 qcow2 command line option, 63
 encrypt.iter-time
 qcow2 command line option, 63
 encrypt.ivgen-alg
 qcow2 command line option, 63
 encrypt.ivgen-hash-alg
 qcow2 command line option, 63
 encrypt.key-secret
 qcow command line option, 64
 qcow2 command line option, 63
 encryption
 qcow command line option, 64
 qcow2 command line option, 62
 extract16 (*C function*), 971
 extract32 (*C function*), 970
 extract64 (*C function*), 971
 extract8 (*C function*), 971

F

filter-drivers command line option
 preallocate, 73
 find_first_bit (*C function*), 969
 find_first_zero_bit (*C function*), 969
 find_last_bit (*C function*), 968
 find_next_bit (*C function*), 969
 find_next_zero_bit (*C function*), 969

H

half_shuffle32 (*C function*), 973
 half_shuffle64 (*C function*), 973
 half_unshuffle32 (*C function*), 973
 half_unshuffle64 (*C function*), 974
 hash-alg
 luks command line option, 65
 hwversion
 image-formats command line option,
 65

I

if=INPUT
 qemu-img-dd command line option, 182
 image-formats command line option
 backing_file, 65
 bochs, 66
 cloop, 66
 compat6, 65
 dmz, 66
 hwversion, 65
 luks, 65
 parallels, 66
 qcow, 64
 qcow2, 62
 qed, 64
 raw, 62
 subformat, 65
 vdi, 65
 VHDX, 66
 vmdk, 65
 vpc, 65
 info [-object OBJECTDEF] [-image-opts]
 [-f FMT] [-output=OFMT]
 [-backing-chain] [-U] FILENAME
 qemu-img command line option, 180
 qemu-img-commands command line
 option, 185
 INTERFACE_CHECK (*C function*), 999
 INTERFACE_CLASS (*C function*), 999
 InterfaceClass (*C type*), 999
 InterfaceInfo (*C type*), 999
 iter-time
 luks command line option, 65
 ivgen-alg
 luks command line option, 65
 ivgen-hash-alg
 luks command line option, 65

K

key-secret
 luks command line option, 65

L

lazy_refcounts
 qcow2 command line option, [63](#)
list BINARY PATTERN...
 qemu-trace-stap command line
 option, [198](#)
log_size
 VHDX command line option, [66](#)
luks
 image-formats command line option,
 [65](#)
luks command line option
 cipher-alg, [65](#)
 cipher-mode, [65](#)
 hash-alg, [65](#)
 iter-time, [65](#)
 ivgen-alg, [65](#)
 ivgen-hash-alg, [65](#)
 key-secret, [65](#)

M

map [-object OBJECTDEF] [-image-opts]
 [-f FMT] [-start-offset=OFFSET]
 [-max-length=LEN]
 [-output=OFMT] [-U] FILENAME
 qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [186](#)
measure [-output=OFMT] [-O OUTPUT_FMT]
 [-o OPTIONS] [-size N
 | [-object OBJECTDEF]
 [-image-opts] [-f FMT] [-l
 SNAPSHOT_PARAM] FILENAME]
 qemu-img command line option, [180](#)
 qemu-img-commands command line
 option, [187](#)
memory_global_after_dirty_log_sync (C
 function), [902](#)
memory_global_dirty_log_start (C function),
 [903](#)
memory_global_dirty_log_stop (C function),
 [903](#)
memory_global_dirty_log_sync (C function),
 [902](#)
memory_listener_register (C function), [902](#)
memory_listener_unregister (C function), [903](#)
memory_region_add_coalescing (C function),
 [899](#)
memory_region_add_eventfd (C function), [900](#)
memory_region_add_subregion (C function),
 [900](#)
memory_region_add_subregion_overlap (C
 function), [901](#)

memory_region_clear_coalescing (C func-
 tion), [899](#)
memory_region_clear_dirty_bitmap (C func-
 tion), [897](#)
memory_region_clear_flush_coalesced (C
 function), [899](#)
memory_region_del_eventfd (C function), [900](#)
memory_region_del_subregion (C function),
 [901](#)
memory_region_dispatch_read (C function),
 [903](#)
memory_region_dispatch_write (C function),
 [903](#)
memory_region_find (C function), [902](#)
memory_region_flush_rom_device (C func-
 tion), [898](#)
memory_region_from_host (C function), [895](#)
memory_region_get_dirty_log_mask (C func-
 tion), [895](#)
memory_region_get_fd (C function), [895](#)
memory_region_get_iommu (C function), [892](#)
memory_region_get_iommu_class_nocheck
 (C function), [892](#)
memory_region_get_ram_addr (C function), [901](#)
memory_region_get_ram_ptr (C function), [896](#)
memory_region_init (C function), [885](#)
memory_region_init_alias (C function), [889](#)
memory_region_init_io (C function), [886](#)
memory_region_init_iommu (C function), [890](#)
memory_region_init_ram (C function), [890](#)
memory_region_init_ram_device_ptr (C
 function), [888](#)
memory_region_init_ram_from_fd (C func-
 tion), [888](#)
memory_region_init_ram_from_file (C func-
 tion), [887](#)
memory_region_init_ram_nomigrate (C func-
 tion), [886](#)
memory_region_init_ram_ptr (C function), [888](#)
memory_region_init_ram_shared_nomigrate
 (C function), [886](#)
memory_region_init_resizeable_ram (C
 function), [887](#)
memory_region_init_rom (C function), [891](#)
memory_region_init_rom_device (C function),
 [891](#)
memory_region_init_rom_device_nomigrate
 (C function), [889](#)
memory_region_init_rom_nomigrate (C func-
 tion), [889](#)
memory_region_iommu_attrs_to_index (C
 function), [894](#)
memory_region_iommu_get_attr (C function),
 [894](#)

- `memory_region_iommu_get_min_page_size` (C function), 892
 - `memory_region_iommu_num_indexes` (C function), 894
 - `memory_region_iommu_replay` (C function), 893
 - `memory_region_iommu_set_page_size_mask` (C function), 894
 - `memory_region_is_logging` (C function), 895
 - `memory_region_is_mapped` (C function), 901
 - `memory_region_is_nonvolatile` (C function), 895
 - `memory_region_is_ram` (C function), 892
 - `memory_region_is_ram_device` (C function), 892
 - `memory_region_is_rom` (C function), 895
 - `memory_region_is_romd` (C function), 892
 - `memory_region_msync` (C function), 896
 - `memory_region_name` (C function), 894
 - `memory_region_notify_iommu` (C function), 893
 - `memory_region_notify_iommu_one` (C function), 893
 - `memory_region_owner` (C function), 891
 - `memory_region_present` (C function), 901
 - `memory_region_ref` (C function), 885
 - `memory_region_register_iommu_notifier` (C function), 893
 - `memory_region_reset_dirty` (C function), 898
 - `memory_region_rom_device_set_romd` (C function), 898
 - `memory_region_set_coalescing` (C function), 899
 - `memory_region_set_dirty` (C function), 896
 - `memory_region_set_flush_coalesced` (C function), 899
 - `memory_region_set_log` (C function), 896
 - `memory_region_set_nonvolatile` (C function), 898
 - `memory_region_set_readonly` (C function), 898
 - `memory_region_size` (C function), 892
 - `memory_region_snapshot_and_clear_dirty` (C function), 897
 - `memory_region_snapshot_get_dirty` (C function), 897
 - `memory_region_transaction_begin` (C function), 902
 - `memory_region_transaction_commit` (C function), 902
 - `memory_region_unref` (C function), 885
 - `memory_region_unregister_iommu_notifier` (C function), 894
 - `memory_region_writeback` (C function), 896
 - `MemoryListener` (C type), 882
 - `MemoryRegionSection` (C type), 884
 - `module_object_class_by_name` (C function), 1006
- N**
- `nocow`
 - qcow2 command line option, 63
- O**
- `OBJECT` (C function), 998
 - `Object` (C type), 994
 - `OBJECT_CHECK` (C function), 998
 - `object_child_foreach` (C function), 1021
 - `object_child_foreach_recursive` (C function), 1021
 - `OBJECT_CLASS` (C function), 998
 - `object_class_by_name` (C function), 1006
 - `OBJECT_CLASS_CHECK` (C function), 998
 - `object_class_dynamic_cast` (C function), 1005
 - `object_class_dynamic_cast_assert` (C function), 1005
 - `object_class_get_list` (C function), 1006
 - `object_class_get_list_sorted` (C function), 1006
 - `object_class_get_name` (C function), 1006
 - `object_class_get_parent` (C function), 1006
 - `object_class_is_abstract` (C function), 1006
 - `object_class_property_find` (C function), 1009
 - `object_class_property_find_err` (C function), 1009
 - `object_class_property_iter_init` (C function), 1010
 - `OBJECT_DECLARE_SIMPLE_TYPE` (C function), 995
 - `OBJECT_DECLARE_TYPE` (C function), 995
 - `OBJECT_DEFINE_ABSTRACT_TYPE` (C function), 997
 - `OBJECT_DEFINE_TYPE` (C function), 996
 - `OBJECT_DEFINE_TYPE_EXTENDED` (C function), 996
 - `OBJECT_DEFINE_TYPE_WITH_INTERFACES` (C function), 996
 - `object_dynamic_cast` (C function), 1003
 - `object_dynamic_cast_assert` (C function), 1004
 - `object_get_canonical_path` (C function), 1015
 - `object_get_canonical_path_component` (C function), 1015
 - `OBJECT_GET_CLASS` (C function), 999
 - `object_get_class` (C function), 1004
 - `object_get_internal_root` (C function), 1015
 - `object_get_objects_root` (C function), 1014
 - `object_get_root` (C function), 1014
 - `object_get_typename` (C function), 1004
 - `object_initialize` (C function), 1002
 - `object_initialize_child` (C function), 1003

`object_initialize_child_with_props` (C function), 1002

`object_initialize_child_with_propsv` (C function), 1003

`object_new` (C function), 1000

`object_new_with_class` (C function), 1000

`object_new_with_props` (C function), 1000

`object_new_with_propv` (C function), 1001

`object_property_add` (C function), 1007

`object_property_add_alias` (C function), 1020

`object_property_add_bool` (C function), 1018

`object_property_add_child` (C function), 1016

`object_property_add_const_link` (C function), 1020

`object_property_add_enum` (C function), 1018

`object_property_add_link` (C function), 1017

`object_property_add_str` (C function), 1018

`object_property_add_tm` (C function), 1018

`object_property_add_uint16_ptr` (C function), 1019

`object_property_add_uint32_ptr` (C function), 1019

`object_property_add_uint64_ptr` (C function), 1020

`object_property_add_uint8_ptr` (C function), 1019

`object_property_allow_set_link` (C function), 1017

`object_property_find` (C function), 1009

`object_property_find_err` (C function), 1009

`object_property_get` (C function), 1010

`object_property_get_bool` (C function), 1012

`object_property_get_enum` (C function), 1013

`object_property_get_int` (C function), 1012

`object_property_get_link` (C function), 1011

`object_property_get_str` (C function), 1011

`object_property_get_type` (C function), 1014

`object_property_get_uint` (C function), 1013

`object_property_help` (C function), 1022

`object_property_iter_init` (C function), 1009

`object_property_iter_next` (C function), 1010

`object_property_parse` (C function), 1014

`object_property_print` (C function), 1014

`object_property_set` (C function), 1013

`object_property_set_bool` (C function), 1012

`object_property_set_default_bool` (C function), 1008

`object_property_set_default_int` (C function), 1008

`object_property_set_default_str` (C function), 1008

`object_property_set_default_uint` (C function), 1008

`object_property_set_description` (C function), 1021

`object_property_set_int` (C function), 1012

`object_property_set_link` (C function), 1011

`object_property_set_str` (C function), 1011

`object_property_set_uint` (C function), 1013

`object_property_try_add` (C function), 1007

`object_property_try_add_child` (C function), 1016

`object_ref` (C function), 1007

`object_resolve_path` (C function), 1015

`object_resolve_path_component` (C function), 1016

`object_resolve_path_type` (C function), 1016

`object_set_props` (C function), 1001

`object_set_propv` (C function), 1002

`object_type_get_instance_size` (C function), 1022

`object_unref` (C function), 1007

`ObjectClass` (C type), 994

`ObjectFree` (C type), 994

`ObjectPropertyAccessor` (C type), 992

`ObjectPropertyInit` (C type), 993

`ObjectPropertyRelease` (C type), 993

`ObjectPropertyResolve` (C type), 992

`ObjectUnparent` (C type), 993

`of=OUTPUT`

- `qemu-img-dd` command line option, 182

P

`parallels`

- `image-formats` command line option, 66

`prealloc-align`

- `preallocate` command line option, 74

`prealloc-size`

- `preallocate` command line option, 74

`preallocate`

- `filter-drivers` command line option, 73

`preallocate` command line option

- `prealloc-align`, 74
- `prealloc-size`, 74

`preallocation`

- `qcow2` command line option, 63
- `raw` command line option, 62

Q

`qcow`

- `image-formats` command line option, 64

`qcow` command line option

- `backing_file`, 64
- `encrypt.format`, 64

```

    encrypt.key-secret, 64
    encryption, 64
qcow2
    image-formats command line option,
        62
qcow2 command line option
    backing_file, 62
    backing_fmt, 62
    cluster_size, 63
    compat, 62
    encrypt.cipher-alg, 63
    encrypt.cipher-mode, 63
    encrypt.format, 62
    encrypt.hash-alg, 63
    encrypt.iter-time, 63
    encrypt.ivgen-alg, 63
    encrypt.ivgen-hash-alg, 63
    encrypt.key-secret, 63
    encryption, 62
    lazy_refcounts, 63
    nocow, 63
    preallocation, 63
qed
    image-formats command line option,
        64
qed command line option
    backing_file, 64
    backing_fmt, 64
    cluster_size, 64
    table_size, 64
qemu-ga command line option
    -D, -dump-conf, 250
    -F, -fsfreeze-hook=PATH, 250
    -V, -version, 250
    -b, -blacklist=LIST, 250
    -d, -daemon, 250
    -f, -pidfile=PATH, 250
    -h, -help, 250
    -l, -logfile=PATH, 250
    -m, -method=METHOD, 250
    -p, -path=PATH, 250
    -t, -statedir=PATH, 250
    -v, -verbose, 250
qemu-img command line option
    -T, -trace [[enable=]PATTERN][,events=FILE], 179
    -V, -version, 179
    -h, -help, 179
    amend [-object OBJECTDEF]
        [-image-opts] [-p] [-q] [-f
        FMT] [-t CACHE] [-force] -o
        OPTIONS FILENAME, 180
    bench [-c COUNT] [-d
        DEPTH] [-f FMT]
        [-flush-interval=FLUSH_INTERVAL]
        [-i AIO] [-n] [-no-drain] [-o
        OFFSET] [-pattern=PATTERN] [-q]
        [-s BUFFER_SIZE] [-S STEP_SIZE]
        [-t CACHE] [-w] [-U] FILENAME,
        180
    bitmap (-merge SOURCE | -add
        | -remove | -clear |
        -enable | -disable)... [-b
        SOURCE_FILE [-F SOURCE_FMT]]
        [-g GRANULARITY] [-object
        OBJECTDEF] [-image-opts | -f
        FMT] FILENAME BITMAP, 180
    check [-object OBJECTDEF]
        [-image-opts] [-q] [-f FMT]
        [-output=OFMT] [-r [leaks
        | all]] [-T SRC_CACHE] [-U]
        FILENAME, 180
    commit [-object OBJECTDEF]
        [-image-opts] [-q] [-f FMT]
        [-t CACHE] [-b BASE] [-r
        RATE_LIMIT] [-d] [-p] FILENAME,
        180
    compare [-object OBJECTDEF]
        [-image-opts] [-f FMT] [-F FMT]
        [-T SRC_CACHE] [-p] [-q] [-s]
        [-U] FILENAME1 FILENAME2, 180
    convert [-object OBJECTDEF]
        [-image-opts]
        [-target-image-opts]
        [-target-is-zero] [-bitmaps]
        [-U] [-C] [-c] [-p] [-q]
        [-n] [-f FMT] [-t CACHE] [-T
        SRC_CACHE] [-O OUTPUT_FMT] [-B
        BACKING_FILE] [-o OPTIONS]
        [-l SNAPSHOT_PARAM] [-S
        SPARSE_SIZE] [-r RATE_LIMIT]
        [-m NUM_COROUTINES] [-W]
        [-salvage] FILENAME [FILENAME2
        ...] OUTPUT_FILENAME, 180
    create [-object OBJECTDEF] [-q]
        [-f FMT] [-b BACKING_FILE] [-F
        BACKING_FMT] [-u] [-o OPTIONS]
        FILENAME [SIZE], 180
    def [-object OBJECTDEF] [-q]
        [-f FMT] [-b BACKING_FILE] [-F
        BACKING_FMT] [-u] [-o OPTIONS]
        FILENAME [SIZE], 180
    diff [-object OBJECTDEF] [-q]
        [-f FMT] [-b BACKING_FILE] [-F
        BACKING_FMT] [-u] [-o OPTIONS]
        FILENAME [SIZE], 180
    info [-object OBJECTDEF]
        [-image-opts] [-f FMT]
        [-output=OFMT] [-backing-chain]
        [-U] FILENAME, 180
    map [-object OBJECTDEF]
        [-image-opts] [-f FMT]

```

```

        [-start-offset=OFFSET]
        [-max-length=LEN]
        [-output=OFMT] [-U] FILENAME,
180
measure [-output=OFMT] [-O
        OUTPUT_FMT] [-o OPTIONS]
        [-size N | [-object OBJECTDEF]
        [-image-opts] [-f FMT] [-l
        SNAPSHOT_PARAM] FILENAME], 180
rebase [-object OBJECTDEF]
        [-image-opts] [-U] [-q] [-f
        FMT] [-t CACHE] [-T SRC_CACHE]
        [-p] [-u] -b BACKING_FILE [-F
        BACKING_FMT] FILENAME, 180
resize [-object OBJECTDEF]
        [-image-opts] [-f FMT]
        [-preallocation=PREALLOC] [-q]
        [-shrink] FILENAME [+ | -]SIZE,
180
snapshot [-object OBJECTDEF]
        [-image-opts] [-U] [-q] [-l |
        -a SNAPSHOT | -c SNAPSHOT | -d
        SNAPSHOT] FILENAME, 180
qemu-img-commands command line option
amend [-object OBJECTDEF]
        [-image-opts] [-p] [-q] [-f
        FMT] [-t CACHE] [-force] -o
        OPTIONS FILENAME, 182
bench [-c COUNT] [-d
        DEPTH] [-f FMT]
        [-flush-interval=FLUSH_INTERVAL]
        [-i AIO] [-n] [-no-drain] [-o
        OFFSET] [-pattern=PATTERN] [-q]
        [-s BUFFER_SIZE] [-S STEP_SIZE]
        [-t CACHE] [-w] [-U] FILENAME,
183
bitmap (-merge SOURCE | -add
        | -remove | -clear |
        -enable | -disable)... [-b
        SOURCE_FILE [-F SOURCE_FMT]]
        [-g GRANULARITY] [-object
        OBJECTDEF] [-image-opts | -f
        FMT] FILENAME BITMAP, 183
check [-object OBJECTDEF]
        [-image-opts] [-q] [-f FMT]
        [-output=OFMT] [-r [leaks
        | all]] [-T SRC_CACHE] [-U]
        FILENAME, 183
commit [-object OBJECTDEF]
        [-image-opts] [-q] [-f FMT]
        [-t CACHE] [-b BASE] [-r
        RATE_LIMIT] [-d] [-p] FILENAME,
184
compare [-object OBJECTDEF]
        [-image-opts] [-f FMT] [-F FMT]
        [-T SRC_CACHE] [-p] [-q] [-s]
        [-U] FILENAME1 FILENAME2, 184
convert [-object OBJECTDEF]
        [-image-opts]
        [-target-image-opts]
        [-target-is-zero] [-bitmaps]
        [-U] [-C] [-c] [-p] [-q]
        [-n] [-f FMT] [-t CACHE] [-T
        SRC_CACHE] [-O OUTPUT_FMT] [-B
        BACKING_FILE] [-o OPTIONS]
        [-l SNAPSHOT_PARAM] [-S
        SPARSE_SIZE] [-r RATE_LIMIT]
        [-m NUM_COROUTINES] [-W]
        FILENAME [FILENAME2 [...]]
        OUTPUT_FILENAME, 184
create [-object OBJECTDEF] [-q]
        [-f FMT] [-b BACKING_FILE] [-F
        BACKING_FMT] [-u] [-o OPTIONS]
        FILENAME [SIZE], 185
dd [-image-opts] [-U] [-f FMT] [-O
        OUTPUT_FMT] [bs=BLOCK_SIZE]
        [count=BLOCKS] [skip=BLOCKS]
        if=INPUT of=OUTPUT, 185
info [-object OBJECTDEF]
        [-image-opts] [-f FMT]
        [-output=OFMT] [-backing-chain]
        [-U] FILENAME, 185
map [-object OBJECTDEF]
        [-image-opts] [-f FMT]
        [-start-offset=OFFSET]
        [-max-length=LEN]
        [-output=OFMT] [-U] FILENAME,
186
measure [-output=OFMT] [-O
        OUTPUT_FMT] [-o OPTIONS]
        [-size N | [-object OBJECTDEF]
        [-image-opts] [-f FMT] [-l
        SNAPSHOT_PARAM] FILENAME], 187
rebase [-object OBJECTDEF]
        [-image-opts] [-U] [-q] [-f
        FMT] [-t CACHE] [-T SRC_CACHE]
        [-p] [-u] -b BACKING_FILE [-F
        BACKING_FMT] FILENAME, 187
resize [-object OBJECTDEF]
        [-image-opts] [-f FMT]
        [-preallocation=PREALLOC] [-q]
        [-shrink] FILENAME [+ | -]SIZE,
188
snapshot [-object OBJECTDEF]
        [-image-opts] [-U] [-q] [-l |
        -a SNAPSHOT | -c SNAPSHOT | -d
        SNAPSHOT] FILENAME, 187
qemu-img-common-opts command line

```

```

    option
    -backing-chain, 181
    -force-share (-U), 181
    -image-opts, 181
    -object OBJECTDEF, 180
    -target-image-opts, 181
    -S SIZE, 181
    -T SRC_CACHE, 181
    -c, 181
    -h, 181
    -p, 181
    -q, 181
    -t CACHE, 181
qemu-img-compare command line option
    -F, 181
    -f, 181
    -s, 181
qemu-img-convert command line option
    -bitmaps, 181
    -salvage, 182
    -target-is-zero, 182
    -C, 182
    -W, 182
    -m, 182
    -n, 182
    -r, 182
qemu-img-dd command line option
    bs=BLOCK_SIZE, 182
    count=BLOCKS, 182
    if=INPUT, 182
    of=OUTPUT, 182
    skip=BLOCKS, 182
qemu-img-snapshot command line option
    -a, 182
    -c, 182
    -d, 182
    -l, 182
    snapshot, 182
qemu-nbd command line option
    -aio=AIO, 194
    -cache=CACHE, 194
    -detect-zeroes=DETECT_ZEROES, 194
    -discard=DISCARD, 194
    -fork, 194
    -image-opts, 193
    -object type,id=ID,...props..., 193
    -pid-file=PATH, 194
    -tls-authz=ID, 194
    -tls-creds=ID, 194
    -A, -allocation-depth, 193
    -B, -bitmap=NAME, 193
    -D, -description=DESCRIPTION, 194
    -L, -list, 194
    -T, -trace [[enable=]PATTERN][,events=FILE][,fi
        195
    -V, -version, 195
    -b, -bind=IFACE, 193
    -c, -connect=DEV, 194
    -d, -disconnect, 194
    -e, -shared=NUM, 194
    -f, -format=FMT, 193
    -h, -help, 195
    -k, -socket=PATH, 193
    -l, -load-snapshot=SNAPSHOT_PARAM,
        193
    -n, -nocache, 194
    -o, -offset=OFFSET, 193
    -p, -port=PORT, 193
    -r, -read-only, 193
    -s, -snapshot, 193
    -t, -persistent, 194
    -v, -verbose, 194
    -x, -export-name=NAME, 194
qemu-pr-helper command line option
    -T, -trace [[enable=]PATTERN][,events=FILE][,fi
        197
    -V, -version, 197
    -d, -daemon, 196
    -f, -pidfile=PATH, 197
    -g, -group=GROUP, 197
    -h, -help, 197
    -k, -socket=PATH, 197
    -q, -quiet, 197
    -u, -user=USER, 197
    -v, -verbose, 197
qemu-storage-daemon command line
    option
    -blockdev BLOCKDEVDEF, 191
    -chardev CHARDEVDEF, 191
    -export [type=]nbd,id=<id>,node-name=<node-name>
        191
    -monitor MONITORDEF, 191
    -nbd-server addr.type=inet,addr.host=<host>,add
        192
    -object help, 192
    -T, -trace [[enable=]PATTERN][,events=FILE][,fi
        191
    -V, -version, 191
    -h, -help, 191
qemu-trace-stap command line option
    -verbose, -v, 198
    list BINARY PATTERN..., 198
    run OPTIONS BINARY PATTERN..., 198
qemu-trace-stap-run command line
    option
    -pid=PID, -p PID, 198
qmp_expect_error_and_unref (C function), 942

```

`qmp_rsp_is_err` (*C function*), 941
`qtest_add` (*C function*), 940
`qtest_add_data_func` (*C function*), 940
`qtest_add_data_func_full` (*C function*), 940
`qtest_add_func` (*C function*), 939
`qtest_big_endian` (*C function*), 939
`qtest_bufread` (*C function*), 938
`qtest_bufwrite` (*C function*), 938
`qtest_cb_for_every_machine` (*C function*), 941
`qtest_clock_set` (*C function*), 939
`qtest_clock_step` (*C function*), 939
`qtest_clock_step_next` (*C function*), 939
`qtest_get_arch` (*C function*), 939
`qtest_get_irq` (*C function*), 934
`qtest_hmp` (*C function*), 933
`qtest_inb` (*C function*), 935
`qtest_init` (*C function*), 930
`qtest_init_with_serial` (*C function*), 930
`qtest_init_without_qmp_handshake` (*C function*), 930
`qtest_initf` (*C function*), 929
`qtest_inl` (*C function*), 935
`qtest_inw` (*C function*), 935
`qtest_irq_intercept_in` (*C function*), 934
`qtest_irq_intercept_out` (*C function*), 934
`qtest_memread` (*C function*), 937
`qtest_memset` (*C function*), 938
`qtest_memwrite` (*C function*), 938
`qtest_outb` (*C function*), 934
`qtest_outl` (*C function*), 935
`qtest_outw` (*C function*), 935
`qtest_probe_child` (*C function*), 942
`qtest_qmp` (*C function*), 931
`qtest_qmp_assert_success` (*C function*), 940
`qtest_qmp_device_add` (*C function*), 941
`qtest_qmp_device_add_qdict` (*C function*), 941
`qtest_qmp_device_del` (*C function*), 941
`qtest_qmp_event_ref` (*C function*), 933
`qtest_qmp_eventwait` (*C function*), 933
`qtest_qmp_eventwait_ref` (*C function*), 933
`qtest_qmp_fds` (*C function*), 930
`qtest_qmp_receive` (*C function*), 932
`qtest_qmp_receive_dict` (*C function*), 932
`qtest_qmp_send` (*C function*), 931
`qtest_qmp_send_raw` (*C function*), 931
`qtest_qmp_vsend` (*C function*), 932
`qtest_qmp_vsend_fds` (*C function*), 932
`qtest_quit` (*C function*), 930
`qtest_readb` (*C function*), 936
`qtest_readl` (*C function*), 937
`qtest_readq` (*C function*), 937
`qtest_readw` (*C function*), 936
`qtest_rtas_call` (*C function*), 937
`qtest_set_expected_status` (*C function*), 942

`qtest_set_irq_in` (*C function*), 934
`qtest_vhmp` (*C function*), 933
`qtest_vinitf` (*C function*), 930
`qtest_vqmp` (*C function*), 932
`qtest_vqmp_fds` (*C function*), 931
`qtest_writeb` (*C function*), 935
`qtest_writel` (*C function*), 936
`qtest_writeq` (*C function*), 936
`qtest_writew` (*C function*), 936

R

`raw`
 image-formats command line option, 62
`raw` command line option
 preallocation, 62
`rebase` [-object OBJECTDEF]
 [-image-opts] [-U] [-q] [-f FMT] [-t CACHE] [-T SRC_CACHE]
 [-p] [-u] -b BACKING_FILE [-F BACKING_FMT] FILENAME
 qemu-img command line option, 180
 qemu-img-commands command line option, 187
`resize` [-object OBJECTDEF]
 [-image-opts] [-f FMT]
 [-preallocation=PREALLOC] [-q]
 [-shrink] FILENAME [+ | -]SIZE
 qemu-img command line option, 180
 qemu-img-commands command line option, 188
`rol16` (*C function*), 969
`rol32` (*C function*), 970
`rol64` (*C function*), 970
`rol8` (*C function*), 969
`ror16` (*C function*), 970
`ror32` (*C function*), 970
`ror64` (*C function*), 970
`ror8` (*C function*), 969
`run` OPTIONS BINARY PATTERN...
 qemu-trace-stap command line option, 198

S

`set_bit` (*C function*), 967
`set_bit_atomic` (*C function*), 967
`sextract32` (*C function*), 971
`sextract64` (*C function*), 972
`skip=BLOCKS`
 qemu-img-dd command line option, 182
`snapshot`
 qemu-img-snapshot command line option, 182

snapshot [-object OBJECTDEF]
 [-image-opts] [-U] [-q] [-l |
 -a SNAPSHOT | -c SNAPSHOT | -d
 SNAPSHOT] FILENAME

qemu-img command line option, 180

qemu-img-commands command line
 option, 187

static

vdi command line option, 65

subformat

image-formats command line option,
 65

VHDX command line option, 66

vpc command line option, 65

T

table_size

qed command line option, 64

test_and_change_bit (*C function*), 968

test_and_clear_bit (*C function*), 968

test_and_set_bit (*C function*), 968

test_bit (*C function*), 968

type_register (*C function*), 1004

type_register_static (*C function*), 1004

type_register_static_array (*C function*),
 1005

TypeInfo (*C type*), 997

V

vdi

image-formats command line option,
 65

vdi command line option

static, 65

VHDX

image-formats command line option,
 66

VHDX command line option

block_size, 66

block_state_zero, 66

log_size, 66

subformat, 66

virtfs-proxy-helper command line
 option

-f, -fd SOCKET_ID, 199

-g, -gid GID, 200

-h, 199

-n, -nodaemon, 200

-p, -path PATH, 199

-s, -socket SOCKET_FILE, 199

-u, -uid UID, 200

virtiofsd command line option

-cache=none|auto|always, 201

-fd=FDNUM, 201

-socket-group=GROUP, 201

-socket-path=PATH, 201

-syslog, 200

-thread-pool-size=NUM, 201

-V, -version, 200

-d, 200

-h, -help, 200

-o OPTION, 200

vmrk

image-formats command line option,
 65

vpc

image-formats command line option,
 65

vpc command line option

subformat, 65