

---

# **Linux Virt Documentation**

**The kernel development community**

**Jun 10, 2024**



# CONTENTS

<b>1</b>	<b>KVM</b>	<b>1</b>
1.1	The Definitive KVM (Kernel-based Virtual Machine) API Documentation	1
1.2	Secure Encrypted Virtualization (SEV)	137
1.3	KVM CPUID bits	142
1.4	The KVM halt polling system	144
1.5	Linux KVM Hypercall	146
1.6	KVM Lock Overview	150
1.7	The x86 kvm shadow mmu	155
1.8	KVM-specific MSRs	164
1.9	Nested VMX	170
1.10	The PPC KVM paravirtual interface	175
1.11	The s390 DIAGNOSE call on KVM	180
1.12	s390 (IBM Z) Ultravisor and Protected VMs	181
1.13	s390 (IBM Z) Boot/IPL of Protected VMs	183
1.14	Timekeeping Virtualization for X86-Based Architectures	185
1.15	KVM VCPU Requests	197
1.16	Review checklist for kvm patches	202
1.17	ARM	203
1.18	Devices	207
1.19	Running nested guests with KVM	241
<b>2</b>	<b>UML HowTo</b>	<b>247</b>
2.1	Introduction	249
2.2	Building a UML instance	250
2.3	Setting Up UML Networking	252
2.4	Running UML	258
2.5	Advanced UML Topics	263
2.6	Contributing to UML and Developing with UML	266
<b>3</b>	<b>Paravirt_ops</b>	<b>271</b>
<b>4</b>	<b>Guest halt polling</b>	<b>273</b>
4.1	Module Parameters	273
4.2	Further Notes	274
<b>5</b>	<b>Nitro Enclaves</b>	<b>275</b>
5.1	Overview	275
	<b>Bibliography</b>	<b>277</b>



## **1.1 The Definitive KVM (Kernel-based Virtual Machine) API Documentation**

### **1.1.1 1. General description**

The kvm API is a set of ioctls that are issued to control various aspects of a virtual machine. The ioctls belong to the following classes:

- System ioctls: These query and set global attributes which affect the whole kvm subsystem. In addition a system ioctl is used to create virtual machines.
- VM ioctls: These query and set attributes that affect an entire virtual machine, for example memory layout. In addition a VM ioctl is used to create virtual cpus (vcpus) and devices.

VM ioctls must be issued from the same process (address space) that was used to create the VM.

- vcpu ioctls: These query and set attributes that control the operation of a single virtual cpu.

vcpu ioctls should be issued from the same thread that was used to create the vcpu, except for asynchronous vcpu ioctl that are marked as such in the documentation. Otherwise, the first ioctl after switching threads could see a performance impact.

- device ioctls: These query and set attributes that control the operation of a single device.

device ioctls must be issued from the same process (address space) that was used to create the VM.

### 1.1.2 2. File descriptors

The kvm API is centered around file descriptors. An initial `open( "/dev/kvm" )` obtains a handle to the kvm subsystem; this handle can be used to issue system `ioctl`s. A `KVM_CREATE_VM` `ioctl` on this handle will create a VM file descriptor which can be used to issue VM `ioctl`s. A `KVM_CREATE_VCPU` or `KVM_CREATE_DEVICE` `ioctl` on a VM fd will create a virtual cpu or device and return a file descriptor pointing to the new resource. Finally, `ioctl`s on a vcpu or device fd can be used to control the vcpu or device. For vcpus, this includes the important task of actually running guest code.

In general file descriptors can be migrated among processes by means of `fork()` and the `SCM_RIGHTS` facility of unix domain socket. These kinds of tricks are explicitly not supported by kvm. While they will not cause harm to the host, their actual behavior is not guaranteed by the API. See “General description” for details on the `ioctl` usage model that is supported by KVM.

It is important to note that although VM `ioctl`s may only be issued from the process that created the VM, a VM’s lifecycle is associated with its file descriptor, not its creator (process). In other words, the VM and its resources, *including the associated address space*, are not freed until the last reference to the VM’s file descriptor has been released. For example, if `fork()` is issued after `ioctl(KVM_CREATE_VM)`, the VM will not be freed until both the parent (original) process and its child have put their references to the VM’s file descriptor.

Because a VM’s resources are not freed until the last reference to its file descriptor is released, creating additional references to a VM via `fork()`, `dup()`, etc ...without careful consideration is strongly discouraged and may have unwanted side effects, e.g. memory allocated by and on behalf of the VM’s process may not be freed/unaccounted when the VM is shut down.

### 1.1.3 3. Extensions

As of Linux 2.6.22, the KVM ABI has been stabilized: no backward incompatible change are allowed. However, there is an extension facility that allows backward-compatible extensions to the API to be queried and used.

The extension mechanism is not based on the Linux version number. Instead, kvm defines extension identifiers and a facility to query whether a particular extension identifier is available. If it is, a set of `ioctl`s is available for application use.

### 1.1.4 4. API description

This section describes `ioctl`s that can be used to control kvm guests. For each `ioctl`, the following information is provided along with a description:

**Capability:**

which KVM extension provides this `ioctl`. Can be ‘basic’, which means that it will be provided by any kernel that supports API version 12 (see section 4.1), a `KVM_CAP_xyz` constant, which means availability needs to be checked with `KVM_CHECK_EXTENSION` (see section 4.4), or ‘none’ which means that while not all kernels

support this ioctl, there's no capability bit to check its availability: for kernels that don't support the ioctl, the ioctl returns -ENOTTY.

**Architectures:**

which instruction set architectures provide this ioctl. x86 includes both i386 and x86\_64.

**Type:**

system, vm, or vcpu.

**Parameters:**

what parameters are accepted by the ioctl.

**Returns:**

the return value. General error numbers (EBADF, ENOMEM, EINVAL) are not detailed, but errors with specific meanings are.

## 4.1 KVM\_GET\_API\_VERSION

**Capability**

basic

**Architectures**

all

**Type**

system ioctl

**Parameters**

none

**Returns**

the constant KVM\_API\_VERSION (=12)

This identifies the API version as the stable kvm API. It is not expected that this number will change. However, Linux 2.6.20 and 2.6.21 report earlier versions; these are not documented and not supported. Applications should refuse to run if KVM\_GET\_API\_VERSION returns a value other than 12. If this check passes, all ioctls described as 'basic' will be available.

## 4.2 KVM\_CREATE\_VM

**Capability**

basic

**Architectures**

all

**Type**

system ioctl

**Parameters**

machine type identifier (KVM\_VM\_\*)

**Returns**

a VM fd that can be used to control the new virtual machine.

The new VM has no virtual cpus and no memory. You probably want to use 0 as machine type.

In order to create user controlled virtual machines on S390, check `KVM_CAP_S390_UCONTROL` and use the flag `KVM_VM_S390_UCONTROL` as privileged user (`CAP_SYS_ADMIN`).

To use hardware assisted virtualization on MIPS (VZ ASE) rather than the default trap & emulate implementation (which changes the virtual memory layout to fit in user mode), check `KVM_CAP_MIPS_VZ` and use the flag `KVM_VM_MIPS_VZ`.

On arm64, the physical address size for a VM (IPA Size limit) is limited to 40bits by default. The limit can be configured if the host supports the extension `KVM_CAP_ARM_VM_IPA_SIZE`. When supported, use `KVM_VM_TYPE_ARM_IPA_SIZE(IPA_Bits)` to set the size in the machine type identifier, where `IPA_Bits` is the maximum width of any physical address used by the VM. The `IPA_Bits` is encoded in bits[7-0] of the machine type identifier.

e.g, to configure a guest to use 48bit physical address size:

```
vm_fd = ioctl(dev_fd, KVM_CREATE_VM, KVM_VM_TYPE_ARM_IPA_SIZE(48));
```

The requested size (`IPA_Bits`) must be:

0	Implies default size, 40bits (for backward compatibility)
N	Implies N bits, where N is a positive integer such that, $32 \leq N \leq \text{Host\_IPA\_Limit}$

`Host_IPA_Limit` is the maximum possible value for `IPA_Bits` on the host and is dependent on the CPU capability and the kernel configuration. The limit can be retrieved using `KVM_CAP_ARM_VM_IPA_SIZE` of the `KVM_CHECK_EXTENSION` `ioctl()` at run-time.

Creation of the VM will fail if the requested IPA size (whether it is implicit or explicit) is unsupported on the host.

Please note that configuring the IPA size does not affect the capability exposed by the guest CPUs in `ID_AA64MMFR0_EL1[PARange]`. It only affects size of the address translated by the stage2 level (guest physical to host physical address translations).

### 4.3 KVM\_GET\_MSR\_INDEX\_LIST, KVM\_GET\_MSR\_FEATURE\_INDEX\_LIST

#### Capability

basic, `KVM_CAP_GET_MSR_FEATURES` for  
`KVM_GET_MSR_FEATURE_INDEX_LIST`

#### Architectures

x86

#### Type

system `ioctl`



**Parameters**

struct kvm\_msr\_list (in/out)

**Returns**

0 on success; -1 on error

Errors:

EFAULT	the msr index list cannot be read from or written to
E2BIG	the msr index list is to be to fit in the array specified by the user.

```
struct kvm_msr_list {
    __u32 nmsrs; /* number of msrs in entries */
    __u32 indices[0];
};
```

The user fills in the size of the indices array in nmsrs, and in return kvm adjusts nmsrs to reflect the actual number of msrs and fills in the indices array with their numbers.

KVM\_GET\_MSR\_INDEX\_LIST returns the guest msrs that are supported. The list varies by kvm version and host processor, but does not change otherwise.

Note: if kvm indicates supports MCE (KVM\_CAP\_MCE), then the MCE bank MSRs are not returned in the MSR list, as different vcpus can have a different number of banks, as set via the KVM\_X86\_SETUP\_MCE ioctl.

KVM\_GET\_MSR\_FEATURE\_INDEX\_LIST returns the list of MSRs that can be passed to the KVM\_GET\_MSRS system ioctl. This lets userspace probe host capabilities and processor features that are exposed via MSRs (e.g., VMX capabilities). This list also varies by kvm version and host processor, but does not change otherwise.

**4.4 KVM\_CHECK\_EXTENSION****Capability**

basic, KVM\_CAP\_CHECK\_EXTENSION\_VM for vm ioctl

**Architectures**

all

**Type**

system ioctl, vm ioctl

**Parameters**

extension identifier (KVM\_CAP\_\*)

**Returns**

0 if unsupported; 1 (or some other positive integer) if supported

The API allows the application to query about extensions to the core kvm API. Userspace passes an extension identifier (an integer) and receives an integer that describes the extension availability. Generally 0 means no and 1 means yes, but some extensions may report additional information in the integer return value.

Based on their initialization different VMs may have different capabilities. It is thus encouraged to use the `vm ioctl` to query for capabilities (available with `KVM_CAP_CHECK_EXTENSION_VM` on the `vm fd`)

### 4.5 KVM\_GET\_VCPU\_MMAP\_SIZE

**Capability**

basic

**Architectures**

all

**Type**

system ioctl

**Parameters**

none

**Returns**

size of vcpu mmap area, in bytes

The `KVM_RUN` ioctl (cf.) communicates with userspace via a shared memory region. This ioctl returns the size of that region. See the `KVM_RUN` documentation for details.

### 4.6 KVM\_SET\_MEMORY\_REGION

**Capability**

basic

**Architectures**

all

**Type**

vm ioctl

**Parameters**

struct `kvm_memory_region` (in)

**Returns**

0 on success, -1 on error

This ioctl is obsolete and has been removed.

### 4.7 KVM\_CREATE\_VCPU

**Capability**

basic

**Architectures**

all

**Type**

vm ioctl

**Parameters**

vcpu id (apic id on x86)

**Returns**

vcpu fd on success, -1 on error

This API adds a vcpu to a virtual machine. No more than `max_vcpus` may be added. The vcpu id is an integer in the range `[0, max_vcpu_id)`.

The recommended `max_vcpus` value can be retrieved using the `KVM_CAP_NR_VCPUS` of the `KVM_CHECK_EXTENSION` `ioctl()` at run-time. The maximum possible value for `max_vcpus` can be retrieved using the `KVM_CAP_MAX_VCPUS` of the `KVM_CHECK_EXTENSION` `ioctl()` at run-time.

If the `KVM_CAP_NR_VCPUS` does not exist, you should assume that `max_vcpus` is 4 cpus max. If the `KVM_CAP_MAX_VCPUS` does not exist, you should assume that `max_vcpus` is same as the value returned from `KVM_CAP_NR_VCPUS`.

The maximum possible value for `max_vcpu_id` can be retrieved using the `KVM_CAP_MAX_VCPU_ID` of the `KVM_CHECK_EXTENSION` `ioctl()` at run-time.

If the `KVM_CAP_MAX_VCPU_ID` does not exist, you should assume that `max_vcpu_id` is the same as the value returned from `KVM_CAP_MAX_VCPUS`.

On powerpc using book3s\_hv mode, the vcpus are mapped onto virtual threads in one or more virtual CPU cores. (This is because the hardware requires all the hardware threads in a CPU core to be in the same partition.) The `KVM_CAP_PPC_SMT` capability indicates the number of vcpus per virtual core (vcore). The vcore id is obtained by dividing the vcpu id by the number of vcpus per vcore. The vcpus in a given vcore will always be in the same physical core as each other (though that might be a different physical core from time to time). Userspace can control the threading (SMT) mode of the guest by its allocation of vcpu ids. For example, if userspace wants single-threaded guest vcpus, it should make all vcpu ids be a multiple of the number of vcpus per vcore.

For virtual cpus that have been created with S390 user controlled virtual machines, the resulting vcpu fd can be memory mapped at page offset `KVM_S390_SIE_PAGE_OFFSET` in order to obtain a memory map of the virtual cpu's hardware control block.

## 4.8 KVM\_GET\_DIRTY\_LOG (vm ioctl)

**Capability**

basic

**Architectures**

all

**Type**

vm ioctl

**Parameters**

struct `kvm_dirty_log` (in/out)

**Returns**

0 on success, -1 on error

```
/* for KVM_GET_DIRTY_LOG */
struct kvm_dirty_log {
    __u32 slot;
    __u32 padding;
    union {
        void __user *dirty_bitmap; /* one bit per page */
        __u64 padding;
    };
};
```

Given a memory slot, return a bitmap containing any pages dirtied since the last call to this ioctl. Bit 0 is the first page in the memory slot. Ensure the entire structure is cleared to avoid padding issues.

If KVM\_CAP\_MULTI\_ADDRESS\_SPACE is available, bits 16-31 specifies the address space for which you want to return the dirty bitmap. They must be less than the value that KVM\_CHECK\_EXTENSION returns for the KVM\_CAP\_MULTI\_ADDRESS\_SPACE capability.

The bits in the dirty bitmap are cleared before the ioctl returns, unless KVM\_CAP\_MANUAL\_DIRTY\_LOG\_PROTECT2 is enabled. For more information, see the description of the capability.

## 4.9 KVM\_SET\_MEMORY\_ALIAS

### Capability

basic

### Architectures

x86

### Type

vm ioctl

### Parameters

struct kvm\_memory\_alias (in)

### Returns

0 (success), -1 (error)

This ioctl is obsolete and has been removed.

## 4.10 KVM\_RUN

### Capability

basic

### Architectures

all

### Type

vcpu ioctl

**Parameters**

none

**Returns**

0 on success, -1 on error

Errors:

EINTR    an unmasked signal is pending
--

This ioctl is used to run a guest virtual cpu. While there are no explicit parameters, there is an implicit parameter block that can be obtained by mmap()ing the vcpu fd at offset 0, with the size given by KVM\_GET\_VCPU\_MMAP\_SIZE. The parameter block is formatted as a 'struct kvm\_run' (see below).

**4.11 KVM\_GET\_REGS****Capability**

basic

**Architectures**

all except ARM, arm64

**Type**

vcpu ioctl

**Parameters**

struct kvm\_regs (out)

**Returns**

0 on success, -1 on error

Reads the general purpose registers from the vcpu.

```
/* x86 */
struct kvm_regs {
    /* out (KVM_GET_REGS) / in (KVM_SET_REGS) */
    __u64 rax, rbx, rcx, rdx;
    __u64 rsi, rdi, rsp, rbp;
    __u64 r8,  r9,  r10, r11;
    __u64 r12, r13, r14, r15;
    __u64 rip, rflags;
};

/* mips */
struct kvm_regs {
    /* out (KVM_GET_REGS) / in (KVM_SET_REGS) */
    __u64 gpr[32];
    __u64 hi;
    __u64 lo;
    __u64 pc;
};
```

## 4.12 KVM\_SET\_REGS

**Capability**

basic

**Architectures**

all except ARM, arm64

**Type**

vcpu ioctl

**Parameters**

struct kvm\_regs (in)

**Returns**

0 on success, -1 on error

Writes the general purpose registers into the vcpu.

See KVM\_GET\_REGS for the data structure.

## 4.13 KVM\_GET\_SREGS

**Capability**

basic

**Architectures**

x86, ppc

**Type**

vcpu ioctl

**Parameters**

struct kvm\_sregs (out)

**Returns**

0 on success, -1 on error

Reads special registers from the vcpu.

```
/* x86 */
struct kvm_sregs {
    struct kvm_segment cs, ds, es, fs, gs, ss;
    struct kvm_segment tr, ldt;
    struct kvm_dtable gdt, idt;
    __u64 cr0, cr2, cr3, cr4, cr8;
    __u64 efer;
    __u64 apic_base;
    __u64 interrupt_bitmap[(KVM_NR_INTERRUPTS + 63) / 64];
};

/* ppc -- see arch/powerpc/include/uapi/asm/kvm.h */
```

interrupt\_bitmap is a bitmap of pending external interrupts. At most one bit may be set. This interrupt has been acknowledged by the APIC but not yet injected into the cpu core.

#### 4.14 KVM\_SET\_SREGS

**Capability**

basic

**Architectures**

x86, ppc

**Type**

vcpu ioctl

**Parameters**

struct kvm\_sregs (in)

**Returns**

0 on success, -1 on error

Writes special registers into the vcpu. See KVM\_GET\_SREGS for the data structures.

#### 4.15 KVM\_TRANSLATE

**Capability**

basic

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_translation (in/out)

**Returns**

0 on success, -1 on error

Translates a virtual address according to the vcpu's current address translation mode.

```
struct kvm_translation {
    /* in */
    __u64 linear_address;

    /* out */
    __u64 physical_address;
    __u8  valid;
    __u8  writeable;
    __u8  usermode;
    __u8  pad[5];
};
```

## 4.16 KVM\_INTERRUPT

**Capability**

basic

**Architectures**

x86, ppc, mips

**Type**

vcpu ioctl

**Parameters**

struct kvm\_interrupt (in)

**Returns**

0 on success, negative on failure.

Queues a hardware interrupt vector to be injected.

```
/* for KVM_INTERRUPT */
struct kvm_interrupt {
    /* in */
    __u32 irq;
};
```

**X86:****Returns**

0	on success,
-EEXIST	if an interrupt is already enqueued
-EINVAL	the irq number is invalid
-ENXIO	if the PIC is in the kernel
-EFAULT	if the pointer is invalid

Note ‘irq’ is an interrupt vector, not an interrupt pin or line. This ioctl is useful if the in-kernel PIC is not used.

**PPC:**

Queues an external interrupt to be injected. This ioctl is overloaded with 3 different irq values:

## a) KVM\_INTERRUPT\_SET

This injects an edge type external interrupt into the guest once it’s ready to receive interrupts. When injected, the interrupt is done.

## b) KVM\_INTERRUPT\_UNSET

This unsets any pending interrupt.

Only available with KVM\_CAP\_PPC\_UNSET\_IRQ.



c) `KVM_INTERRUPT_SET_LEVEL`

This injects a level type external interrupt into the guest context. The interrupt stays pending until a specific ioctl with `KVM_INTERRUPT_UNSET` is triggered.

Only available with `KVM_CAP_PPC_IRQ_LEVEL`.

Note that any value for ‘irq’ other than the ones stated above is invalid and incurs unexpected behavior.

This is an asynchronous vcpu ioctl and can be invoked from any thread.

## **MIPS:**

Queues an external interrupt to be injected into the virtual CPU. A negative interrupt number dequeues the interrupt.

This is an asynchronous vcpu ioctl and can be invoked from any thread.

### **4.17 KVM\_DEBUG\_GUEST**

**Capability**

basic

**Architectures**

none

**Type**

vcpu ioctl

**Parameters**

none)

**Returns**

-1 on error

Support for this has been removed. Use `KVM_SET_GUEST_DEBUG` instead.

### **4.18 KVM\_GET\_MSRS**

**Capability**

basic (vcpu), `KVM_CAP_GET_MSR_FEATURES` (system)

**Architectures**

x86

**Type**

system ioctl, vcpu ioctl

**Parameters**

struct `kvm_msrs` (in/out)

**Returns**

number of msrs successfully returned; -1 on error

When used as a system ioctl: Reads the values of MSR-based features that are available for the VM. This is similar to KVM\_GET\_SUPPORTED\_CPUID, but it returns MSR indices and values. The list of msr-based features can be obtained using KVM\_GET\_MSR\_FEATURE\_INDEX\_LIST in a system ioctl.

When used as a vcpu ioctl: Reads model-specific registers from the vcpu. Supported msr indices can be obtained using KVM\_GET\_MSR\_INDEX\_LIST in a system ioctl.

```
struct kvm_msrs {
    __u32 nmsrs; /* number of msrs in entries */
    __u32 pad;

    struct kvm_msr_entry entries[0];
};

struct kvm_msr_entry {
    __u32 index;
    __u32 reserved;
    __u64 data;
};
```

Application code should set the ‘nmsrs’ member (which indicates the size of the entries array) and the ‘index’ member of each array entry. kvm will fill in the ‘data’ member.

## 4.19 KVM\_SET\_MSRS

### Capability

basic

### Architectures

x86

### Type

vcpu ioctl

### Parameters

struct kvm\_msrs (in)

### Returns

number of msrs successfully set (see below), -1 on error

Writes model-specific registers to the vcpu. See KVM\_GET\_MSRS for the data structures.

Application code should set the ‘nmsrs’ member (which indicates the size of the entries array), and the ‘index’ and ‘data’ members of each array entry.

It tries to set the MSRs in array entries[] one by one. If setting an MSR fails, e.g., due to setting reserved bits, the MSR isn’t supported/emulated by KVM, etc..., it stops processing the MSR list and returns the number of MSRs that have been set successfully.

## 4.20 KVM\_SET\_CPUID

**Capability**

basic

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_cpuid (in)

**Returns**

0 on success, -1 on error

Defines the vcpu responses to the cpuid instruction. Applications should use the KVM\_SET\_CPUID2 ioctl if available.

Note, when this IOCTL fails, KVM gives no guarantees that previous valid CPUID configuration (if there is) is not corrupted. Userspace can get a copy of the resulting CPUID configuration through KVM\_GET\_CPUID2 in case.

```
struct kvm_cpuid_entry {
    __u32 function;
    __u32 eax;
    __u32 ebx;
    __u32 ecx;
    __u32 edx;
    __u32 padding;
};

/* for KVM_SET_CPUID */
struct kvm_cpuid {
    __u32 nent;
    __u32 padding;
    struct kvm_cpuid_entry entries[0];
};
```

## 4.21 KVM\_SET\_SIGNAL\_MASK

**Capability**

basic

**Architectures**

all

**Type**

vcpu ioctl

**Parameters**

struct kvm\_signal\_mask (in)

**Returns**

0 on success, -1 on error

Defines which signals are blocked during execution of KVM\_RUN. This signal mask temporarily overrides the threads signal mask. Any unblocked signal received (except SIGKILL and SIGSTOP, which retain their traditional behaviour) will cause KVM\_RUN to return with -EINTR.

Note the signal will only be delivered if not blocked by the original signal mask.

```
/* for KVM_SET_SIGNAL_MASK */
struct kvm_signal_mask {
    __u32 len;
    __u8  sigset[0];
};
```

## 4.22 KVM\_GET\_FPU

**Capability**

basic

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_fpu (out)

**Returns**

0 on success, -1 on error

Reads the floating point state from the vcpu.

```
/* for KVM_GET_FPU and KVM_SET_FPU */
struct kvm_fpu {
    __u8  fpr[8][16];
    __u16 fcw;
    __u16 fsw;
    __u8  ftwx; /* in fxsave format */
    __u8  pad1;
    __u16 last_opcode;
    __u64 last_ip;
    __u64 last_dp;
    __u8  xmm[16][16];
    __u32 mxcsr;
    __u32 pad2;
};
```

## 4.23 KVM\_SET\_FPU

### Capability

basic

### Architectures

x86

### Type

vcpu ioctl

### Parameters

struct kvm\_fpu (in)

### Returns

0 on success, -1 on error

Writes the floating point state to the vcpu.

```
/* for KVM_GET_FPU and KVM_SET_FPU */
struct kvm_fpu {
    __u8  fpr[8][16];
    __u16 fcw;
    __u16 fsw;
    __u8  ftwx; /* in fxsave format */
    __u8  pad1;
    __u16 last_opcode;
    __u64 last_ip;
    __u64 last_dp;
    __u8  xmm[16][16];
    __u32 mxcsr;
    __u32 pad2;
};
```

## 4.24 KVM\_CREATE\_IRQCHIP

### Capability

KVM\_CAP\_IRQCHIP, KVM\_CAP\_S390\_IRQCHIP (s390)

### Architectures

x86, ARM, arm64, s390

### Type

vm ioctl

### Parameters

none

### Returns

0 on success, -1 on error

Creates an interrupt controller model in the kernel. On x86, creates a virtual ioapic, a virtual PIC (two PICs, nested), and sets up future vcpus to have a local APIC. IRQ routing for GSIs 0-15 is set to both PIC and IOAPIC; GSI 16-23 only go to the IOAPIC. On ARM/arm64, a GICv2 is created. Any other GIC versions require

the usage of `KVM_CREATE_DEVICE`, which also supports creating a GICv2. Using `KVM_CREATE_DEVICE` is preferred over `KVM_CREATE_IRQCHIP` for GICv2. On s390, a dummy irq routing table is created.

Note that on s390 the `KVM_CAP_S390_IRQCHIP` vm capability needs to be enabled before `KVM_CREATE_IRQCHIP` can be used.

## 4.25 KVM\_IRQ\_LINE

### Capability

`KVM_CAP_IRQCHIP`

### Architectures

x86, arm, arm64

### Type

vm ioctl

### Parameters

`struct kvm_irq_level`

### Returns

0 on success, -1 on error

Sets the level of a GSI input to the interrupt controller model in the kernel. On some architectures it is required that an interrupt controller model has been previously created with `KVM_CREATE_IRQCHIP`. Note that edge-triggered interrupts require the level to be set to 1 and then back to 0.

On real hardware, interrupt pins can be active-low or active-high. This does not matter for the level field of `struct kvm_irq_level`: 1 always means active (asserted), 0 means inactive (deasserted).

x86 allows the operating system to program the interrupt polarity (active-low/active-high) for level-triggered interrupts, and KVM used to consider the polarity. However, due to bitrot in the handling of active-low interrupts, the above convention is now valid on x86 too. This is signaled by `KVM_CAP_X86_IOAPIC_POLARITY_IGNORED`. Userspace should not present interrupts to the guest as active-low unless this capability is present (or unless it is not using the in-kernel irqchip, of course).

ARM/arm64 can signal an interrupt either at the CPU level, or at the in-kernel irqchip (GIC), and for in-kernel irqchip can tell the GIC to use PPIs designated for specific cpus. The `irq` field is interpreted like this:

bits:	31 ... 28   27 ... 24   23 ... 16   15 ... 0
field:	vcpu2_index   irq_type   vcpu_index   irq_id

The `irq_type` field has the following values:

- **`irq_type[0]`:**  
out-of-kernel GIC: `irq_id` 0 is IRQ, `irq_id` 1 is FIQ
- **`irq_type[1]`:**  
in-kernel GIC: SPI, `irq_id` between 32 and 1019 (incl.) (the `vcpu_index` field is ignored)

- **irq\_type[2]:**

in-kernel GIC: PPI, irq\_id between 16 and 31 (incl.)

(The irq\_id field thus corresponds nicely to the IRQ ID in the ARM GIC specs)

In both cases, level is used to assert/deassert the line.

When KVM\_CAP\_ARM\_IRQ\_LINE\_LAYOUT\_2 is supported, the target vcpu is identified as (256 \* vcpu2\_index + vcpu\_index). Otherwise, vcpu2\_index must be zero.

Note that on arm/arm64, the KVM\_CAP\_IRQCHIP capability only conditions injection of interrupts for the in-kernel irqchip. KVM\_IRQ\_LINE can always be used for a userspace interrupt controller.

```
struct kvm_irq_level {
    union {
        __u32 irq;      /* GSI */
        __s32 status;   /* not used for KVM_IRQ_LEVEL */
    };
    __u32 level;        /* 0 or 1 */
};
```

## 4.26 KVM\_GET\_IRQCHIP

### Capability

KVM\_CAP\_IRQCHIP

### Architectures

x86

### Type

vm ioctl

### Parameters

struct kvm\_irqchip (in/out)

### Returns

0 on success, -1 on error

Reads the state of a kernel interrupt controller created with KVM\_CREATE\_IRQCHIP into a buffer provided by the caller.

```
struct kvm_irqchip {
    __u32 chip_id; /* 0 = PIC1, 1 = PIC2, 2 = IOAPIC */
    __u32 pad;
    union {
        char dummy[512]; /* reserving space */
        struct kvm_pic_state pic;
        struct kvm_ioapic_state ioapic;
    } chip;
};
```

## 4.27 KVM\_SET\_IRQCHIP

**Capability**

KVM\_CAP\_IRQCHIP

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct kvm\_irqchip (in)

**Returns**

0 on success, -1 on error

Sets the state of a kernel interrupt controller created with KVM\_CREATE\_IRQCHIP from a buffer provided by the caller.

```
struct kvm_irqchip {
    __u32 chip_id; /* 0 = PIC1, 1 = PIC2, 2 = IOAPIC */
    __u32 pad;
    union {
        char dummy[512]; /* reserving space */
        struct kvm_pic_state pic;
        struct kvm_ioapic_state ioapic;
    } chip;
};
```

## 4.28 KVM\_XEN\_HVM\_CONFIG

**Capability**

KVM\_CAP\_XEN\_HVM

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct kvm\_xen\_hvm\_config (in)

**Returns**

0 on success, -1 on error

Sets the MSR that the Xen HVM guest uses to initialize its hypercall page, and provides the starting address and size of the hypercall blobs in userspace. When the guest writes the MSR, kvm copies one page of a blob (32- or 64-bit, depending on the vcpu mode) to guest memory.

```
struct kvm_xen_hvm_config {
    __u32 flags;
    __u32 msr;
```

(continues on next page)



(continued from previous page)

```
__u64 blob_addr_32;  
__u64 blob_addr_64;  
__u8 blob_size_32;  
__u8 blob_size_64;  
__u8 pad2[30];  
};
```

## 4.29 KVM\_GET\_CLOCK

### Capability

KVM\_CAP\_ADJUST\_CLOCK

### Architectures

x86

### Type

vm ioctl

### Parameters

struct kvm\_clock\_data (out)

### Returns

0 on success, -1 on error

Gets the current timestamp of kvmclock as seen by the current guest. In conjunction with KVM\_SET\_CLOCK, it is used to ensure monotonicity on scenarios such as migration.

When KVM\_CAP\_ADJUST\_CLOCK is passed to KVM\_CHECK\_EXTENSION, it returns the set of bits that KVM can return in struct kvm\_clock\_data's flag member.

The only flag defined now is KVM\_CLOCK\_TSC\_STABLE. If set, the returned value is the exact kvmclock value seen by all VCPUs at the instant when KVM\_GET\_CLOCK was called. If clear, the returned value is simply CLOCK\_MONOTONIC plus a constant offset; the offset can be modified with KVM\_SET\_CLOCK. KVM will try to make all VCPUs follow this clock, but the exact value read by each VCPU could differ, because the host TSC is not stable.

```
struct kvm_clock_data {  
    __u64 clock; /* kvmclock current value */  
    __u32 flags;  
    __u32 pad[9];  
};
```

### 4.30 KVM\_SET\_CLOCK

**Capability**

KVM\_CAP\_ADJUST\_CLOCK

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct kvm\_clock\_data (in)

**Returns**

0 on success, -1 on error

Sets the current timestamp of kvmclock to the value specified in its parameter. In conjunction with KVM\_GET\_CLOCK, it is used to ensure monotonicity on scenarios such as migration.

```
struct kvm_clock_data {
    __u64 clock; /* kvmclock current value */
    __u32 flags;
    __u32 pad[9];
};
```

### 4.31 KVM\_GET\_VCPU\_EVENTS

**Capability**

KVM\_CAP\_VCPU\_EVENTS

**Extended by**

KVM\_CAP\_INTR\_SHADOW

**Architectures**

x86, arm, arm64

**Type**

vcpu ioctl

**Parameters**

struct kvm\_vcpu\_event (out)

**Returns**

0 on success, -1 on error

**X86:**

Gets currently pending exceptions, interrupts, and NMIs as well as related states of the vcpu.

```
struct kvm_vcpu_events {
    struct {
        __u8 injected;
        __u8 nr;
        __u8 has_error_code;
        __u8 pending;
        __u32 error_code;
    } exception;
    struct {
        __u8 injected;
        __u8 nr;
        __u8 soft;
        __u8 shadow;
    } interrupt;
    struct {
        __u8 injected;
        __u8 pending;
        __u8 masked;
        __u8 pad;
    } nmi;
    __u32 sipi_vector;
    __u32 flags;
    struct {
        __u8 smm;
        __u8 pending;
        __u8 smm_inside_nmi;
        __u8 latched_init;
    } smi;
    __u8 reserved[27];
    __u8 exception_has_payload;
    __u64 exception_payload;
};
```

The following bits are defined in the flags field:

- `KVM_VCPUEVENT_VALID_SHADOW` may be set to signal that `interrupt.shadow` contains a valid state.
- `KVM_VCPUEVENT_VALID_SMM` may be set to signal that `smi` contains a valid state.
- `KVM_VCPUEVENT_VALID_PAYLOAD` may be set to signal that the `exception_has_payload`, `exception_payload`, and `exception.pending` fields contain a valid state. This bit will be set whenever `KVM_CAP_EXCEPTION_PAYLOAD` is enabled.

**ARM/ARM64:**

If the guest accesses a device that is being emulated by the host kernel in such a way that a real device would generate a physical SError, KVM may make a virtual SError pending for that VCPU. This system error interrupt remains pending until the guest takes the exception by unmasking PSTATE.A.

Running the VCPU may cause it to take a pending SError, or make an access that causes an SError to become pending. The event's description is only valid while the VCPU is not running.

This API provides a way to read and write the pending 'event' state that is not visible to the guest. To save, restore or migrate a VCPU the struct representing the state can be read then written using this GET/SET API, along with the other guest-visible registers. It is not possible to 'cancel' an SError that has been made pending.

A device being emulated in user-space may also wish to generate an SError. To do this the events structure can be populated by user-space. The current state should be read first, to ensure no existing SError is pending. If an existing SError is pending, the architecture's 'Multiple SError interrupts' rules should be followed. (2.5.3 of DDI0587.a "ARM Reliability, Availability, and Serviceability (RAS) Specification").

SError exceptions always have an ESR value. Some CPUs have the ability to specify what the virtual SError's ESR value should be. These systems will advertise KVM\_CAP\_ARM\_INJECT\_SERROR\_ESR. In this case exception.has\_esr will always have a non-zero value when read, and the agent making an SError pending should specify the ISS field in the lower 24 bits of exception.error\_esr. If the system supports KVM\_CAP\_ARM\_INJECT\_SERROR\_ESR, but user-space sets the events with exception.has\_esr as zero, KVM will choose an ESR.

Specifying exception.has\_esr on a system that does not support it will return -EINVAL. Setting anything other than the lower 24bits of exception.error\_esr will return -EINVAL.

It is not possible to read back a pending external abort (injected via KVM\_SET\_VCPU\_EVENTS or otherwise) because such an exception is always delivered directly to the virtual CPU).

```
struct kvm_vcpu_events {
    struct {
        __u8 serror_pending;
        __u8 serror_has_esr;
        __u8 ext_dabt_pending;
        /* Align it to 8 bytes */
        __u8 pad[5];
        __u64 serror_esr;
    } exception;
    __u32 reserved[12];
};
```

## 4.32 KVM\_SET\_VCPU\_EVENTS

### Capability

KVM\_CAP\_VCPU\_EVENTS

### Extended by

KVM\_CAP\_INTR\_SHADOW

### Architectures

x86, arm, arm64

### Type

vcpu ioctl

### Parameters

struct kvm\_vcpu\_event (in)

### Returns

0 on success, -1 on error

### X86:

Set pending exceptions, interrupts, and NMIs as well as related states of the vcpu.

See KVM\_GET\_VCPU\_EVENTS for the data structure.

Fields that may be modified asynchronously by running VCPUs can be excluded from the update. These fields are nmi.pending, sipi\_vector, smi.smm, smi.pending. Keep the corresponding bits in the flags field cleared to suppress overwriting the current in-kernel state. The bits are:

KVM_VCPUEVENT_VALID_NMI_PENDING	transfer nmi.pending to the kernel
KVM_VCPUEVENT_VALID_SIPI_VECTOR	transfer sipi_vector
KVM_VCPUEVENT_VALID_SMM	transfer the smi sub-struct.

If KVM\_CAP\_INTR\_SHADOW is available, KVM\_VCPUEVENT\_VALID\_SHADOW can be set in the flags field to signal that interrupt.shadow contains a valid state and shall be written into the VCPU.

KVM\_VCPUEVENT\_VALID\_SMM can only be set if KVM\_CAP\_X86\_SMM is available.

If KVM\_CAP\_EXCEPTION\_PAYLOAD is enabled, KVM\_VCPUEVENT\_VALID\_PAYLOAD can be set in the flags field to signal that the exception\_has\_payload, exception\_payload, and exception.pending fields contain a valid state and shall be written into the VCPU.

**ARM/ARM64:**

User space may need to inject several types of events to the guest.

Set the pending SError exception state for this VCPU. It is not possible to ‘cancel’ an SError that has been made pending.

If the guest performed an access to I/O memory which could not be handled by userspace, for example because of missing instruction syndrome decode information or because there is no device mapped at the accessed IPA, then userspace can ask the kernel to inject an external abort using the address from the exiting fault on the VCPU. It is a programming error to set `ext_dabt_pending` after an exit which was not either `KVM_EXIT_MMIO` or `KVM_EXIT_ARM_NISV`. This feature is only available if the system supports `KVM_CAP_ARM_INJECT_EXT_DABT`. This is a helper which provides commonality in how userspace reports accesses for the above cases to guests, across different userspace implementations. Nevertheless, userspace can still emulate all Arm exceptions by manipulating individual registers using the `KVM_SET_ONE_REG` API.

See `KVM_GET_VCPU_EVENTS` for the data structure.

**4.33 KVM\_GET\_DEBUGREGS****Capability**

`KVM_CAP_DEBUGREGS`

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct `kvm_debugregs` (out)

**Returns**

0 on success, -1 on error

Reads debug registers from the vcpu.

```
struct kvm_debugregs {
    __u64 db[4];
    __u64 dr6;
    __u64 dr7;
    __u64 flags;
    __u64 reserved[9];
};
```

#### 4.34 KVM\_SET\_DEBUGREGS

**Capability**

KVM\_CAP\_DEBUGREGS

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct kvm\_debugregs (in)

**Returns**

0 on success, -1 on error

Writes debug registers into the vcpu.

See KVM\_GET\_DEBUGREGS for the data structure. The flags field is unused yet and must be cleared on entry.

#### 4.35 KVM\_SET\_USER\_MEMORY\_REGION

**Capability**

KVM\_CAP\_USER\_MEMORY

**Architectures**

all

**Type**

vm ioctl

**Parameters**

struct kvm\_userspace\_memory\_region (in)

**Returns**

0 on success, -1 on error

```
struct kvm_userspace_memory_region {
    __u32 slot;
    __u32 flags;
    __u64 guest_phys_addr;
    __u64 memory_size; /* bytes */
    __u64 userspace_addr; /* start of the userspace allocated
    ↪memory */
};

/* for kvm_memory_region::flags */
#define KVM_MEM_LOG_DIRTY_PAGES    (1UL << 0)
#define KVM_MEM_READONLY          (1UL << 1)
```

This ioctl allows the user to create, modify or delete a guest physical memory slot. Bits 0-15 of “slot” specify the slot id and this value should be less than the maximum number of user memory slots supported per VM. The maximum allowed

slots can be queried using `KVM_CAP_NR_MEMSLOTS`. Slots may not overlap in guest physical address space.

If `KVM_CAP_MULTI_ADDRESS_SPACE` is available, bits 16-31 of “slot” specifies the address space which is being modified. They must be less than the value that `KVM_CHECK_EXTENSION` returns for the `KVM_CAP_MULTI_ADDRESS_SPACE` capability. Slots in separate address spaces are unrelated; the restriction on overlapping slots only applies within each address space.

Deleting a slot is done by passing zero for `memory_size`. When changing an existing slot, it may be moved in the guest physical memory space, or its flags may be modified, but it may not be resized.

Memory for the region is taken starting at the address denoted by the field `userspace_addr`, which must point at user addressable memory for the entire memory slot size. Any object may back this memory, including anonymous memory, ordinary files, and hugetlbfs.

On architectures that support a form of address tagging, `userspace_addr` must be an untagged address.

It is recommended that the lower 21 bits of `guest_phys_addr` and `userspace_addr` be identical. This allows large pages in the guest to be backed by large pages in the host.

The `flags` field supports two flags: `KVM_MEM_LOG_DIRTY_PAGES` and `KVM_MEM_READONLY`. The former can be set to instruct KVM to keep track of writes to memory within the slot. See `KVM_GET_DIRTY_LOG` ioctl to know how to use it. The latter can be set, if `KVM_CAP_READONLY_MEM` capability allows it, to make a new slot read-only. In this case, writes to this memory will be posted to userspace as `KVM_EXIT_MMIO` exits.

When the `KVM_CAP_SYNC_MMU` capability is available, changes in the backing of the memory region are automatically reflected into the guest. For example, an `mmap()` that affects the region will be made visible immediately. Another example is `madvise(MADV_DROP)`.

It is recommended to use this API instead of the `KVM_SET_MEMORY_REGION` ioctl. The `KVM_SET_MEMORY_REGION` does not allow fine grained control over memory allocation and is deprecated.

### **4.36 KVM\_SET\_TSS\_ADDR**

#### **Capability**

`KVM_CAP_SET_TSS_ADDR`

#### **Architectures**

x86

#### **Type**

vm ioctl

#### **Parameters**

unsigned long `tss_address` (in)

#### **Returns**

0 on success, -1 on error



This ioctl defines the physical address of a three-page region in the guest physical address space. The region must be within the first 4GB of the guest physical address space and must not conflict with any memory slot or any mmio address. The guest may malfunction if it accesses this memory region.

This ioctl is required on Intel-based hosts. This is needed on Intel hardware because of a quirk in the virtualization implementation (see the internals documentation when it pops into existence).

#### 4.37 KVM\_ENABLE\_CAP

**Capability**

KVM\_CAP\_ENABLE\_CAP

**Architectures**

mips, ppc, s390

**Type**

vcpu ioctl

**Parameters**

struct kvm\_enable\_cap (in)

**Returns**

0 on success; -1 on error

**Capability**

KVM\_CAP\_ENABLE\_CAP\_VM

**Architectures**

all

**Type**

vm ioctl

**Parameters**

struct kvm\_enable\_cap (in)

**Returns**

0 on success; -1 on error

---

**Note:** Not all extensions are enabled by default. Using this ioctl the application can enable an extension, making it available to the guest.

---

On systems that do not support this ioctl, it always fails. On systems that do support it, it only works for extensions that are supported for enablement.

To check if a capability can be enabled, the KVM\_CHECK\_EXTENSION ioctl should be used.

```
struct kvm_enable_cap {
    /* in */
    __u32 cap;
```

The capability that is supposed to get enabled.

```
__u32 flags;
```

A bitfield indicating future enhancements. Has to be 0 for now.

```
__u64 args[4];
```

Arguments for enabling a feature. If a feature needs initial values to function properly, this is the place to put them.

```
    __u8  pad[64];  
};
```

The vcpu ioctl should be used for vcpu-specific capabilities, the vm ioctl for vm-wide capabilities.

### 4.38 KVM\_GET\_MP\_STATE

#### Capability

KVM\_CAP\_MP\_STATE

#### Architectures

x86, s390, arm, arm64

#### Type

vcpu ioctl

#### Parameters

struct kvm\_mp\_state (out)

#### Returns

0 on success; -1 on error

```
struct kvm_mp_state {  
    __u32 mp_state;  
};
```

Returns the vcpu's current "multiprocessing state" (though also valid on uniprocessor guests).

Possible values are:

KVM_MP_STATE_RUNNING	the vcpu is currently running [x86,arm/arm64]
KVM_MP_STATE_UNINITIALIZED	the vcpu is an application processor (AP) which has not yet received an INIT signal [x86]
KVM_MP_STATE_INIT_PENDING	the vcpu has received an INIT signal, and is now ready for a SIPI [x86]
KVM_MP_STATE_HALTED	the vcpu has executed a HLT instruction and is waiting for an interrupt [x86]
KVM_MP_STATE_SIPI_RECEIVED	the vcpu has just received a SIPI (vector accessible via KVM_GET_VCPU_EVENTS) [x86]
KVM_MP_STATE_STOPPED	the vcpu is stopped [s390,arm/arm64]
KVM_MP_STATE_ERROR	the vcpu is in a special error state [s390]
KVM_MP_STATE_OPERATING	the vcpu is operating (running or halted) [s390]
KVM_MP_STATE_LOAD_STARTUP	the vcpu is in a special load/startup state [s390]

On x86, this ioctl is only useful after KVM\_CREATE\_IRQCHIP. Without an in-kernel irqchip, the multiprocessing state must be maintained by userspace on these architectures.

#### For arm/arm64:

The only states that are valid are KVM\_MP\_STATE\_STOPPED and KVM\_MP\_STATE\_RUNNABLE which reflect if the vcpu is paused or not.

### 4.39 KVM\_SET\_MP\_STATE

#### Capability

KVM\_CAP\_MP\_STATE

#### Architectures

x86, s390, arm, arm64

#### Type

vcpu ioctl

#### Parameters

struct kvm\_mp\_state (in)

#### Returns

0 on success; -1 on error

Sets the vcpu's current "multiprocessing state" ; see KVM\_GET\_MP\_STATE for arguments.

On x86, this ioctl is only useful after KVM\_CREATE\_IRQCHIP. Without an in-kernel irqchip, the multiprocessing state must be maintained by userspace on these architectures.

### For arm/arm64:

The only states that are valid are `KVM_MP_STATE_STOPPED` and `KVM_MP_STATE_RUNNABLE` which reflect if the vcpu should be paused or not.

#### 4.40 KVM\_SET\_IDENTITY\_MAP\_ADDR

##### Capability

`KVM_CAP_SET_IDENTITY_MAP_ADDR`

##### Architectures

x86

##### Type

vm ioctl

##### Parameters

unsigned long identity (in)

##### Returns

0 on success, -1 on error

This ioctl defines the physical address of a one-page region in the guest physical address space. The region must be within the first 4GB of the guest physical address space and must not conflict with any memory slot or any mmio address. The guest may malfunction if it accesses this memory region.

Setting the address to 0 will result in resetting the address to its default (0xffffbc000).

This ioctl is required on Intel-based hosts. This is needed on Intel hardware because of a quirk in the virtualization implementation (see the internals documentation when it pops into existence).

Fails if any VCPU has already been created.

#### 4.41 KVM\_SET\_BOOT\_CPU\_ID

##### Capability

`KVM_CAP_SET_BOOT_CPU_ID`

##### Architectures

x86

##### Type

vm ioctl

##### Parameters

unsigned long vcpu\_id

##### Returns

0 on success, -1 on error

Define which vcpu is the Bootstrap Processor (BSP). Values are the same as the vcpu id in `KVM_CREATE_VCPU`. If this ioctl is not called, the default is vcpu 0.

#### 4.42 KVM\_GET\_XSAVE

**Capability**

KVM\_CAP\_XSAVE

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_xsaves (out)

**Returns**

0 on success, -1 on error

```
struct kvm_xsaves {
    __u32 region[1024];
};
```

This ioctl would copy current vcpu's xsaves struct to the userspace.

#### 4.43 KVM\_SET\_XSAVE

**Capability**

KVM\_CAP\_XSAVE

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_xsaves (in)

**Returns**

0 on success, -1 on error

```
struct kvm_xsaves {
    __u32 region[1024];
};
```

This ioctl would copy userspace's xsaves struct to the kernel.

#### 4.44 KVM\_GET\_XCRS

**Capability**

KVM\_CAP\_XCRS

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_xcrs (out)

**Returns**

0 on success, -1 on error

```
struct kvm_xcr {
    __u32 xcr;
    __u32 reserved;
    __u64 value;
};

struct kvm_xcrs {
    __u32 nr_xcrs;
    __u32 flags;
    struct kvm_xcr xcrs[KVM_MAX_XCRS];
    __u64 padding[16];
};
```

This ioctl would copy current vcpu' s xcrs to the userspace.

#### 4.45 KVM\_SET\_XCRS

**Capability**

KVM\_CAP\_XCRS

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_xcrs (in)

**Returns**

0 on success, -1 on error

```
struct kvm_xcr {
    __u32 xcr;
    __u32 reserved;
    __u64 value;
};
```

(continues on next page)

(continued from previous page)

```

struct kvm_xcrs {
    __u32 nr_xcrs;
    __u32 flags;
    struct kvm_xcr xcrs[KVM_MAX_XCRS];
    __u64 padding[16];
};

```

This ioctl would set vcpu' s xcr to the value userspace specified.

#### 4.46 KVM\_GET\_SUPPORTED\_CPUID

##### Capability

KVM\_CAP\_EXT\_CPUID

##### Architectures

x86

##### Type

system ioctl

##### Parameters

struct kvm\_cpuid2 (in/out)

##### Returns

0 on success, -1 on error

```

struct kvm_cpuid2 {
    __u32 nent;
    __u32 padding;
    struct kvm_cpuid_entry2 entries[0];
};

#define KVM_CPUID_FLAG_SIGNIFCANT_INDEX          BIT(0)
#define KVM_CPUID_FLAG_STATEFUL_FUNC            BIT(1) /* deprecated */
/*
#define KVM_CPUID_FLAG_STATE_READ_NEXT          BIT(2) /*
/* deprecated */

struct kvm_cpuid_entry2 {
    __u32 function;
    __u32 index;
    __u32 flags;
    __u32 eax;
    __u32 ebx;
    __u32 ecx;
    __u32 edx;
    __u32 padding[3];
};

```

This ioctl returns x86 cpuid features which are supported by both the hardware and kvm in its default configuration. Userspace can use the information returned

by this `ioctl` to construct `cpuid` information (for `KVM_SET_CPUID2`) that is consistent with hardware, kernel, and userspace capabilities, and with user requirements (for example, the user may wish to constrain `cpuid` to emulate older hardware, or for feature consistency across a cluster).

Note that certain capabilities, such as `KVM_CAP_X86_DISABLE_EXITS`, may expose `cpuid` features (e.g. `MONITOR`) which are not supported by `kvm` in its default configuration. If userspace enables such capabilities, it is responsible for modifying the results of this `ioctl` appropriately.

Userspace invokes `KVM_GET_SUPPORTED_CPUID` by passing a `kvm_cpuid2` structure with the `'nent'` field indicating the number of entries in the variable-size array `'entries'`. If the number of entries is too low to describe the `cpu` capabilities, an error (`E2BIG`) is returned. If the number is too high, the `'nent'` field is adjusted and an error (`ENOMEM`) is returned. If the number is just right, the `'nent'` field is adjusted to the number of valid entries in the `'entries'` array, which is then filled.

The entries returned are the host `cpuid` as returned by the `cpuid` instruction, with unknown or unsupported features masked out. Some features (for example, `x2apic`), may not be present in the host `cpu`, but are exposed by `kvm` if it can emulate them efficiently. The fields in each entry are defined as follows:

**function:**

the `eax` value used to obtain the entry

**index:**

the `ecx` value used to obtain the entry (for entries that are affected by `ecx`)

**flags:**

an OR of zero or more of the following:

**KVM\_CPUID\_FLAG\_SIGNIFCANT\_INDEX:**

if the `index` field is valid

**eax, ebx, ecx, edx:**

the values returned by the `cpuid` instruction for this function/index combination

The TSC deadline timer feature (`CPUID` leaf 1, `ecx[24]`) is always returned as false, since the feature depends on `KVM_CREATE_IRQCHIP` for local APIC support. Instead it is reported via:

`ioctl(KVM_CHECK_EXTENSION, KVM_CAP_TSC_DEADLINE_TIMER)`

if that returns true and you use `KVM_CREATE_IRQCHIP`, or if you emulate the feature in userspace, then you can enable the feature for `KVM_SET_CPUID2`.



## 4.47 KVM\_PPC\_GET\_PVINFO

### Capability

KVM\_CAP\_PPC\_GET\_PVINFO

### Architectures

ppc

### Type

vm ioctl

### Parameters

struct kvm\_ppc\_pvinfos (out)

### Returns

0 on success, !0 on error

```
struct kvm_ppc_pvinfos {
    __u32 flags;
    __u32 hcall[4];
    __u8  pad[108];
};
```

This ioctl fetches PV specific information that need to be passed to the guest using the device tree or other means from vm context.

The hcall array defines 4 instructions that make up a hypercall.

If any additional field gets added to this structure later on, a bit for that additional piece of information will be set in the flags bitmap.

The flags bitmap is defined as:

```
/* the host supports the ePAPR idle hcall
#define KVM_PPC_PVINFO_FLAGS_EV_IDLE    (1<<0)
```

## 4.52 KVM\_SET\_GSI\_ROUTING

### Capability

KVM\_CAP\_IRQ\_ROUTING

### Architectures

x86 s390 arm arm64

### Type

vm ioctl

### Parameters

struct kvm\_irq\_routing (in)

### Returns

0 on success, -1 on error

Sets the GSI routing table entries, overwriting any previously set entries.

On arm/arm64, GSI routing has the following limitation:

- GSI routing does not apply to KVM\_IRQ\_LINE but only to KVM\_IRQFD.

```
struct kvm_irq_routing {
    __u32 nr;
    __u32 flags;
    struct kvm_irq_routing_entry entries[0];
};
```

No flags are specified so far, the corresponding field must be set to zero.

```
struct kvm_irq_routing_entry {
    __u32 gsi;
    __u32 type;
    __u32 flags;
    __u32 pad;
    union {
        struct kvm_irq_routing_irqchip irqchip;
        struct kvm_irq_routing_msi msi;
        struct kvm_irq_routing_s390_adapter adapter;
        struct kvm_irq_routing_hv_sint hv_sint;
        __u32 pad[8];
    } u;
};

/* gsi routing entry types */
#define KVM_IRQ_ROUTING_IRQCHIP 1
#define KVM_IRQ_ROUTING_MSI 2
#define KVM_IRQ_ROUTING_S390_ADAPTER 3
#define KVM_IRQ_ROUTING_HV_SINT 4
```

flags:

- **KVM\_MSI\_VALID\_DEVID**: used along with **KVM\_IRQ\_ROUTING\_MSI** routing entry type, specifies that the **devid** field contains a valid value. The per-VM **KVM\_CAP\_MSI\_DEVID** capability advertises the requirement to provide the device ID. If this capability is not available, userspace should never set the **KVM\_MSI\_VALID\_DEVID** flag as the **ioctl** might fail.
- zero otherwise

```
struct kvm_irq_routing_irqchip {
    __u32 irqchip;
    __u32 pin;
};

struct kvm_irq_routing_msi {
    __u32 address_lo;
    __u32 address_hi;
    __u32 data;
    union {
        __u32 pad;
        __u32 devid;
    };
};
```

If `KVM_MSI_VALID_DEVID` is set, `devid` contains a unique device identifier for the device that wrote the MSI message. For PCI, this is usually a BFD identifier in the lower 16 bits.

On x86, `address_hi` is ignored unless the `KVM_X2APIC_API_USE_32BIT_IDS` feature of `KVM_CAP_X2APIC_API` capability is enabled. If it is enabled, `address_hi` bits 31-8 provide bits 31-8 of the destination id. Bits 7-0 of `address_hi` must be zero.

```
struct kvm_irq_routing_s390_adapter {
    __u64 ind_addr;
    __u64 summary_addr;
    __u64 ind_offset;
    __u32 summary_offset;
    __u32 adapter_id;
};

struct kvm_irq_routing_hv_sint {
    __u32 vcpu;
    __u32 sint;
};
```

## 4.55 KVM\_SET\_TSC\_KHZ

### Capability

`KVM_CAP_TSC_CONTROL`

### Architectures

x86

### Type

`vcpu ioctl`

### Parameters

`virtual tsc_khz`

### Returns

0 on success, -1 on error

Specifies the tsc frequency for the virtual machine. The unit of the frequency is KHz.

## 4.56 KVM\_GET\_TSC\_KHZ

### Capability

`KVM_CAP_GET_TSC_KHZ`

### Architectures

x86

### Type

`vcpu ioctl`

**Parameters**

none

**Returns**

virtual tsc-khz on success, negative value on error

Returns the tsc frequency of the guest. The unit of the return value is KHz. If the host has unstable tsc this ioctl returns -EIO instead as an error.

**4.57 KVM\_GET\_LAPIC****Capability**

KVM\_CAP\_IRQCHIP

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_lapic\_state (out)

**Returns**

0 on success, -1 on error

```
#define KVM_APIC_REG_SIZE 0x400
struct kvm_lapic_state {
    char regs[KVM_APIC_REG_SIZE];
};
```

Reads the Local APIC registers and copies them into the input argument. The data format and layout are the same as documented in the architecture manual.

If KVM\_X2APIC\_API\_USE\_32BIT\_IDS feature of KVM\_CAP\_X2APIC\_API is enabled, then the format of APIC\_ID register depends on the APIC mode (reported by MSR\_IA32\_APICBASE) of its VCPU. x2APIC stores APIC ID in the APIC\_ID register (bytes 32-35). xAPIC only allows an 8-bit APIC ID which is stored in bits 31-24 of the APIC register, or equivalently in byte 35 of struct kvm\_lapic\_state's regs field. KVM\_GET\_LAPIC must then be called after MSR\_IA32\_APICBASE has been set with KVM\_SET\_MSR.

If KVM\_X2APIC\_API\_USE\_32BIT\_IDS feature is disabled, struct kvm\_lapic\_state always uses xAPIC format.

## 4.58 KVM\_SET\_LAPIC

**Capability**

KVM\_CAP\_IRQCHIP

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_lapic\_state (in)

**Returns**

0 on success, -1 on error

```
#define KVM_APIC_REG_SIZE 0x400
struct kvm_lapic_state {
    char regs[KVM_APIC_REG_SIZE];
};
```

Copies the input argument into the Local APIC registers. The data format and layout are the same as documented in the architecture manual.

The format of the APIC ID register (bytes 32-35 of struct kvm\_lapic\_state's regs field) depends on the state of the KVM\_CAP\_X2APIC\_API capability. See the note in KVM\_GET\_LAPIC.

## 4.59 KVM\_IOEVENTFD

**Capability**

KVM\_CAP\_IOEVENTFD

**Architectures**

all

**Type**

vm ioctl

**Parameters**

struct kvm\_ioeventfd (in)

**Returns**

0 on success, !0 on error

This ioctl attaches or detaches an ioeventfd to a legal pio/mmio address within the guest. A guest write in the registered address will signal the provided event instead of triggering an exit.

```
struct kvm_ioeventfd {
    __u64 datamatch;
    __u64 addr;           /* legal pio/mmio address */
    __u32 len;           /* 0, 1, 2, 4, or 8 bytes */
    __s32 fd;
```

(continues on next page)

(continued from previous page)

```
    __u32 flags;
    __u8  pad[36];
};
```

For the special case of virtio-ccw devices on s390, the ioevent is matched to a subchannel/virtqueue tuple instead.

The following flags are defined:

```
#define KVM_IOEVENTFD_FLAG_DATAMATCH (1 << kvm_ioeventfd_flag_nr_
↳datamatch)
#define KVM_IOEVENTFD_FLAG_PIO      (1 << kvm_ioeventfd_flag_nr_
↳pio)
#define KVM_IOEVENTFD_FLAG_DEASSIGN (1 << kvm_ioeventfd_flag_nr_
↳deassign)
#define KVM_IOEVENTFD_FLAG_VIRTIO_CCW_NOTIFY \
    (1 << kvm_ioeventfd_flag_nr_virtio_ccw_notify)
```

If datamatch flag is set, the event will be signaled only if the written value to the registered address is equal to datamatch in struct kvm\_ioeventfd.

For virtio-ccw devices, addr contains the subchannel id and datamatch the virtqueue index.

With KVM\_CAP\_IOEVENTFD\_ANY\_LENGTH, a zero length ioeventfd is allowed, and the kernel will ignore the length of guest write and may get a faster vmexit. The speedup may only apply to specific architectures, but the ioeventfd will work anyway.

## 4.60 KVM\_DIRTY\_TLB

### Capability

KVM\_CAP\_SW\_TLB

### Architectures

ppc

### Type

vcpu ioctl

### Parameters

struct kvm\_dirty\_tlb (in)

### Returns

0 on success, -1 on error

```
struct kvm_dirty_tlb {
    __u64 bitmap;
    __u32 num_dirty;
};
```

This must be called whenever userspace has changed an entry in the shared TLB, prior to calling KVM\_RUN on the associated vcpu.

The “bitmap” field is the userspace address of an array. This array consists of a number of bits, equal to the total number of TLB entries as determined by the last successful call to `KVM_CONFIG_TLB`, rounded up to the nearest multiple of 64.

Each bit corresponds to one TLB entry, ordered the same as in the shared TLB array.

The array is little-endian: the bit 0 is the least significant bit of the first byte, bit 8 is the least significant bit of the second byte, etc. This avoids any complications with differing word sizes.

The “num\_dirty” field is a performance hint for KVM to determine whether it should skip processing the bitmap and just invalidate everything. It must be set to the number of set bits in the bitmap.

## 4.62 KVM\_CREATE\_SPAPR\_TCE

### Capability

`KVM_CAP_SPAPR_TCE`

### Architectures

powerpc

### Type

vm ioctl

### Parameters

`struct kvm_create_spapr_tce` (in)

### Returns

file descriptor for manipulating the created TCE table

This creates a virtual TCE (translation control entry) table, which is an IOMMU for PAPR-style virtual I/O. It is used to translate logical addresses used in virtual I/O into guest physical addresses, and provides a scatter/gather capability for PAPR virtual I/O.

```
/* for KVM_CAP_SPAPR_TCE */
struct kvm_create_spapr_tce {
    __u64 liobn;
    __u32 window_size;
};
```

The `liobn` field gives the logical IO bus number for which to create a TCE table. The `window_size` field specifies the size of the DMA window which this TCE table will translate - the table will contain one 64 bit TCE entry for every 4kiB of the DMA window.

When the guest issues an `H_PUT_TCE` hcall on a `liobn` for which a TCE table has been created using this `ioctl()`, the kernel will handle it in real mode, updating the TCE table. `H_PUT_TCE` calls for other `liobns` will cause a vm exit and must be handled by userspace.

The return value is a file descriptor which can be passed to `mmap(2)` to map the created TCE table into userspace. This lets userspace read the entries written by

kernel-handled `H_PUT_TCE` calls, and also lets userspace update the TCE table directly which is useful in some circumstances.

### 4.63 KVM\_ALLOCATE\_RMA

**Capability**

`KVM_CAP_PPC_RMA`

**Architectures**

powerpc

**Type**

vm ioctl

**Parameters**

`struct kvm_allocate_rma (out)`

**Returns**

file descriptor for mapping the allocated RMA

This allocates a Real Mode Area (RMA) from the pool allocated at boot time by the kernel. An RMA is a physically-contiguous, aligned region of memory used on older POWER processors to provide the memory which will be accessed by real-mode (MMU off) accesses in a KVM guest. POWER processors support a set of sizes for the RMA that usually includes 64MB, 128MB, 256MB and some larger powers of two.

```
/* for KVM_ALLOCATE_RMA */
struct kvm_allocate_rma {
    __u64 rma_size;
};
```

The return value is a file descriptor which can be passed to `mmap(2)` to map the allocated RMA into userspace. The mapped area can then be passed to the `KVM_SET_USER_MEMORY_REGION` ioctl to establish it as the RMA for a virtual machine. The size of the RMA in bytes (which is fixed at host kernel boot time) is returned in the `rma_size` field of the argument structure.

The `KVM_CAP_PPC_RMA` capability is 1 or 2 if the `KVM_ALLOCATE_RMA` ioctl is supported; 2 if the processor requires all virtual machines to have an RMA, or 1 if the processor can use an RMA but doesn't require it, because it supports the Virtual RMA (VRMA) facility.

### 4.64 KVM\_NMI

**Capability**

`KVM_CAP_USER_NMI`

**Architectures**

x86

**Type**

vcpu ioctl



**Parameters**

none

**Returns**

0 on success, -1 on error

Queues an NMI on the thread's vcpu. Note this is well defined only when KVM\_CREATE\_IRQCHIP has not been called, since this is an interface between the virtual cpu core and virtual local APIC. After KVM\_CREATE\_IRQCHIP has been called, this interface is completely emulated within the kernel.

To use this to emulate the LINT1 input with KVM\_CREATE\_IRQCHIP, use the following algorithm:

- pause the vcpu
- read the local APIC's state (KVM\_GET\_LAPIC)
- check whether changing LINT1 will queue an NMI (see the LVT entry for LINT1)
- if so, issue KVM\_NMI
- resume the vcpu

Some guests configure the LINT1 NMI input to cause a panic, aiding in debugging.

**4.65 KVM\_S390\_UCAS\_MAP****Capability**

KVM\_CAP\_S390\_UCONTROL

**Architectures**

s390

**Type**

vcpu ioctl

**Parameters**

struct kvm\_s390\_ucas\_mapping (in)

**Returns**

0 in case of success

The parameter is defined like this:

```
struct kvm_s390_ucas_mapping {
    __u64 user_addr;
    __u64 vcpu_addr;
    __u64 length;
};
```

This ioctl maps the memory at “user\_addr” with the length “length” to the vcpu's address space starting at “vcpu\_addr”. All parameters need to be aligned by 1 megabyte.

## 4.66 KVM\_S390\_UCAS\_UNMAP

**Capability**

KVM\_CAP\_S390\_UCONTROL

**Architectures**

s390

**Type**

vcpu ioctl

**Parameters**

struct kvm\_s390\_ucas\_mapping (in)

**Returns**

0 in case of success

The parameter is defined like this:

```
struct kvm_s390_ucas_mapping {
    __u64 user_addr;
    __u64 vcpu_addr;
    __u64 length;
};
```

This ioctl unmaps the memory in the vcpu's address space starting at "vcpu\_addr" with the length "length". The field "user\_addr" is ignored. All parameters need to be aligned by 1 megabyte.

## 4.67 KVM\_S390\_VCPU\_FAULT

**Capability**

KVM\_CAP\_S390\_UCONTROL

**Architectures**

s390

**Type**

vcpu ioctl

**Parameters**

vcpu absolute address (in)

**Returns**

0 in case of success

This call creates a page table entry on the virtual cpu's address space (for user controlled virtual machines) or the virtual machine's address space (for regular virtual machines). This only works for minor faults, thus it's recommended to access subject memory page via the user page table upfront. This is useful to handle validity intercepts for user controlled virtual machines to fault in the virtual cpu's lowcore pages prior to calling the KVM\_RUN ioctl.

## 4.68 KVM\_SET\_ONE\_REG

### Capability

KVM\_CAP\_ONE\_REG

### Architectures

all

### Type

vcpu ioctl

### Parameters

struct kvm\_one\_reg (in)

### Returns

0 on success, negative value on failure

Errors:

ENOEL	no such register
EIN-	invalid register ID, or no such register or used with VMs in
VAL	protected virtualization mode on s390
EPERN	(arm64) register access not allowed before vcpu finalization

(These error codes are indicative only: do not rely on a specific error code being returned in a specific situation.)

```
struct kvm_one_reg {
    __u64 id;
    __u64 addr;
};
```

Using this ioctl, a single vcpu register can be set to a specific value defined by user space with the passed in struct `kvm_one_reg`, where `id` refers to the register identifier as described below and `addr` is a pointer to a variable with the respective size. There can be architecture agnostic and architecture specific registers. Each have their own range of operation and their own constants and width. To keep track of the implemented registers, find a list below:

Arch	Register	Width (bits)
PPC	KVM_REG_PPC_HIOR	64
PPC	KVM_REG_PPC_IAC1	64
PPC	KVM_REG_PPC_IAC2	64
PPC	KVM_REG_PPC_IAC3	64
PPC	KVM_REG_PPC_IAC4	64
PPC	KVM_REG_PPC_DAC1	64
PPC	KVM_REG_PPC_DAC2	64
PPC	KVM_REG_PPC_DABR	64
PPC	KVM_REG_PPC_DSCR	64
PPC	KVM_REG_PPC_PURR	64

continues on next page

Table 1 - continued from previous page

Arch	Register	Width (bits)
PPC	KVM_REG_PPC_SPURR	64
PPC	KVM_REG_PPC_DAR	64
PPC	KVM_REG_PPC_DSISR	32
PPC	KVM_REG_PPC_AMR	64
PPC	KVM_REG_PPC_UAMOR	64
PPC	KVM_REG_PPC_MMCR0	64
PPC	KVM_REG_PPC_MMCR1	64
PPC	KVM_REG_PPC_MMCR2	64
PPC	KVM_REG_PPC_MMCR3	64
PPC	KVM_REG_PPC_SIAR	64
PPC	KVM_REG_PPC_SDAR	64
PPC	KVM_REG_PPC_SIER	64
PPC	KVM_REG_PPC_SIER2	64
PPC	KVM_REG_PPC_SIER3	64
PPC	KVM_REG_PPC_PMC1	32
PPC	KVM_REG_PPC_PMC2	32
PPC	KVM_REG_PPC_PMC3	32
PPC	KVM_REG_PPC_PMC4	32
PPC	KVM_REG_PPC_PMC5	32
PPC	KVM_REG_PPC_PMC6	32
PPC	KVM_REG_PPC_PMC7	32
PPC	KVM_REG_PPC_PMC8	32
PPC	KVM_REG_PPC_FPR0	64
...		
PPC	KVM_REG_PPC_FPR31	64
PPC	KVM_REG_PPC_VR0	128
...		
PPC	KVM_REG_PPC_VR31	128
PPC	KVM_REG_PPC_VSR0	128
...		
PPC	KVM_REG_PPC_VSR31	128
PPC	KVM_REG_PPC_FPSCR	64
PPC	KVM_REG_PPC_VSCR	32
PPC	KVM_REG_PPC_VPA_ADDR	64
PPC	KVM_REG_PPC_VPA_SLB	128
PPC	KVM_REG_PPC_VPA_DTL	128
PPC	KVM_REG_PPC_EPCR	32
PPC	KVM_REG_PPC_EPR	32
PPC	KVM_REG_PPC_TCR	32
PPC	KVM_REG_PPC_TSR	32
PPC	KVM_REG_PPC_OR_TSR	32
PPC	KVM_REG_PPC_CLEAR_TSR	32
PPC	KVM_REG_PPC_MAS0	32
PPC	KVM_REG_PPC_MAS1	32
PPC	KVM_REG_PPC_MAS2	64
PPC	KVM_REG_PPC_MAS7_3	64

continues on next page

Table 1 - continued from previous page

Arch	Register	Width (bits)
PPC	KVM_REG_PPC_MAS4	32
PPC	KVM_REG_PPC_MAS6	32
PPC	KVM_REG_PPC_MMUCFG	32
PPC	KVM_REG_PPC_TLB0CFG	32
PPC	KVM_REG_PPC_TLB1CFG	32
PPC	KVM_REG_PPC_TLB2CFG	32
PPC	KVM_REG_PPC_TLB3CFG	32
PPC	KVM_REG_PPC_TLB0PS	32
PPC	KVM_REG_PPC_TLB1PS	32
PPC	KVM_REG_PPC_TLB2PS	32
PPC	KVM_REG_PPC_TLB3PS	32
PPC	KVM_REG_PPC_EPTCFG	32
PPC	KVM_REG_PPC_ICP_STATE	64
PPC	KVM_REG_PPC_VP_STATE	128
PPC	KVM_REG_PPC_TB_OFFSET	64
PPC	KVM_REG_PPC_SPMC1	32
PPC	KVM_REG_PPC_SPMC2	32
PPC	KVM_REG_PPC_IAMR	64
PPC	KVM_REG_PPC_TFHAR	64
PPC	KVM_REG_PPC_TFIAR	64
PPC	KVM_REG_PPC_TEXASR	64
PPC	KVM_REG_PPC_FSCR	64
PPC	KVM_REG_PPC_PSPB	32
PPC	KVM_REG_PPC_EBBHR	64
PPC	KVM_REG_PPC_EBBRR	64
PPC	KVM_REG_PPC_BESCR	64
PPC	KVM_REG_PPC_TAR	64
PPC	KVM_REG_PPC_DPDES	64
PPC	KVM_REG_PPC_DAWR	64
PPC	KVM_REG_PPC_DAWRX	64
PPC	KVM_REG_PPC_CIABR	64
PPC	KVM_REG_PPC_IC	64
PPC	KVM_REG_PPC_VTB	64
PPC	KVM_REG_PPC_CSIGR	64
PPC	KVM_REG_PPC_TACR	64
PPC	KVM_REG_PPC_TCSCR	64
PPC	KVM_REG_PPC_PID	64
PPC	KVM_REG_PPC_ACOP	64
PPC	KVM_REG_PPC_VRSAVE	32
PPC	KVM_REG_PPC_LPCR	32
PPC	KVM_REG_PPC_LPCR_64	64
PPC	KVM_REG_PPC_PPR	64
PPC	KVM_REG_PPC_ARCH_COMPAT	32
PPC	KVM_REG_PPC_DABRX	32
PPC	KVM_REG_PPC_WORT	64
PPC	KVM_REG_PPC_SPRG9	64
PPC	KVM_REG_PPC_DBSR	32
PPC	KVM_REG_PPC_TIDR	64

continues on next page

Table 1 - continued from previous page

Arch	Register	Width (bits)
PPC	KVM_REG_PPC_PSSCR	64
PPC	KVM_REG_PPC_DEC_EXPIRY	64
PPC	KVM_REG_PPC_PTCR	64
PPC	KVM_REG_PPC_TM_GPR0	64
...		
PPC	KVM_REG_PPC_TM_GPR31	64
PPC	KVM_REG_PPC_TM_VSR0	128
...		
PPC	KVM_REG_PPC_TM_VSR63	128
PPC	KVM_REG_PPC_TM_CR	64
PPC	KVM_REG_PPC_TM_LR	64
PPC	KVM_REG_PPC_TM_CTR	64
PPC	KVM_REG_PPC_TM_FPSCR	64
PPC	KVM_REG_PPC_TM_AMR	64
PPC	KVM_REG_PPC_TM_PPR	64
PPC	KVM_REG_PPC_TM_VRSAVE	64
PPC	KVM_REG_PPC_TM_VSCR	32
PPC	KVM_REG_PPC_TM_DSCR	64
PPC	KVM_REG_PPC_TM_TAR	64
PPC	KVM_REG_PPC_TM_XER	64
MIPS	KVM_REG_MIPS_R0	64
...		
MIPS	KVM_REG_MIPS_R31	64
MIPS	KVM_REG_MIPS_HI	64
MIPS	KVM_REG_MIPS_LO	64
MIPS	KVM_REG_MIPS_PC	64
MIPS	KVM_REG_MIPS_CP0_INDEX	32
MIPS	KVM_REG_MIPS_CP0_ENTRYLO0	64
MIPS	KVM_REG_MIPS_CP0_ENTRYLO1	64
MIPS	KVM_REG_MIPS_CP0_CONTEXT	64
MIPS	KVM_REG_MIPS_CP0_CONTEXTCONFIG	32
MIPS	KVM_REG_MIPS_CP0_USERLOCAL	64
MIPS	KVM_REG_MIPS_CP0_XCONTEXTCONFIG	64
MIPS	KVM_REG_MIPS_CP0_PAGEMASK	32
MIPS	KVM_REG_MIPS_CP0_PAGEGRAIN	32
MIPS	KVM_REG_MIPS_CP0_SEGCTL0	64
MIPS	KVM_REG_MIPS_CP0_SEGCTL1	64
MIPS	KVM_REG_MIPS_CP0_SEGCTL2	64
MIPS	KVM_REG_MIPS_CP0_PWBASE	64
MIPS	KVM_REG_MIPS_CP0_PWFIELD	64
MIPS	KVM_REG_MIPS_CP0_PWSIZE	64
MIPS	KVM_REG_MIPS_CP0_WIRED	32
MIPS	KVM_REG_MIPS_CP0_PWCTL	32
MIPS	KVM_REG_MIPS_CP0_HWRENA	32
MIPS	KVM_REG_MIPS_CP0_BADVADDR	64
MIPS	KVM_REG_MIPS_CP0_BADINSTR	32
MIPS	KVM_REG_MIPS_CP0_BADINSTRP	32
MIPS	KVM_REG_MIPS_CP0_COUNT	32

continues on next page

Table 1 - continued from previous page

Arch	Register	Width (bits)
MIPS	KVM_REG_MIPS_CP0_ENTRYHI	64
MIPS	KVM_REG_MIPS_CP0_COMPARE	32
MIPS	KVM_REG_MIPS_CP0_STATUS	32
MIPS	KVM_REG_MIPS_CP0_INTCTL	32
MIPS	KVM_REG_MIPS_CP0_CAUSE	32
MIPS	KVM_REG_MIPS_CP0_EPC	64
MIPS	KVM_REG_MIPS_CP0_PRID	32
MIPS	KVM_REG_MIPS_CP0_EBASE	64
MIPS	KVM_REG_MIPS_CP0_CONFIG	32
MIPS	KVM_REG_MIPS_CP0_CONFIG1	32
MIPS	KVM_REG_MIPS_CP0_CONFIG2	32
MIPS	KVM_REG_MIPS_CP0_CONFIG3	32
MIPS	KVM_REG_MIPS_CP0_CONFIG4	32
MIPS	KVM_REG_MIPS_CP0_CONFIG5	32
MIPS	KVM_REG_MIPS_CP0_CONFIG7	32
MIPS	KVM_REG_MIPS_CP0_XCONTEXT	64
MIPS	KVM_REG_MIPS_CP0_ERROREPC	64
MIPS	KVM_REG_MIPS_CP0_KSCRATCH1	64
MIPS	KVM_REG_MIPS_CP0_KSCRATCH2	64
MIPS	KVM_REG_MIPS_CP0_KSCRATCH3	64
MIPS	KVM_REG_MIPS_CP0_KSCRATCH4	64
MIPS	KVM_REG_MIPS_CP0_KSCRATCH5	64
MIPS	KVM_REG_MIPS_CP0_KSCRATCH6	64
MIPS	KVM_REG_MIPS_CP0_MAAR(0..63)	64
MIPS	KVM_REG_MIPS_COUNT_CTL	64
MIPS	KVM_REG_MIPS_COUNT_RESUME	64
MIPS	KVM_REG_MIPS_COUNT_HZ	64
MIPS	KVM_REG_MIPS_FPR_32(0..31)	32
MIPS	KVM_REG_MIPS_FPR_64(0..31)	64
MIPS	KVM_REG_MIPS_VEC_128(0..31)	128
MIPS	KVM_REG_MIPS_FCR_IR	32
MIPS	KVM_REG_MIPS_FCR_CSR	32
MIPS	KVM_REG_MIPS_MSA_IR	32
MIPS	KVM_REG_MIPS_MSA_CSR	32

ARM registers are mapped using the lower 32 bits. The upper 16 of that is the register group type, or coprocessor number:

ARM core registers have the following id bit patterns:

```
0x4020 0000 0010 <index into the kvm_regs struct:16>
```

ARM 32-bit CP15 registers have the following id bit patterns:

```
0x4020 0000 000F <zero:1> <crn:4> <crm:4> <opc1:4> <opc2:3>
```

ARM 64-bit CP15 registers have the following id bit patterns:

```
0x4030 0000 000F <zero:1> <zero:4> <crm:4> <opcl:4> <zero:3>
```

ARM CCSIDR registers are demultiplexed by CSSELR value:

```
0x4020 0000 0011 00 <csselr:8>
```

ARM 32-bit VFP control registers have the following id bit patterns:

```
0x4020 0000 0012 1 <regno:12>
```

ARM 64-bit FP registers have the following id bit patterns:

```
0x4030 0000 0012 0 <regno:12>
```

ARM firmware pseudo-registers have the following bit pattern:

```
0x4030 0000 0014 <regno:16>
```

arm64 registers are mapped using the lower 32 bits. The upper 16 of that is the register group type, or coprocessor number:

arm64 core/FP-SIMD registers have the following id bit patterns. Note that the size of the access is variable, as the `kvm_regs` structure contains elements ranging from 32 to 128 bits. The index is a 32bit value in the `kvm_regs` structure seen as a 32bit array:

```
0x60x0 0000 0010 <index into the kvm_regs struct:16>
```

Specifically:



Encoding	Register	Bits	kvm_regs member
0x6030 0000 0010 X0 0000		64	regs.reg[0]
0x6030 0000 0010 X1 0002		64	regs.reg[1]
...			
0x6030 0000 0010 X30 003c		64	regs.reg[30]
0x6030 0000 0010 SP 003e		64	regs.sp
0x6030 0000 0010 PC 0040		64	regs.pc
0x6030 0000 0010 PSTATE 0042		64	regs.pstate
0x6030 0000 0010 SP_EL1 0044		64	sp_el1
0x6030 0000 0010 ELR_EL1 0046		64	elr_el1
0x6030 0000 0010 SPSR_EL1 0048		64	spsr[KVM_SPSR_EL1] (alias SPSR_SVC)
0x6030 0000 0010 SPSR_ABT 004a		64	spsr[KVM_SPSR_ABT]
0x6030 0000 0010 SPSR_UND 004c		64	spsr[KVM_SPSR_UND]
0x6030 0000 0010 SPSR_IRQ 004e		64	spsr[KVM_SPSR_IRQ]
0x6060 0000 0010 SPSR_FIQ 0050		64	spsr[KVM_SPSR_FIQ]
0x6040 0000 0010 V0 0054		128	fp_regs.vregs[0] <sup>1</sup>
0x6040 0000 0010 V1 0058		128	fp_regs.vregs[1] <sup>Page 53, 1</sup>
...			
0x6040 0000 0010 V31 00d0		128	fp_regs.vregs[31] <sup>1</sup>
0x6020 0000 0010 FPSR 00d4		32	fp_regs.fpsr
0x6020 0000 0010 FPCR 00d5		32	fp_regs.fpcr

arm64 CCSIDR registers are demultiplexed by CSSELR value:

```
0x6020 0000 0011 00 <csselr:8>
```

arm64 system registers have the following id bit patterns:

<sup>1</sup> These encodings are not accepted for SVE-enabled vcpus. See KVM\_ARM\_VCPU\_INIT.

The equivalent register content can be accessed via bits [127:0] of the corresponding SVE Zn registers instead for vcpus that have SVE enabled (see below).

```
0x6030 0000 0013 <op0:2> <op1:3> <crn:4> <crm:4> <op2:3>
```

**Warning:** Two system register IDs do not follow the specified pattern. These are KVM\_REG\_ARM\_TIMER\_CVAL and KVM\_REG\_ARM\_TIMER\_CNT, which map to system registers CNTV\_CVAL\_EL0 and CNTVCT\_EL0 respectively. These two had their values accidentally swapped, which means TIMER\_CVAL is derived from the register encoding for CNTVCT\_EL0 and TIMER\_CNT is derived from the register encoding for CNTV\_CVAL\_EL0. As this is API, it must remain this way.

arm64 firmware pseudo-registers have the following bit pattern:

```
0x6030 0000 0014 <regno:16>
```

arm64 SVE registers have the following bit patterns:

```
0x6080 0000 0015 00 <n:5> <slice:5>   Zn bits[2048*slice + 2047 : ↵
↵2048*slice]
0x6050 0000 0015 04 <n:4> <slice:5>   Pn bits[256*slice + 255 : ↵
↵256*slice]
0x6050 0000 0015 060 <slice:5>         FFR bits[256*slice + 255 : ↵
↵256*slice]
0x6060 0000 0015 ffff                   KVM_REG_ARM64_SVE_VLS pseudo-
↵register
```

Access to register IDs where  $2048 * \text{slice} \geq 128 * \text{max\_vq}$  will fail with ENOENT.  $\text{max\_vq}$  is the vcpu's maximum supported vector length in 128-bit quadwords: see<sup>2</sup> below.

These registers are only accessible on vcpus for which SVE is enabled. See KVM\_ARM\_VCPU\_INIT for details.

In addition, except for KVM\_REG\_ARM64\_SVE\_VLS, these registers are not accessible until the vcpu's SVE configuration has been finalized using KVM\_ARM\_VCPU\_FINALIZE(KVM\_ARM\_VCPU\_SVE). See KVM\_ARM\_VCPU\_INIT and KVM\_ARM\_VCPU\_FINALIZE for more information about this procedure.

KVM\_REG\_ARM64\_SVE\_VLS is a pseudo-register that allows the set of vector lengths supported by the vcpu to be discovered and configured by userspace. When transferred to or from user memory via KVM\_GET\_ONE\_REG or KVM\_SET\_ONE\_REG, the value of this register is of type `_u64[KVM_ARM64_SVE_VLS_WORDS]`, and encodes the set of vector lengths as follows:

```
_u64 vector_lengths[KVM_ARM64_SVE_VLS_WORDS];
```

(continues on next page)

<sup>2</sup> The maximum value  $\text{vq}$  for which the above condition is true is  $\text{max\_vq}$ . This is the maximum vector length available to the guest on this vcpu, and determines which register slices are visible through this ioctl interface.

(continued from previous page)

```

if (vq >= SVE_VQ_MIN && vq <= SVE_VQ_MAX &&
    ((vector_lengths[(vq - KVM_ARM64_SVE_VQ_MIN) / 64] >>
      ((vq - KVM_ARM64_SVE_VQ_MIN) % 64)) & 1))
    /* Vector length vq * 16 bytes supported */
else
    /* Vector length vq * 16 bytes not supported */

```

(See Documentation/arm64/sve.rst for an explanation of the “vq” nomenclature.)

KVM\_REG\_ARM64\_SVE\_VLS is only accessible after KVM\_ARM\_VCPU\_INIT. KVM\_ARM\_VCPU\_INIT initialises it to the best set of vector lengths that the host supports.

Userspace may subsequently modify it if desired until the vcpu’s SVE configuration is finalized using KVM\_ARM\_VCPU\_FINALIZE(KVM\_ARM\_VCPU\_SVE).

Apart from simply removing all vector lengths from the host set that exceed some value, support for arbitrarily chosen sets of vector lengths is hardware-dependent and may not be available. Attempting to configure an invalid set of vector lengths via KVM\_SET\_ONE\_REG will fail with EINVAL.

After the vcpu’s SVE configuration is finalized, further attempts to write this register will fail with EPERM.

MIPS registers are mapped using the lower 32 bits. The upper 16 of that is the register group type:

MIPS core registers (see above) have the following id bit patterns:

```
0x7030 0000 0000 <reg:16>
```

MIPS CP0 registers (see KVM\_REG\_MIPS\_CP0\_\* above) have the following id bit patterns depending on whether they’re 32-bit or 64-bit registers:

```

0x7020 0000 0001 00 <reg:5> <sel:3>    (32-bit)
0x7030 0000 0001 00 <reg:5> <sel:3>    (64-bit)

```

Note: KVM\_REG\_MIPS\_CP0\_ENTRYLO0 and KVM\_REG\_MIPS\_CP0\_ENTRYLO1 are the MIPS64 versions of the EntryLo registers regardless of the word size of the host hardware, host kernel, guest, and whether XPA is present in the guest, i.e. with the RI and XI bits (if they exist) in bits 63 and 62 respectively, and the PFNX field starting at bit 30.

MIPS MAARs (see KVM\_REG\_MIPS\_CP0\_MAAR(\*) above) have the following id bit patterns:

```
0x7030 0000 0001 01 <reg:8>
```

MIPS KVM control registers (see above) have the following id bit patterns:

```
0x7030 0000 0002 <reg:16>
```

MIPS FPU registers (see KVM\_REG\_MIPS\_FPR\_{32,64}() above) have the following id bit patterns depending on the size of the register being accessed. They are

always accessed according to the current guest FPU mode (Status.FR and Config5.FRE), i.e. as the guest would see them, and they become unpredictable if the guest FPU mode is changed. MIPS SIMD Architecture (MSA) vector registers (see KVM\_REG\_MIPS\_VEC\_128() above) have similar patterns as they overlap the FPU registers:

```
0x7020 0000 0003 00 <0:3> <reg:5> (32-bit FPU registers)
0x7030 0000 0003 00 <0:3> <reg:5> (64-bit FPU registers)
0x7040 0000 0003 00 <0:3> <reg:5> (128-bit MSA vector registers)
```

MIPS FPU control registers (see KVM\_REG\_MIPS\_FCR\_{IR,CSR} above) have the following id bit patterns:

```
0x7020 0000 0003 01 <0:3> <reg:5>
```

MIPS MSA control registers (see KVM\_REG\_MIPS\_MSA\_{IR,CSR} above) have the following id bit patterns:

```
0x7020 0000 0003 02 <0:3> <reg:5>
```

## 4.69 KVM\_GET\_ONE\_REG

### Capability

KVM\_CAP\_ONE\_REG

### Architectures

all

### Type

vcpu ioctl

### Parameters

struct kvm\_one\_reg (in and out)

### Returns

0 on success, negative value on failure

Errors include:

ENOENT	no such register
EIN- VAL	invalid register ID, or no such register or used with VMs in protected virtualization mode on s390
EPERM	(arm64) register access not allowed before vcpu finalization

(These error codes are indicative only: do not rely on a specific error code being returned in a specific situation.)

This ioctl allows to receive the value of a single register implemented in a vcpu. The register to read is indicated by the “id” field of the kvm\_one\_reg struct passed in. On success, the register value can be found at the memory location pointed to by “addr” .

The list of registers accessible using this interface is identical to the list in 4.68.

## 4.70 KVM\_KVMCLOCK\_CTRL

**Capability**

KVM\_CAP\_KVMCLOCK\_CTRL

**Architectures**

Any that implement pvclocks (currently x86 only)

**Type**

vcpu ioctl

**Parameters**

None

**Returns**

0 on success, -1 on error

This ioctl sets a flag accessible to the guest indicating that the specified vCPU has been paused by the host userspace.

The host will set a flag in the pvclock structure that is checked from the soft lockup watchdog. The flag is part of the pvclock structure that is shared between guest and host, specifically the second bit of the flags field of the pvclock\_vcpu\_time\_info structure. It will be set exclusively by the host and read/cleared exclusively by the guest. The guest operation of checking and clearing the flag must be an atomic operation so load-link/store-conditional, or equivalent must be used. There are two cases where the guest will clear the flag: when the soft lockup watchdog timer resets itself or when a soft lockup is detected. This ioctl can be called any time after pausing the vcpu, but before it is resumed.

## 4.71 KVM\_SIGNAL\_MSI

**Capability**

KVM\_CAP\_SIGNAL\_MSI

**Architectures**

x86 arm arm64

**Type**

vm ioctl

**Parameters**

struct kvm\_msi (in)

**Returns**

&gt;0 on delivery, 0 if guest blocked the MSI, and -1 on error

Directly inject a MSI message. Only valid with in-kernel irqchip that handles MSI messages.

```
struct kvm_msi {
    __u32 address_lo;
    __u32 address_hi;
    __u32 data;
    __u32 flags;
```

(continues on next page)

(continued from previous page)

```
    __u32 devid;
    __u8  pad[12];
};
```

**flags:**

**KVM\_MSI\_VALID\_DEVID:** devid contains a valid value. The per-VM **KVM\_CAP\_MSI\_DEVID** capability advertises the requirement to provide the device ID. If this capability is not available, userspace should never set the **KVM\_MSI\_VALID\_DEVID** flag as the ioctl might fail.

If **KVM\_MSI\_VALID\_DEVID** is set, devid contains a unique device identifier for the device that wrote the MSI message. For PCI, this is usually a BFD identifier in the lower 16 bits.

On x86, address\_hi is ignored unless the **KVM\_X2APIC\_API\_USE\_32BIT\_IDS** feature of **KVM\_CAP\_X2APIC\_API** capability is enabled. If it is enabled, address\_hi bits 31-8 provide bits 31-8 of the destination id. Bits 7-0 of address\_hi must be zero.

## 4.71 KVM\_CREATE\_PIT2

**Capability**

**KVM\_CAP\_PIT2**

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct kvm\_pit\_config (in)

**Returns**

0 on success, -1 on error

Creates an in-kernel device model for the i8254 PIT. This call is only valid after enabling in-kernel irqchip support via **KVM\_CREATE\_IRQCHIP**. The following parameters have to be passed:

```
struct kvm_pit_config {
    __u32 flags;
    __u32 pad[15];
};
```

Valid flags are:

```
#define KVM_PIT_SPEAKER_DUMMY    1 /* emulate speaker port stub */
```

PIT timer interrupts may use a per-VM kernel thread for injection. If it exists, this thread will have a name of the following pattern:

```
kvm-pit/<owner-process-pid>
```

When running a guest with elevated priorities, the scheduling parameters of this thread may have to be adjusted accordingly.

This IOCTL replaces the obsolete KVM\_CREATE\_PIT.

#### 4.72 KVM\_GET\_PIT2

##### Capability

KVM\_CAP\_PIT\_STATE2

##### Architectures

x86

##### Type

vm ioctl

##### Parameters

struct kvm\_pit\_state2 (out)

##### Returns

0 on success, -1 on error

Retrieves the state of the in-kernel PIT model. Only valid after KVM\_CREATE\_PIT2. The state is returned in the following structure:

```
struct kvm_pit_state2 {
    struct kvm_pit_channel_state channels[3];
    __u32 flags;
    __u32 reserved[9];
};
```

Valid flags are:

```
/* disable PIT in HPET legacy mode */
#define KVM_PIT_FLAGS_HPET_LEGACY 0x00000001
```

This IOCTL replaces the obsolete KVM\_GET\_PIT.

#### 4.73 KVM\_SET\_PIT2

##### Capability

KVM\_CAP\_PIT\_STATE2

##### Architectures

x86

##### Type

vm ioctl

##### Parameters

struct kvm\_pit\_state2 (in)

**Returns**

0 on success, -1 on error

Sets the state of the in-kernel PIT model. Only valid after KVM\_CREATE\_PIT2. See KVM\_GET\_PIT2 for details on struct kvm\_pit\_state2.

This IOCTL replaces the obsolete KVM\_SET\_PIT.

**4.74 KVM\_PPC\_GET\_SMMU\_INFO****Capability**

KVM\_CAP\_PPC\_GET\_SMMU\_INFO

**Architectures**

powerpc

**Type**

vm ioctl

**Parameters**

None

**Returns**

0 on success, -1 on error

This populates and returns a structure describing the features of the “Server” class MMU emulation supported by KVM. This can in turn be used by userspace to generate the appropriate device-tree properties for the guest operating system.

The structure contains some global information, followed by an array of supported segment page sizes:

```
struct kvm_ppc_smmu_info {
    __u64 flags;
    __u32 slb_size;
    __u32 pad;
    struct kvm_ppc_one_seg_page_size sps[KVM_PPC_PAGE_SIZES_MAX_
↪SZ];
};
```

The supported flags are:

- **KVM\_PPC\_PAGE\_SIZES\_REAL:**

When that flag is set, guest page sizes must “fit” the backing store page sizes. When not set, any page size in the list can be used regardless of how they are backed by userspace.

- **KVM\_PPC\_1T\_SEGMENTS**

The emulated MMU supports 1T segments in addition to the standard 256M ones.

- **KVM\_PPC\_NO\_HASH**

This flag indicates that HPT guests are not supported by KVM, thus all guests must use radix MMU mode.

The “slb\_size” field indicates how many SLB entries are supported



The “sps” array contains 8 entries indicating the supported base page sizes for a segment in increasing order. Each entry is defined as follow:

```
struct kvm_ppc_one_seg_page_size {
    __u32 page_shift;      /* Base page shift of segment (or 0) */
    __u32 slb_enc;        /* SLB encoding for BookS */
    struct kvm_ppc_one_page_size enc[KVM_PPC_PAGE_SIZES_MAX_SZ];
};
```

An entry with a “page\_shift” of 0 is unused. Because the array is organized in increasing order, a lookup can stop when encountering such an entry.

The “slb\_enc” field provides the encoding to use in the SLB for the page size. The bits are in positions such as the value can directly be OR’ed into the “vsid” argument of the slbmte instruction.

The “enc” array is a list which for each of those segment base page size provides the list of supported actual page sizes (which can be only larger or equal to the base page size), along with the corresponding encoding in the hash PTE. Similarly, the array is 8 entries sorted by increasing sizes and an entry with a “0” shift is an empty entry and a terminator:

```
struct kvm_ppc_one_page_size {
    __u32 page_shift;      /* Page shift (or 0) */
    __u32 pte_enc;        /* Encoding in the HPTE (>>12) */
};
```

The “pte\_enc” field provides a value that can OR’ed into the hash PTE’s RPN field (ie, it needs to be shifted left by 12 to OR it into the hash PTE second double word).

## 4.75 KVM\_IRQFD

### Capability

KVM\_CAP\_IRQFD

### Architectures

x86 s390 arm arm64

### Type

vm ioctl

### Parameters

struct kvm\_irqfd (in)

### Returns

0 on success, -1 on error

Allows setting an eventfd to directly trigger a guest interrupt. `kvm_irqfd.fd` specifies the file descriptor to use as the eventfd and `kvm_irqfd.gsi` specifies the irqchip pin toggled by this event. When an event is triggered on the eventfd, an interrupt is injected into the guest using the specified gsi pin. The irqfd is removed using the `KVM_IRQFD_FLAG_DEASSIGN` flag, specifying both `kvm_irqfd.fd` and `kvm_irqfd.gsi`.

With `KVM_CAP_IRQFD_RESAMPLE`, `KVM_IRQFD` supports a de-assert and notify mechanism allowing emulation of level-triggered, irqfd-based interrupts. When `KVM_IRQFD_FLAG_RESAMPLE` is set the user must pass an additional eventfd in the `kvm_irqfd.resamplefd` field. When operating in resample mode, posting of an interrupt through `kvm_irqfd` asserts the specified gsi in the irqchip. When the irqchip is resampled, such as from an EOI, the gsi is de-asserted and the user is notified via `kvm_irqfd.resamplefd`. It is the user's responsibility to re-queue the interrupt if the device making use of it still requires service. Note that closing the resamplefd is not sufficient to disable the irqfd. The `KVM_IRQFD_FLAG_RESAMPLE` is only necessary on assignment and need not be specified with `KVM_IRQFD_FLAG_DEASSIGN`.

On arm/arm64, gsi routing being supported, the following can happen:

- in case no routing entry is associated to this gsi, injection fails
- in case the gsi is associated to an irqchip routing entry, `irqchip.pin + 32` corresponds to the injected SPI ID.
- in case the gsi is associated to an MSI routing entry, the MSI message and device ID are translated into an LPI (support restricted to GICv3 ITS in-kernel emulation).

#### **4.76 KVM\_PPC\_ALLOCATE\_HTAB**

##### **Capability**

`KVM_CAP_PPC_ALLOC_HTAB`

##### **Architectures**

powerpc

##### **Type**

vm ioctl

##### **Parameters**

Pointer to u32 containing hash table order (in/out)

##### **Returns**

0 on success, -1 on error

This requests the host kernel to allocate an MMU hash table for a guest using the PAPR paravirtualization interface. This only does anything if the kernel is configured to use the Book 3S HV style of virtualization. Otherwise the capability doesn't exist and the ioctl returns an `ENOTTY` error. The rest of this description assumes Book 3S HV.

There must be no vcpus running when this ioctl is called; if there are, it will do nothing and return an `EBUSY` error.

The parameter is a pointer to a 32-bit unsigned integer variable containing the order (log base 2) of the desired size of the hash table, which must be between 18 and 46. On successful return from the ioctl, the value will not be changed by the kernel.

If no hash table has been allocated when any vcpu is asked to run (with the `KVM_RUN` ioctl), the host kernel will allocate a default-sized hash table (16 MB).

If this ioctl is called when a hash table has already been allocated, with a different order from the existing hash table, the existing hash table will be freed and a new one allocated. If this is ioctl is called when a hash table has already been allocated of the same order as specified, the kernel will clear out the existing hash table (zero all HPTEs). In either case, if the guest is using the virtualized real-mode area (VRMA) facility, the kernel will re-create the VMRA HPTEs on the next KVM\_RUN of any vcpu.

## 4.77 KVM\_S390\_INTERRUPT

### Capability

basic

### Architectures

s390

### Type

vm ioctl, vcpu ioctl

### Parameters

struct kvm\_s390\_interrupt (in)

### Returns

0 on success, -1 on error

Allows to inject an interrupt to the guest. Interrupts can be floating (vm ioctl) or per cpu (vcpu ioctl), depending on the interrupt type.

Interrupt parameters are passed via kvm\_s390\_interrupt:

```
struct kvm_s390_interrupt {
    __u32 type;
    __u32 parm;
    __u64 parm64;
};
```

type can be one of the following:

### KVM\_S390\_SIGP\_STOP (vcpu)

- sigp stop; optional flags in parm

### KVM\_S390\_PROGRAM\_INT (vcpu)

- program check; code in parm

### KVM\_S390\_SIGP\_SET\_PREFIX (vcpu)

- sigp set prefix; prefix address in parm

### KVM\_S390\_RESTART (vcpu)

- restart

### KVM\_S390\_INT\_CLOCK\_COMP (vcpu)

- clock comparator interrupt

### KVM\_S390\_INT\_CPU\_TIMER (vcpu)

- CPU timer interrupt

**KVM\_S390\_INT\_VIRTIO (vm)**

- virtio external interrupt; external interrupt parameters in parm and parm64

**KVM\_S390\_INT\_SERVICE (vm)**

- sclp external interrupt; sclp parameter in parm

**KVM\_S390\_INT\_EMERGENCY (vcpu)**

- sigp emergency; source cpu in parm

**KVM\_S390\_INT\_EXTERNAL\_CALL (vcpu)**

- sigp external call; source cpu in parm

**KVM\_S390\_INT\_IO(ai,cssid,ssid,schid) (vm)**

- compound value to indicate an I/O interrupt (ai - adapter interrupt; cssid,ssid,schid - subchannel); I/O interruption parameters in parm (subchannel) and parm64 (intparm, interruption subclass)

**KVM\_S390\_MCHK (vm, vcpu)**

- machine check interrupt; cr 14 bits in parm, machine check interrupt code in parm64 (note that machine checks needing further payload are not supported by this ioctl)

This is an asynchronous vcpu ioctl and can be invoked from any thread.

**4.78 KVM\_PPC\_GET\_HTAB\_FD****Capability**

KVM\_CAP\_PPC\_HTAB\_FD

**Architectures**

powerpc

**Type**

vm ioctl

**Parameters**

Pointer to struct kvm\_get\_htab\_fd (in)

**Returns**

file descriptor number ( $\geq 0$ ) on success, -1 on error

This returns a file descriptor that can be used either to read out the entries in the guest's hashed page table (HPT), or to write entries to initialize the HPT. The returned fd can only be written to if the KVM\_GET\_HTAB\_WRITE bit is set in the flags field of the argument, and can only be read if that bit is clear. The argument struct looks like this:

```
/* For KVM_PPC_GET_HTAB_FD */
struct kvm_get_htab_fd {
    __u64    flags;
```

(continues on next page)

(continued from previous page)

```

        __u64    start_index;
        __u64    reserved[2];
};

/* Values for kvm_get_htab_fd.flags */
#define KVM_GET_HTAB_BOLTED_ONLY    ((__u64)0x1)
#define KVM_GET_HTAB_WRITE         ((__u64)0x2)

```

The ‘start\_index’ field gives the index in the HPT of the entry at which to start reading. It is ignored when writing.

Reads on the fd will initially supply information about all “interesting” HPT entries. Interesting entries are those with the bolted bit set, if the KVM\_GET\_HTAB\_BOLTED\_ONLY bit is set, otherwise all entries. When the end of the HPT is reached, the read() will return. If read() is called again on the fd, it will start again from the beginning of the HPT, but will only return HPT entries that have changed since they were last read.

Data read or written is structured as a header (8 bytes) followed by a series of valid HPT entries (16 bytes) each. The header indicates how many valid HPT entries there are and how many invalid entries follow the valid entries. The invalid entries are not represented explicitly in the stream. The header format is:

```

struct kvm_get_htab_header {
    __u32    index;
    __u16    n_valid;
    __u16    n_invalid;
};

```

Writes to the fd create HPT entries starting at the index given in the header; first ‘n\_valid’ valid entries with contents from the data written, then ‘n\_invalid’ invalid entries, invalidating any previously valid entries found.

## 4.79 KVM\_CREATE\_DEVICE

### Capability

KVM\_CAP\_DEVICE\_CTRL

### Type

vm ioctl

### Parameters

struct kvm\_create\_device (in/out)

### Returns

0 on success, -1 on error

Errors:

EN-ODEV	The device type is unknown or unsupported
EEX-IST	Device already created, and this type of device may not be instantiated multiple times

Other error conditions may be defined by individual device types or have their standard meanings.

Creates an emulated device in the kernel. The file descriptor returned in `fd` can be used with `KVM_SET/GET/HAS_DEVICE_ATTR`.

If the `KVM_CREATE_DEVICE_TEST` flag is set, only test whether the device type is supported (not necessarily whether it can be created in the current vm).

Individual devices should not define flags. Attributes should be used for specifying any behavior that is not implied by the device type number.

```
struct kvm_create_device {
    __u32  type;    /* in: KVM_DEV_TYPE_XXX */
    __u32  fd;      /* out: device handle */
    __u32  flags;   /* in: KVM_CREATE_DEVICE_XXX */
};
```

## 4.80 KVM\_SET\_DEVICE\_ATTR/KVM\_GET\_DEVICE\_ATTR

### Capability

`KVM_CAP_DEVICE_CTRL`, `KVM_CAP_VM_ATTRIBUTES` for vm device, `KVM_CAP_VCPU_ATTRIBUTES` for vcpu device

### Type

device ioctl, vm ioctl, vcpu ioctl

### Parameters

`struct kvm_device_attr`

### Returns

0 on success, -1 on error

Errors:

ENX	The group or attribute is unknown/unsupported for this device or hardware support is missing.
EPE	The attribute cannot (currently) be accessed this way (e.g. read-only attribute, or attribute that only makes sense when the device is in a different state)

Other error conditions may be defined by individual device types.

Gets/sets a specified piece of device configuration and/or state. The semantics are device-specific. See individual device documentation in the “devices” directory. As with `ONE_REG`, the size of the data transferred is defined by the particular attribute.

```

struct kvm_device_attr {
    __u32    flags;           /* no flags currently defined */
    __u32    group;          /* device-defined */
    __u64    attr;           /* group-defined */
    __u64    addr;           /* userspace address of attr data */
};

```

#### 4.81 KVM\_HAS\_DEVICE\_ATTR

##### Capability

KVM\_CAP\_DEVICE\_CTRL, KVM\_CAP\_VM\_ATTRIBUTES for vm device, KVM\_CAP\_VCPU\_ATTRIBUTES for vcpu device

##### Type

device ioctl, vm ioctl, vcpu ioctl

##### Parameters

struct kvm\_device\_attr

##### Returns

0 on success, -1 on error

Errors:

ENXIO The group or attribute is unknown/unsupported for this device or hardware support is missing.

Tests whether a device supports a particular attribute. A successful return indicates the attribute is implemented. It does not necessarily indicate that the attribute can be read or written in the device's current state. "addr" is ignored.

#### 4.82 KVM\_ARM\_VCPU\_INIT

##### Capability

basic

##### Architectures

arm, arm64

##### Type

vcpu ioctl

##### Parameters

struct kvm\_vcpu\_init (in)

##### Returns

0 on success; -1 on error

Errors:

EINVAL	the target is unknown, or the combination of features is invalid.
ENOENT	a features bit specified is unknown.

This tells KVM what type of CPU to present to the guest, and what optional features it should have. This will cause a reset of the cpu registers to their initial values. If this is not called, KVM\_RUN will return ENOEXEC for that vcpu.

Note that because some registers reflect machine topology, all vcpus should be created before this ioctl is invoked.

Userspace can call this function multiple times for a given vcpu, including after the vcpu has been run. This will reset the vcpu to its initial state. All calls to this function after the initial call must use the same target and same set of feature flags, otherwise EINVAL will be returned.

Possible features:

- KVM\_ARM\_VCPU\_POWER\_OFF: Starts the CPU in a power-off state. Depends on KVM\_CAP\_ARM\_PSCI. If not set, the CPU will be powered on and execute guest code when KVM\_RUN is called.
- KVM\_ARM\_VCPU\_EL1\_32BIT: Starts the CPU in a 32bit mode. Depends on KVM\_CAP\_ARM\_EL1\_32BIT (arm64 only).
- KVM\_ARM\_VCPU\_PSCI\_0\_2: Emulate PSCI v0.2 (or a future revision backward compatible with v0.2) for the CPU. Depends on KVM\_CAP\_ARM\_PSCI\_0\_2.
- KVM\_ARM\_VCPU\_PMU\_V3: Emulate PMUv3 for the CPU. Depends on KVM\_CAP\_ARM\_PMU\_V3.
- KVM\_ARM\_VCPU\_PTRAUTH\_ADDRESS: Enables Address Pointer authentication for arm64 only. Depends on KVM\_CAP\_ARM\_PTRAUTH\_ADDRESS. If KVM\_CAP\_ARM\_PTRAUTH\_ADDRESS and KVM\_CAP\_ARM\_PTRAUTH\_GENERIC are both present, then both KVM\_ARM\_VCPU\_PTRAUTH\_ADDRESS and KVM\_ARM\_VCPU\_PTRAUTH\_GENERIC must be requested or neither must be requested.
- KVM\_ARM\_VCPU\_PTRAUTH\_GENERIC: Enables Generic Pointer authentication for arm64 only. Depends on KVM\_CAP\_ARM\_PTRAUTH\_GENERIC. If KVM\_CAP\_ARM\_PTRAUTH\_ADDRESS and KVM\_CAP\_ARM\_PTRAUTH\_GENERIC are both present, then both KVM\_ARM\_VCPU\_PTRAUTH\_ADDRESS and KVM\_ARM\_VCPU\_PTRAUTH\_GENERIC must be requested or neither must be requested.
- KVM\_ARM\_VCPU\_SVE: Enables SVE for the CPU (arm64 only). Depends on KVM\_CAP\_ARM\_SVE. Requires KVM\_ARM\_VCPU\_FINALIZE(KVM\_ARM\_VCPU\_SVE):
  - After KVM\_ARM\_VCPU\_INIT:
    - \* KVM\_REG\_ARM64\_SVE\_VLS may be read using KVM\_GET\_ONE\_REG: the initial value of this pseudo-register



indicates the best set of vector lengths possible for a vcpu on this host.

- Before `KVM_ARM_VCPU_FINALIZE(KVM_ARM_VCPU_SVE)`:
  - \* `KVM_RUN` and `KVM_GET_REG_LIST` are not available;
  - \* `KVM_GET_ONE_REG` and `KVM_SET_ONE_REG` cannot be used to access the scalable architectural SVE registers `KVM_REG_ARM64_SVE_ZREG()`, `KVM_REG_ARM64_SVE_PREG()` or `KVM_REG_ARM64_SVE_FFR`;
  - \* `KVM_REG_ARM64_SVE_VLS` may optionally be written using `KVM_SET_ONE_REG`, to modify the set of vector lengths available for the vcpu.
- After `KVM_ARM_VCPU_FINALIZE(KVM_ARM_VCPU_SVE)`:
  - \* the `KVM_REG_ARM64_SVE_VLS` pseudo-register is immutable, and can no longer be written using `KVM_SET_ONE_REG`.

### 4.83 KVM\_ARM\_PREFERRED\_TARGET

#### Capability

basic

#### Architectures

arm, arm64

#### Type

vm ioctl

#### Parameters

struct `kvm_vcpu_init` (out)

#### Returns

0 on success; -1 on error

Errors:

ENODEV    no preferred target available for the host
--

This queries KVM for preferred CPU target type which can be emulated by KVM on underlying host.

The ioctl returns struct `kvm_vcpu_init` instance containing information about preferred CPU target type and recommended features for it. The `kvm_vcpu_init->features` bitmap returned will have feature bits set if the preferred target recommends setting these features, but this is not mandatory.

The information returned by this ioctl can be used to prepare an instance of struct `kvm_vcpu_init` for `KVM_ARM_VCPU_INIT` ioctl which will result in VCPU matching underlying host.

## 4.84 KVM\_GET\_REG\_LIST

**Capability**

basic

**Architectures**

arm, arm64, mips

**Type**

vcpu ioctl

**Parameters**

struct kvm\_reg\_list (in/out)

**Returns**

0 on success; -1 on error

Errors:

E2BI the reg index list is too big to fit in the array specified by the user (the number required will be written into n).

```
struct kvm_reg_list {
    __u64 n; /* number of registers in reg[] */
    __u64 reg[0];
};
```

This ioctl returns the guest registers that are supported for the KVM\_GET\_ONE\_REG/KVM\_SET\_ONE\_REG calls.

## 4.85 KVM\_ARM\_SET\_DEVICE\_ADDR (deprecated)

**Capability**

KVM\_CAP\_ARM\_SET\_DEVICE\_ADDR

**Architectures**

arm, arm64

**Type**

vm ioctl

**Parameters**

struct kvm\_arm\_device\_address (in)

**Returns**

0 on success, -1 on error

Errors:

ENODEV	The device id is unknown
ENXIO	Device not supported on current system
EEXIST	Address already set
E2BIG	Address outside guest physical address space
EBUSY	Address overlaps with other device range

```
struct kvm_arm_device_addr {
    __u64 id;
    __u64 addr;
};
```

Specify a device address in the guest's physical address space where guests can access emulated or directly exposed devices, which the host kernel needs to know about. The id field is an architecture specific identifier for a specific device.

ARM/arm64 divides the id field into two parts, a device id and an address type id specific to the individual device:

bits:	63	...	32	31	...	16	15	...	0
field:		0x00000000			device id			addr type id	

ARM/arm64 currently only require this when using the in-kernel GIC support for the hardware VGIC features, using `KVM_ARM_DEVICE_VGIC_V2` as the device id. When setting the base address for the guest's mapping of the VGIC virtual CPU and distributor interface, the ioctl must be called after calling `KVM_CREATE_IRQCHIP`, but before calling `KVM_RUN` on any of the VCPUs. Calling this ioctl twice for any of the base addresses will return `-EEXIST`.

Note, this IOCTL is deprecated and the more flexible `SET/GET_DEVICE_ATTR` API should be used instead.

## 4.86 KVM\_PPC\_RTAS\_DEFINE\_TOKEN

### Capability

`KVM_CAP_PPC_RTAS`

### Architectures

ppc

### Type

vm ioctl

### Parameters

`struct kvm_rtas_token_args`

### Returns

0 on success, -1 on error

Defines a token value for a RTAS (Run Time Abstraction Services) service in order to allow it to be handled in the kernel. The argument struct gives the name of the service, which must be the name of a service that has a kernel-side implementation. If the token value is non-zero, it will be associated with that service, and subsequent RTAS calls by the guest specifying that token will be handled by the

kernel. If the token value is 0, then any token associated with the service will be forgotten, and subsequent RTAS calls by the guest for that service will be passed to userspace to be handled.

## **4.87 KVM\_SET\_GUEST\_DEBUG**

### **Capability**

KVM\_CAP\_SET\_GUEST\_DEBUG

### **Architectures**

x86, s390, ppc, arm64

### **Type**

vcpu ioctl

### **Parameters**

struct kvm\_guest\_debug (in)

### **Returns**

0 on success; -1 on error

```
struct kvm_guest_debug {
    __u32 control;
    __u32 pad;
    struct kvm_guest_debug_arch arch;
};
```

Set up the processor specific debug registers and configure vcpu for handling guest debug events. There are two parts to the structure, the first a control bitfield indicates the type of debug events to handle when running. Common control bits are:

- KVM\_GUESTDBG\_ENABLE: guest debugging is enabled
- KVM\_GUESTDBG\_SINGLESTEP: the next run should single-step

The top 16 bits of the control field are architecture specific control flags which can include the following:

- KVM\_GUESTDBG\_USE\_SW\_BP: using software breakpoints [x86, arm64]
- KVM\_GUESTDBG\_USE\_HW\_BP: using hardware breakpoints [x86, s390, arm64]
- KVM\_GUESTDBG\_INJECT\_DB: inject DB type exception [x86]
- KVM\_GUESTDBG\_INJECT\_BP: inject BP type exception [x86]
- KVM\_GUESTDBG\_EXIT\_PENDING: trigger an immediate guest exit [s390]

For example KVM\_GUESTDBG\_USE\_SW\_BP indicates that software breakpoints are enabled in memory so we need to ensure breakpoint exceptions are correctly trapped and the KVM run loop exits at the breakpoint and not running off into the normal guest vector. For KVM\_GUESTDBG\_USE\_HW\_BP we need to ensure the guest vCPUs architecture specific registers are updated to the correct (supplied) values.

The second part of the structure is architecture specific and typically contains a set of debug registers.

For arm64 the number of debug registers is implementation defined and can be determined by querying the KVM\_CAP\_GUEST\_DEBUG\_HW\_BPS and KVM\_CAP\_GUEST\_DEBUG\_HW\_WPS capabilities which return a positive number indicating the number of supported registers.

For ppc, the KVM\_CAP\_PPC\_GUEST\_DEBUG\_SSTEP capability indicates whether the single-step debug event (KVM\_GUESTDBG\_SINGLESTEP) is supported.

When debug events exit the main run loop with the reason KVM\_EXIT\_DEBUG with the kvm\_debug\_exit\_arch part of the kvm\_run structure containing architecture specific debug information.

#### 4.88 KVM\_GET\_EMULATED\_CPUID

##### Capability

KVM\_CAP\_EXT\_EMUL\_CPUID

##### Architectures

x86

##### Type

system ioctl

##### Parameters

struct kvm\_cpuid2 (in/out)

##### Returns

0 on success, -1 on error

```
struct kvm_cpuid2 {
    __u32 nent;
    __u32 flags;
    struct kvm_cpuid_entry2 entries[0];
};
```

The member 'flags' is used for passing flags from userspace.

```
#define KVM_CPUID_FLAG_SIGNIFCANT_INDEX        BIT(0)
#define KVM_CPUID_FLAG_STATEFUL_FUNC          BIT(1) /* deprecated
↳ */
#define KVM_CPUID_FLAG_STATE_READ_NEXT        BIT(2) /*
↳ deprecated */

struct kvm_cpuid_entry2 {
    __u32 function;
    __u32 index;
    __u32 flags;
    __u32 eax;
    __u32 ebx;
    __u32 ecx;
    __u32 edx;
```

(continues on next page)

(continued from previous page)

```
    __u32 padding[3];  
};
```

This ioctl returns x86 cpuid features which are emulated by kvm. Userspace can use the information returned by this ioctl to query which features are emulated by kvm instead of being present natively.

Userspace invokes KVM\_GET\_EMULATED\_CPUID by passing a kvm\_cpuid2 structure with the 'nent' field indicating the number of entries in the variable-size array 'entries'. If the number of entries is too low to describe the cpu capabilities, an error (E2BIG) is returned. If the number is too high, the 'nent' field is adjusted and an error (ENOMEM) is returned. If the number is just right, the 'nent' field is adjusted to the number of valid entries in the 'entries' array, which is then filled.

The entries returned are the set CPUID bits of the respective features which kvm emulates, as returned by the CPUID instruction, with unknown or unsupported feature bits cleared.

Features like x2apic, for example, may not be present in the host cpu but are exposed by kvm in KVM\_GET\_SUPPORTED\_CPUID because they can be emulated efficiently and thus not included here.

The fields in each entry are defined as follows:

**function:**

the eax value used to obtain the entry

**index:**

the ecx value used to obtain the entry (for entries that are affected by ecx)

**flags:**

an OR of zero or more of the following:

**KVM\_CPUID\_FLAG\_SIGNIFICANT\_INDEX:**

if the index field is valid

eax, ebx, ecx, edx:

the values returned by the cpuid instruction for this function/index combination

## 4.89 KVM\_S390\_MEM\_OP

**Capability**

KVM\_CAP\_S390\_MEM\_OP

**Architectures**

s390

**Type**

vcpu ioctl

**Parameters**

struct kvm\_s390\_mem\_op (in)

**Returns**

= 0 on success, < 0 on generic error (e.g. -EFAULT or -ENOMEM),  
 > 0 if an exception occurred while walking the page tables

Read or write data from/to the logical (virtual) memory of a VCPU.

Parameters are specified via the following structure:

```
struct kvm_s390_mem_op {
    __u64 gaddr;          /* the guest address */
    __u64 flags;          /* flags */
    __u32 size;           /* amount of bytes */
    __u32 op;            /* type of operation */
    __u64 buf;            /* buffer in userspace */
    __u8 ar;              /* the access register number */
    __u8 reserved[31];    /* should be set to 0 */
};
```

The type of operation is specified in the “op” field. It is either `KVM_S390_MEMOP_LOGICAL_READ` for reading from logical memory space or `KVM_S390_MEMOP_LOGICAL_WRITE` for writing to logical memory space. The `KVM_S390_MEMOP_F_CHECK_ONLY` flag can be set in the “flags” field to check whether the corresponding memory access would create an access exception (without touching the data in the memory at the destination). In case an access exception occurred while walking the MMU tables of the guest, the ioctl returns a positive error number to indicate the type of exception. This exception is also raised directly at the corresponding VCPU if the flag `KVM_S390_MEMOP_F_INJECT_EXCEPTION` is set in the “flags” field.

The start address of the memory region has to be specified in the “gaddr” field, and the length of the region in the “size” field (which must not be 0). The maximum value for “size” can be obtained by checking the `KVM_CAP_S390_MEM_OP` capability. “buf” is the buffer supplied by the userspace application where the read data should be written to for `KVM_S390_MEMOP_LOGICAL_READ`, or where the data that should be written is stored for a `KVM_S390_MEMOP_LOGICAL_WRITE`. When `KVM_S390_MEMOP_F_CHECK_ONLY` is specified, “buf” is unused and can be NULL. “ar” designates the access register number to be used; the valid range is 0..15.

The “reserved” field is meant for future extensions. It is not used by KVM with the currently defined set of flags.

**4.90 KVM\_S390\_GET\_SKEYS****Capability**

`KVM_CAP_S390_SKEYS`

**Architectures**

s390

**Type**

vm ioctl

**Parameters**

struct kvm\_s390\_keys

**Returns**

0 on success, KVM\_S390\_GET\_KEYS\_NONE if guest is not using storage keys, negative value on error

This ioctl is used to get guest storage key values on the s390 architecture. The ioctl takes parameters via the kvm\_s390\_keys struct:

```
struct kvm_s390_keys {
    __u64 start_gfn;
    __u64 count;
    __u64 skeydata_addr;
    __u32 flags;
    __u32 reserved[9];
};
```

The start\_gfn field is the number of the first guest frame whose storage keys you want to get.

The count field is the number of consecutive frames (starting from start\_gfn) whose storage keys to get. The count field must be at least 1 and the maximum allowed value is defined as KVM\_S390\_SKEYS\_ALLOC\_MAX. Values outside this range will cause the ioctl to return -EINVAL.

The skeydata\_addr field is the address to a buffer large enough to hold count bytes. This buffer will be filled with storage key data by the ioctl.

## **4.91 KVM\_S390\_SET\_SKEYS**

**Capability**

KVM\_CAP\_S390\_SKEYS

**Architectures**

s390

**Type**

vm ioctl

**Parameters**

struct kvm\_s390\_keys

**Returns**

0 on success, negative value on error

This ioctl is used to set guest storage key values on the s390 architecture. The ioctl takes parameters via the kvm\_s390\_keys struct. See section on KVM\_S390\_GET\_SKEYS for struct definition.

The start\_gfn field is the number of the first guest frame whose storage keys you want to set.

The count field is the number of consecutive frames (starting from start\_gfn) whose storage keys to get. The count field must be at least 1 and the maximum allowed



value is defined as `KVM_S390_SKEYS_ALLOC_MAX`. Values outside this range will cause the `ioctl` to return `-EINVAL`.

The `skeydata_addr` field is the address to a buffer containing count bytes of storage keys. Each byte in the buffer will be set as the storage key for a single frame starting at `start_gfn` for count frames.

Note: If any architecturally invalid key value is found in the given data then the `ioctl` will return `-EINVAL`.

## 4.92 KVM\_S390\_IRQ

### Capability

`KVM_CAP_S390_INJECT_IRQ`

### Architectures

s390

### Type

`vcpu ioctl`

### Parameters

`struct kvm_s390_irq` (in)

### Returns

0 on success, -1 on error

Errors:

<p>EIN interrupt type is invalid type is <code>KVM_S390_SIGP_STOP</code> and flag parameter is invalid value, type is <code>KVM_S390_INT_EXTERNAL_CALL</code> and code is bigger than the maximum of VCPUs</p> <p>EBU type is <code>KVM_S390_SIGP_SET_PREFIX</code> and vcpu is not stopped, type is <code>KVM_S390_SIGP_STOP</code> and a stop irq is already pending, type is <code>KVM_S390_INT_EXTERNAL_CALL</code> and an external call interrupt is already pending</p>
---

Allows to inject an interrupt to the guest.

Using `struct kvm_s390_irq` as a parameter allows to inject additional payload which is not possible via `KVM_S390_INTERRUPT`.

Interrupt parameters are passed via `kvm_s390_irq`:

```
struct kvm_s390_irq {
    __u64 type;
    union {
        struct kvm_s390_io_info io;
        struct kvm_s390_ext_info ext;
        struct kvm_s390_pgm_info pgm;
        struct kvm_s390_emerg_info emerg;
        struct kvm_s390_extcall_info extcall;
    };
};
```

(continues on next page)

(continued from previous page)

```

        struct kvm_s390_prefix_info prefix;
        struct kvm_s390_stop_info stop;
        struct kvm_s390_mchk_info mchk;
        char reserved[64];
    } u;
};

```

type can be one of the following:

- KVM\_S390\_SIGP\_STOP - sigp stop; parameter in .stop
- KVM\_S390\_PROGRAM\_INT - program check; parameters in .pgm
- KVM\_S390\_SIGP\_SET\_PREFIX - sigp set prefix; parameters in .prefix
- KVM\_S390\_RESTART - restart; no parameters
- KVM\_S390\_INT\_CLOCK\_COMP - clock comparator interrupt; no parameters
- KVM\_S390\_INT\_CPU\_TIMER - CPU timer interrupt; no parameters
- KVM\_S390\_INT\_EMERGENCY - sigp emergency; parameters in .emerg
- KVM\_S390\_INT\_EXTERNAL\_CALL - sigp external call; parameters in .extcall
- KVM\_S390\_MCHK - machine check interrupt; parameters in .mchk

This is an asynchronous vcpu ioctl and can be invoked from any thread.

#### 4.94 KVM\_S390\_GET\_IRQ\_STATE

##### Capability

KVM\_CAP\_S390\_IRQ\_STATE

##### Architectures

s390

##### Type

vcpu ioctl

##### Parameters

struct kvm\_s390\_irq\_state (out)

##### Returns

>= number of bytes copied into buffer, -EINVAL if buffer size is 0,  
 -ENOBUFS if buffer size is too small to fit all pending interrupts,  
 -EFAULT if the buffer address was invalid

This ioctl allows userspace to retrieve the complete state of all currently pending interrupts in a single buffer. Use cases include migration and introspection. The parameter structure contains the address of a userspace buffer and its length:

```

struct kvm_s390_irq_state {
    __u64 buf;
    __u32 flags;          /* will stay unused for compatibility_
↪ reasons */

```

(continues on next page)

(continued from previous page)

```

    __u32 len;
    __u32 reserved[4]; /* will stay unused for compatibility
↳ reasons */
};

```

Userspace passes in the above struct and for each pending interrupt a struct `kvm_s390_irq` is copied to the provided buffer.

The structure contains a flags and a reserved field for future extensions. As the kernel never checked for `flags == 0` and QEMU never pre-zeroed flags and reserved, these fields can not be used in the future without breaking compatibility.

If `-ENOBUFFS` is returned the buffer provided was too small and userspace may retry with a bigger buffer.

#### 4.95 KVM\_S390\_SET\_IRQ\_STATE

##### Capability

`KVM_CAP_S390_IRQ_STATE`

##### Architectures

s390

##### Type

`vcpu ioctl`

##### Parameters

`struct kvm_s390_irq_state (in)`

##### Returns

0 on success, `-EFAULT` if the buffer address was invalid, `-EINVAL` for an invalid buffer length (see below), `-EBUSY` if there were already interrupts pending, errors occurring when actually injecting the interrupt. See `KVM_S390_IRQ`.

This `ioctl` allows userspace to set the complete state of all cpu-local interrupts currently pending for the vcpu. It is intended for restoring interrupt state after a migration. The input parameter is a userspace buffer containing a struct `kvm_s390_irq_state`:

```

struct kvm_s390_irq_state {
    __u64 buf;
    __u32 flags;          /* will stay unused for compatibility
↳ reasons */
    __u32 len;
    __u32 reserved[4]; /* will stay unused for compatibility
↳ reasons */
};

```

The restrictions for flags and reserved apply as well. (see `KVM_S390_GET_IRQ_STATE`)

The userspace memory referenced by `buf` contains a struct `kvm_s390_irq` for each interrupt to be injected into the guest. If one of the interrupts could not be injected

for some reason the ioctl aborts.

len must be a multiple of sizeof(struct kvm\_s390\_irq). It must be > 0 and it must not exceed (max\_vcpus + 32) \* sizeof(struct kvm\_s390\_irq), which is the maximum number of possibly pending cpu-local interrupts.

#### **4.96 KVM\_SMI**

**Capability**

KVM\_CAP\_X86\_SMM

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

none

**Returns**

0 on success, -1 on error

Queues an SMI on the thread's vcpu.

#### **4.97 KVM\_CAP\_PPC\_MULTITCE**

**Capability**

KVM\_CAP\_PPC\_MULTITCE

**Architectures**

ppc

**Type**

vm

This capability means the kernel is capable of handling hypercalls H\_PUT\_TCE\_INDIRECT and H\_STUFF\_TCE without passing those into the user space. This significantly accelerates DMA operations for PPC KVM guests. User space should expect that its handlers for these hypercalls are not going to be called if user space previously registered LIOBN in KVM (via KVM\_CREATE\_SPAPR\_TCE or similar calls).

In order to enable H\_PUT\_TCE\_INDIRECT and H\_STUFF\_TCE use in the guest, user space might have to advertise it for the guest. For example, IBM pSeries (sPAPR) guest starts using them if "hcall-multi-tce" is present in the "ibm,hypertas-functions" device-tree property.

The hypercalls mentioned above may or may not be processed successfully in the kernel based fast path. If they can not be handled by the kernel, they will get passed on to user space. So user space still has to have an implementation for these despite the in kernel acceleration.

This capability is always enabled.

## 4.98 KVM\_CREATE\_SPAPR\_TCE\_64

### Capability

KVM\_CAP\_SPAPR\_TCE\_64

### Architectures

powerpc

### Type

vm ioctl

### Parameters

struct kvm\_create\_spapr\_tce\_64 (in)

### Returns

file descriptor for manipulating the created TCE table

This is an extension for KVM\_CAP\_SPAPR\_TCE which only supports 32bit windows, described in 4.62 KVM\_CREATE\_SPAPR\_TCE

This capability uses extended struct in ioctl interface:

```
/* for KVM_CAP_SPAPR_TCE_64 */
struct kvm_create_spapr_tce_64 {
    __u64 liobn;
    __u32 page_shift;
    __u32 flags;
    __u64 offset;    /* in pages */
    __u64 size;      /* in pages */
};
```

The aim of extension is to support an additional bigger DMA window with a variable page size. KVM\_CREATE\_SPAPR\_TCE\_64 receives a 64bit window size, an IOMMU page shift and a bus offset of the corresponding DMA window, @size and @offset are numbers of IOMMU pages.

@flags are not used at the moment.

The rest of functionality is identical to KVM\_CREATE\_SPAPR\_TCE.

## 4.99 KVM\_REINJECT\_CONTROL

### Capability

KVM\_CAP\_REINJECT\_CONTROL

### Architectures

x86

### Type

vm ioctl

### Parameters

struct kvm\_reinject\_control (in)

### Returns

0 on success, -EFAULT if struct kvm\_reinject\_control cannot be read,

-ENXIO if KVM\_CREATE\_PIT or KVM\_CREATE\_PIT2 didn't succeed earlier.

i8254 (PIT) has two modes, reinject and !reinject. The default is reinject, where KVM queues elapsed i8254 ticks and monitors completion of interrupt from vector(s) that i8254 injects. Reinject mode dequeues a tick and injects its interrupt whenever there isn't a pending interrupt from i8254. !reinject mode injects an interrupt as soon as a tick arrives.

```
struct kvm_reinject_control {
    __u8 pit_reinject;
    __u8 reserved[31];
};
```

pit\_reinject = 0 (!reinject mode) is recommended, unless running an old operating system that uses the PIT for timing (e.g. Linux 2.4.x).

#### **4.100 KVM\_PPC\_CONFIGURE\_V3\_MMU**

##### **Capability**

KVM\_CAP\_PPC\_RADIX\_MMU or KVM\_CAP\_PPC\_HASH\_MMU\_V3

##### **Architectures**

ppc

##### **Type**

vm ioctl

##### **Parameters**

struct kvm\_ppc\_mmuv3\_cfg (in)

##### **Returns**

0 on success, -EFAULT if struct kvm\_ppc\_mmuv3\_cfg cannot be read, -EINVAL if the configuration is invalid

This ioctl controls whether the guest will use radix or HPT (hashed page table) translation, and sets the pointer to the process table for the guest.

```
struct kvm_ppc_mmuv3_cfg {
    __u64 flags;
    __u64 process_table;
};
```

There are two bits that can be set in flags; KVM\_PPC\_MMUV3\_RADIX and KVM\_PPC\_MMUV3\_GTSE. KVM\_PPC\_MMUV3\_RADIX, if set, configures the guest to use radix tree translation, and if clear, to use HPT translation. KVM\_PPC\_MMUV3\_GTSE, if set and if KVM permits it, configures the guest to be able to use the global TLB and SLB invalidation instructions; if clear, the guest may not use these instructions.

The process\_table field specifies the address and size of the guest process table, which is in the guest's space. This field is formatted as the second doubleword of the partition table entry, as defined in the Power ISA V3.00, Book III section 5.7.6.1.

### 4.101 KVM\_PPC\_GET\_RMMU\_INFO

**Capability**

KVM\_CAP\_PPC\_RADIX\_MMU

**Architectures**

ppc

**Type**

vm ioctl

**Parameters**

struct kvm\_ppc\_rmmu\_info (out)

**Returns**

0 on success, -EFAULT if struct kvm\_ppc\_rmmu\_info cannot be written, -EINVAL if no useful information can be returned

This ioctl returns a structure containing two things: (a) a list containing supported radix tree geometries, and (b) a list that maps page sizes to put in the “AP” (actual page size) field for the tlbie (TLB invalidate entry) instruction.

```
struct kvm_ppc_rmmu_info {
    struct kvm_ppc_radix_geom {
        __u8    page_shift;
        __u8    level_bits[4];
        __u8    pad[3];
    }          geometries[8];
    __u32      ap_encodings[8];
};
```

The geometries[] field gives up to 8 supported geometries for the radix page table, in terms of the log base 2 of the smallest page size, and the number of bits indexed at each level of the tree, from the PTE level up to the PGD level in that order. Any unused entries will have 0 in the page\_shift field.

The ap\_encodings gives the supported page sizes and their AP field encodings, encoded with the AP value in the top 3 bits and the log base 2 of the page size in the bottom 6 bits.

### 4.102 KVM\_PPC\_RESIZE\_HPT\_PREPARE

**Capability**

KVM\_CAP\_SPAPR\_RESIZE\_HPT

**Architectures**

powerpc

**Type**

vm ioctl

**Parameters**

struct kvm\_ppc\_resize\_hpt (in)

**Returns**

0 on successful completion, >0 if a new HPT is being prepared, the

value is an estimated number of milliseconds until preparation is complete, -EFAULT if struct kvm\_reinject\_control cannot be read, -EINVAL if the supplied shift or flags are invalid, -ENOMEM if unable to allocate the new HPT, -ENOSPC if there was a hash collision

```
struct kvm_ppc_rmmu_info {
    struct kvm_ppc_radix_geom {
        __u8    page_shift;
        __u8    level_bits[4];
        __u8    pad[3];
    }          geometries[8];
    __u32      ap_encodings[8];
};
```

The geometries[] field gives up to 8 supported geometries for the radix page table, in terms of the log base 2 of the smallest page size, and the number of bits indexed at each level of the tree, from the PTE level up to the PGD level in that order. Any unused entries will have 0 in the page\_shift field.

The ap\_encodings gives the supported page sizes and their AP field encodings, encoded with the AP value in the top 3 bits and the log base 2 of the page size in the bottom 6 bits.

#### 4.102 KVM\_PPC\_RESIZE\_HPT\_PREPARE

##### Capability

KVM\_CAP\_SPAPR\_RESIZE\_HPT

##### Architectures

powerpc

##### Type

vm ioctl

##### Parameters

struct kvm\_ppc\_resize\_hpt (in)

##### Returns

0 on successful completion, >0 if a new HPT is being prepared, the value is an estimated number of milliseconds until preparation is complete, -EFAULT if struct kvm\_reinject\_control cannot be read, -EINVAL if the supplied shift or flags are invalid, when moving existing HPT entries to the new HPT, -EIO on other error conditions

Used to implement the PAPR extension for runtime resizing of a guest's Hashed Page Table (HPT). Specifically this starts, stops or monitors the preparation of a new potential HPT for the guest, essentially implementing the H\_RESIZE\_HPT\_PREPARE hypercall.

If called with shift > 0 when there is no pending HPT for the guest, this begins preparation of a new pending HPT of size 2<sup>(shift)</sup> bytes. It then returns a positive integer with the estimated number of milliseconds until preparation is complete.

If called when there is a pending HPT whose size does not match that requested in the parameters, discards the existing pending HPT and creates a new one as



above.

If called when there is a pending HPT of the size requested, will:

- If preparation of the pending HPT is already complete, return 0
- If preparation of the pending HPT has failed, return an error code, then discard the pending HPT.
- If preparation of the pending HPT is still in progress, return an estimated number of milliseconds until preparation is complete.

If called with `shift == 0`, discards any currently pending HPT and returns 0 (i.e. cancels any in-progress preparation).

`flags` is reserved for future expansion, currently setting any bits in `flags` will result in an `-EINVAL`.

Normally this will be called repeatedly with the same parameters until it returns `<= 0`. The first call will initiate preparation, subsequent ones will monitor preparation until it completes or fails.

```
struct kvm_ppc_resize_hpt {
    __u64 flags;
    __u32 shift;
    __u32 pad;
};
```

#### 4.103 KVM\_PPC\_RESIZE\_HPT\_COMMIT

##### Capability

`KVM_CAP_SPAPR_RESIZE_HPT`

##### Architectures

powerpc

##### Type

vm ioctl

##### Parameters

`struct kvm_ppc_resize_hpt` (in)

##### Returns

0 on successful completion, `-EFAULT` if `struct kvm_reinject_control` cannot be read, `-EINVAL` if the supplied `shift` or `flags` are invalid, `-ENXIO` if there is no pending HPT, or the pending HPT doesn't have the requested size, `-EBUSY` if the pending HPT is not fully prepared, `-ENOSPC` if there was a hash collision when moving existing HPT entries to the new HPT, `-EIO` on other error conditions

Used to implement the PAPR extension for runtime resizing of a guest's Hashed Page Table (HPT). Specifically this requests that the guest be transferred to working with the new HPT, essentially implementing the `H_RESIZE_HPT_COMMIT` hypercall.

This should only be called after `KVM_PPC_RESIZE_HPT_PREPARE` has returned 0 with the same parameters. In other cases `KVM_PPC_RESIZE_HPT_COMMIT` will

return an error (usually -ENXIO or -EBUSY, though others may be possible if the preparation was started, but failed).

This will have undefined effects on the guest if it has not already placed itself in a quiescent state where no vcpu will make MMU enabled memory accesses.

On successful completion, the pending HPT will become the guest's active HPT and the previous HPT will be discarded.

On failure, the guest will still be operating on its previous HPT.

```
struct kvm_ppc_resize_hpt {
    __u64 flags;
    __u32 shift;
    __u32 pad;
};
```

### 4.104 KVM\_X86\_GET\_MCE\_CAP\_SUPPORTED

**Capability**

KVM\_CAP\_MCE

**Architectures**

x86

**Type**

system ioctl

**Parameters**

u64 mce\_cap (out)

**Returns**

0 on success, -1 on error

Returns supported MCE capabilities. The u64 mce\_cap parameter has the same format as the MSR\_IA32\_MCG\_CAP register. Supported capabilities will have the corresponding bits set.

### 4.105 KVM\_X86\_SETUP\_MCE

**Capability**

KVM\_CAP\_MCE

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

u64 mcg\_cap (in)

**Returns**

0 on success, -EFAULT if u64 mcg\_cap cannot be read, -EINVAL if the requested number of banks is invalid, -EINVAL if requested MCE capability is not supported.

Initializes MCE support for use. The u64 `mcg_cap` parameter has the same format as the `MSR_IA32_MCG_CAP` register and specifies which capabilities should be enabled. The maximum supported number of error-reporting banks can be retrieved when checking for `KVM_CAP_MCE`. The supported capabilities can be retrieved with `KVM_X86_GET_MCE_CAP_SUPPORTED`.

#### 4.106 KVM\_X86\_SET\_MCE

##### Capability

`KVM_CAP_MCE`

##### Architectures

x86

##### Type

`vcpu ioctl`

##### Parameters

`struct kvm_x86_mce` (in)

##### Returns

0 on success, `-EFAULT` if `struct kvm_x86_mce` cannot be read, `-EINVAL` if the bank number is invalid, `-EINVAL` if `VAL` bit is not set in status field.

Inject a machine check error (MCE) into the guest. The input parameter is:

```
struct kvm_x86_mce {
    __u64 status;
    __u64 addr;
    __u64 misc;
    __u64 mcg_status;
    __u8 bank;
    __u8 pad1[7];
    __u64 pad2[3];
};
```

If the MCE being reported is an uncorrected error, KVM will inject it as an MCE exception into the guest. If the guest `MCG_STATUS` register reports that an MCE is in progress, KVM causes an `KVM_EXIT_SHUTDOWN` vmexit.

Otherwise, if the MCE is a corrected error, KVM will just store it in the corresponding bank (provided this bank is not holding a previously reported uncorrected error).

## 4.107 KVM\_S390\_GET\_CMMA\_BITS

**Capability**

KVM\_CAP\_S390\_CMMA\_MIGRATION

**Architectures**

s390

**Type**

vm ioctl

**Parameters**

struct kvm\_s390\_cmma\_log (in, out)

**Returns**

0 on success, a negative value on error

Errors:

ENOMEM	not enough memory can be allocated to complete the task
EINVAL	if CMMMA is not enabled
EIN- VAL	if KVM_S390_CMMA_PEEK is not set but migration mode was not enabled
EIN- VAL	if KVM_S390_CMMA_PEEK is not set but dirty tracking has been disabled (and thus migration mode was automatically disabled)
EFAUL	if the userspace address is invalid or if no page table is present for the addresses (e.g. when using hugepages).

This ioctl is used to get the values of the CMMMA bits on the s390 architecture. It is meant to be used in two scenarios:

- During live migration to save the CMMMA values. Live migration needs to be enabled via the KVM\_REQ\_START\_MIGRATION VM property.
- To non-destructively peek at the CMMMA values, with the flag KVM\_S390\_CMMA\_PEEK set.

The ioctl takes parameters via the `kvm_s390_cmma_log` struct. The desired values are written to a buffer whose location is indicated via the “values” member in the `kvm_s390_cmma_log` struct. The values in the input struct are also updated as needed.

Each CMMMA value takes up one byte.

```
struct kvm_s390_cmma_log {
    __u64 start_gfn;
    __u32 count;
    __u32 flags;
    union {
        __u64 remaining;
        __u64 mask;
    };
    __u64 values;
};
```

`start_gfn` is the number of the first guest frame whose CMMA values are to be retrieved,

`count` is the length of the buffer in bytes,

`values` points to the buffer where the result will be written to.

If `count` is greater than `KVM_S390_SKEYS_MAX`, then it is considered to be `KVM_S390_SKEYS_MAX`. `KVM_S390_SKEYS_MAX` is re-used for consistency with other ioctls.

The result is written in the buffer pointed to by the field `values`, and the values of the input parameter are updated as follows.

Depending on the flags, different actions are performed. The only supported flag so far is `KVM_S390_CMMA_PEEK`.

The default behaviour if `KVM_S390_CMMA_PEEK` is not set is: `start_gfn` will indicate the first page frame whose CMMA bits were dirty. It is not necessarily the same as the one passed as input, as clean pages are skipped.

`count` will indicate the number of bytes actually written in the buffer. It can (and very often will) be smaller than the input value, since the buffer is only filled until 16 bytes of clean values are found (which are then not copied in the buffer). Since a CMMA migration block needs the base address and the length, for a total of 16 bytes, we will send back some clean data if there is some dirty data afterwards, as long as the size of the clean data does not exceed the size of the header. This allows to minimize the amount of data to be saved or transferred over the network at the expense of more roundtrips to userspace. The next invocation of the ioctl will skip over all the clean values, saving potentially more than just the 16 bytes we found.

If `KVM_S390_CMMA_PEEK` is set: the existing storage attributes are read even when not in migration mode, and no other action is performed;

the output `start_gfn` will be equal to the input `start_gfn`,

the output `count` will be equal to the input `count`, except if the end of memory has been reached.

In both cases: the field “remaining” will indicate the total number of dirty CMMA values still remaining, or 0 if `KVM_S390_CMMA_PEEK` is set and migration mode is not enabled.

`mask` is unused.

`values` points to the userspace buffer where the result will be stored.

## 4.108 KVM\_S390\_SET\_CMMA\_BITS

**Capability**

KVM\_CAP\_S390\_CMMA\_MIGRATION

**Architectures**

s390

**Type**

vm ioctl

**Parameters**

struct kvm\_s390\_cmma\_log (in)

**Returns**

0 on success, a negative value on error

This ioctl is used to set the values of the CMMA bits on the s390 architecture. It is meant to be used during live migration to restore the CMMA values, but there are no restrictions on its use. The ioctl takes parameters via the `kvm_s390_cmma_values` struct. Each CMMA value takes up one byte.

```
struct kvm_s390_cmma_log {
    __u64 start_gfn;
    __u32 count;
    __u32 flags;
    union {
        __u64 remaining;
        __u64 mask;
    };
    __u64 values;
};
```

`start_gfn` indicates the starting guest frame number,

`count` indicates how many values are to be considered in the buffer,

`flags` is not used and must be 0.

`mask` indicates which PGSTE bits are to be considered.

`remaining` is not used.

`values` points to the buffer in userspace where to store the values.

This ioctl can fail with `-ENOMEM` if not enough memory can be allocated to complete the task, with `-ENXIO` if CMMA is not enabled, with `-EINVAL` if the count field is too large (e.g. more than `KVM_S390_CMMA_SIZE_MAX`) or if the flags field was not 0, with `-EFAULT` if the userspace address is invalid, if invalid pages are written to (e.g. after the end of memory) or if no page table is present for the addresses (e.g. when using hugepages).

## 4.109 KVM\_PPC\_GET\_CPU\_CHAR

### Capability

KVM\_CAP\_PPC\_GET\_CPU\_CHAR

### Architectures

powerpc

### Type

vm ioctl

### Parameters

struct kvm\_ppc\_cpu\_char (out)

### Returns

0 on successful completion, -EFAULT if struct kvm\_ppc\_cpu\_char cannot be written

This ioctl gives userspace information about certain characteristics of the CPU relating to speculative execution of instructions and possible information leakage resulting from speculative execution (see CVE-2017-5715, CVE-2017-5753 and CVE-2017-5754). The information is returned in struct kvm\_ppc\_cpu\_char, which looks like this:

```
struct kvm_ppc_cpu_char {
    __u64    character;           /* characteristics of the CPU
    ↪ */
    __u64    behaviour;          /* recommended software
    ↪ behaviour */
    __u64    character_mask;      /* valid bits in character */
    __u64    behaviour_mask;      /* valid bits in behaviour */
};
```

For extensibility, the character\_mask and behaviour\_mask fields indicate which bits of character and behaviour have been filled in by the kernel. If the set of defined bits is extended in future then userspace will be able to tell whether it is running on a kernel that knows about the new bits.

The character field describes attributes of the CPU which can help with preventing inadvertent information disclosure - specifically, whether there is an instruction to flash-invalidate the L1 data cache (ori 30,30,0 or mtspr SPRN\_TRIG2,rN), whether the L1 data cache is set to a mode where entries can only be used by the thread that created them, whether the bcctr[l] instruction prevents speculation, and whether a speculation barrier instruction (ori 31,31,0) is provided.

The behaviour field describes actions that software should take to prevent inadvertent information disclosure, and thus describes which vulnerabilities the hardware is subject to; specifically whether the L1 data cache should be flushed when returning to user mode from the kernel, and whether a speculation barrier should be placed between an array bounds check and the array access.

These fields use the same bit definitions as the new H\_GET\_CPU\_CHARACTERISTICS hypercall.

#### **4.110 KVM\_MEMORY\_ENCRYPT\_OP**

**Capability**

basic

**Architectures**

x86

**Type**

vm

**Parameters**

an opaque platform specific structure (in/out)

**Returns**

0 on success; -1 on error

If the platform supports creating encrypted VMs then this ioctl can be used for issuing platform-specific memory encryption commands to manage those encrypted VMs.

Currently, this ioctl is used for issuing Secure Encrypted Virtualization (SEV) commands on AMD Processors. The SEV commands are defined in *Secure Encrypted Virtualization (SEV)*.

#### **4.111 KVM\_MEMORY\_ENCRYPT\_REG\_REGION**

**Capability**

basic

**Architectures**

x86

**Type**

system

**Parameters**

struct kvm\_enc\_region (in)

**Returns**

0 on success; -1 on error

This ioctl can be used to register a guest memory region which may contain encrypted data (e.g. guest RAM, SMRAM etc).

It is used in the SEV-enabled guest. When encryption is enabled, a guest memory region may contain encrypted data. The SEV memory encryption engine uses a tweak such that two identical plaintext pages, each at different locations will have differing ciphertexts. So swapping or moving ciphertext of those pages will not result in plaintext being swapped. So relocating (or migrating) physical backing pages for the SEV guest will require some additional steps.

Note: The current SEV key management spec does not provide commands to swap or migrate (move) ciphertext pages. Hence, for now we pin the guest memory region registered with the ioctl.



### 4.112 KVM\_MEMORY\_ENCRYPT\_UNREG\_REGION

**Capability**

basic

**Architectures**

x86

**Type**

system

**Parameters**

struct kvm\_enc\_region (in)

**Returns**

0 on success; -1 on error

This ioctl can be used to unregister the guest memory region registered with KVM\_MEMORY\_ENCRYPT\_REG\_REGION ioctl above.

### 4.113 KVM\_HYPERV\_EVENTFD

**Capability**

KVM\_CAP\_HYPERV\_EVENTFD

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct kvm\_hyperv\_eventfd (in)

This ioctl (un)registers an eventfd to receive notifications from the guest on the specified Hyper-V connection id through the SIGNAL\_EVENT hypercall, without causing a user exit. SIGNAL\_EVENT hypercall with non-zero event flag number (bits 24-31) still triggers a KVM\_EXIT\_HYPERV\_HCALL user exit.

```
struct kvm_hyperv_eventfd {
    __u32 conn_id;
    __s32 fd;
    __u32 flags;
    __u32 padding[3];
};
```

The conn\_id field should fit within 24 bits:

```
#define KVM_HYPERV_CONN_ID_MASK          0x00ffffff
```

The acceptable values for the flags field are:

```
#define KVM_HYPERV_EVENTFD_DEASSIGN      (1 << 0)
```

**Returns**

0 on success, -EINVAL if conn\_id or flags is outside the allowed

range, -ENOENT on deassign if the conn\_id isn't registered, -EEXIST on assign if the conn\_id is already registered

#### 4.114 KVM\_GET\_NESTED\_STATE

**Capability**

KVM\_CAP\_NESTED\_STATE

**Architectures**

x86

**Type**

vcpu ioctl

**Parameters**

struct kvm\_nested\_state (in/out)

**Returns**

0 on success, -1 on error

Errors:

E2BI the total state size exceeds the value of 'size' specified by the user; the size required will be written into size.

```
struct kvm_nested_state {
    __u16 flags;
    __u16 format;
    __u32 size;

    union {
        struct kvm_vmx_nested_state_hdr vmx;
        struct kvm_svm_nested_state_hdr svm;

        /* Pad the header to 128 bytes. */
        __u8 pad[120];
    } hdr;

    union {
        struct kvm_vmx_nested_state_data vmx[0];
        struct kvm_svm_nested_state_data svm[0];
    } data;
};

#define KVM_STATE_NESTED_GUEST_MODE          0x00000001
#define KVM_STATE_NESTED_RUN_PENDING        0x00000002
#define KVM_STATE_NESTED_EVMCS              0x00000004

#define KVM_STATE_NESTED_FORMAT_VMX         0
#define KVM_STATE_NESTED_FORMAT_SVM         1
```

(continues on next page)

(continued from previous page)

```

#define KVM_STATE_NESTED_VMX_VMCS_SIZE      0x1000

#define KVM_STATE_NESTED_VMX_SMM_GUEST_MODE  0x00000001
#define KVM_STATE_NESTED_VMX_SMM_VMXON      0x00000002

#define KVM_STATE_VMX_PREEMPTION_TIMER_DEADLINE 0x00000001

struct kvm_vmx_nested_state_hdr {
    __u64 vmxon_pa;
    __u64 vmcs12_pa;

    struct {
        __u16 flags;
    } smm;

    __u32 flags;
    __u64 preemption_timer_deadline;
};

struct kvm_vmx_nested_state_data {
    __u8 vmcs12[KVM_STATE_NESTED_VMX_VMCS_SIZE];
    __u8 shadow_vmcs12[KVM_STATE_NESTED_VMX_VMCS_SIZE];
};

```

This ioctl copies the vcpu's nested virtualization state from the kernel to userspace. The maximum size of the state can be retrieved by passing KVM\_CAP\_NESTED\_STATE to the KVM\_CHECK\_EXTENSION ioctl().

#### 4.115 KVM\_SET\_NESTED\_STATE

##### Capability

KVM\_CAP\_NESTED\_STATE

##### Architectures

x86

##### Type

vcpu ioctl

##### Parameters

struct kvm\_nested\_state (in)

##### Returns

0 on success, -1 on error

This copies the vcpu's kvm\_nested\_state struct from userspace to the kernel. For the definition of struct kvm\_nested\_state, see KVM\_GET\_NESTED\_STATE.

#### 4.116 KVM\_(UN)REGISTER\_COALESCED\_MMIO

**Capability**

KVM\_CAP\_COALESCED\_MMIO (for coalesced mmio)  
KVM\_CAP\_COALESCED\_PIO (for coalesced pio)

**Architectures**

all

**Type**

vm ioctl

**Parameters**

struct kvm\_coalesced\_mmio\_zone

**Returns**

0 on success, < 0 on error

Coalesced I/O is a performance optimization that defers hardware register write emulation so that userspace exits are avoided. It is typically used to reduce the overhead of emulating frequently accessed hardware registers.

When a hardware register is configured for coalesced I/O, write accesses do not exit to userspace and their value is recorded in a ring buffer that is shared between kernel and userspace.

Coalesced I/O is used if one or more write accesses to a hardware register can be deferred until a read or a write to another hardware register on the same device. This last access will cause a vmexit and userspace will process accesses from the ring buffer before emulating it. That will avoid exiting to userspace on repeated writes.

Coalesced pio is based on coalesced mmio. There is little difference between coalesced mmio and pio except that coalesced pio records accesses to I/O ports.

#### 4.117 KVM\_CLEAR\_DIRTY\_LOG (vm ioctl)

**Capability**

KVM\_CAP\_MANUAL\_DIRTY\_LOG\_PROTECT2

**Architectures**

x86, arm, arm64, mips

**Type**

vm ioctl

**Parameters**

struct kvm\_dirty\_log (in)

**Returns**

0 on success, -1 on error

```
/* for KVM_CLEAR_DIRTY_LOG */
struct kvm_clear_dirty_log {
    __u32 slot;
    __u32 num_pages;
```

(continues on next page)

(continued from previous page)

```

__u64 first_page;
union {
    void __user *dirty_bitmap; /* one bit per page */
    __u64 padding;
};
};

```

The `ioctl` clears the dirty status of pages in a memory slot, according to the bitmap that is passed in `struct kvm_clear_dirty_log`'s `dirty_bitmap` field. Bit 0 of the bitmap corresponds to page “`first_page`” in the memory slot, and `num_pages` is the size in bits of the input bitmap. `first_page` must be a multiple of 64; `num_pages` must also be a multiple of 64 unless `first_page + num_pages` is the size of the memory slot. For each bit that is set in the input bitmap, the corresponding page is marked “clean” in KVM's dirty bitmap, and dirty tracking is re-enabled for that page (for example via write-protection, or by clearing the dirty bit in a page table entry).

If `KVM_CAP_MULTI_ADDRESS_SPACE` is available, bits 16-31 specifies the address space for which you want to return the dirty bitmap. They must be less than the value that `KVM_CHECK_EXTENSION` returns for the `KVM_CAP_MULTI_ADDRESS_SPACE` capability.

This `ioctl` is mostly useful when `KVM_CAP_MANUAL_DIRTY_LOG_PROTECT2` is enabled; for more information, see the description of the capability. However, it can always be used as long as `KVM_CHECK_EXTENSION` confirms that `KVM_CAP_MANUAL_DIRTY_LOG_PROTECT2` is present.

#### 4.118 KVM\_GET\_SUPPORTED\_HV\_CPUID

##### Capability

`KVM_CAP_HYPERV_CPUID`

##### Architectures

x86

##### Type

`vcpu ioctl`

##### Parameters

`struct kvm_cpuid2` (in/out)

##### Returns

0 on success, -1 on error

```

struct kvm_cpuid2 {
    __u32 nent;
    __u32 padding;
    struct kvm_cpuid_entry2 entries[0];
};

struct kvm_cpuid_entry2 {
    __u32 function;
    __u32 index;
};

```

(continues on next page)

(continued from previous page)

```
__u32 flags;
__u32 eax;
__u32 ebx;
__u32 ecx;
__u32 edx;
__u32 padding[3];
};
```

This ioctl returns x86 cpuid features leaves related to Hyper-V emulation in KVM. Userspace can use the information returned by this ioctl to construct cpuid information presented to guests consuming Hyper-V enlightenments (e.g. Windows or Hyper-V guests).

CPUID feature leaves returned by this ioctl are defined by Hyper-V Top Level Functional Specification (TLFS). These leaves can't be obtained with KVM\_GET\_SUPPORTED\_CPUID ioctl because some of them intersect with KVM feature leaves (0x40000000, 0x40000001).

**Currently, the following list of CPUID leaves are returned:**

- HYPERV\_CPUID\_VENDOR\_AND\_MAX\_FUNCTIONS
- HYPERV\_CPUID\_INTERFACE
- HYPERV\_CPUID\_VERSION
- HYPERV\_CPUID\_FEATURES
- HYPERV\_CPUID\_ENLIGHTMENT\_INFO
- HYPERV\_CPUID\_IMPLEMENT\_LIMITS
- HYPERV\_CPUID\_NESTED\_FEATURES
- HYPERV\_CPUID\_SYNDBG\_VENDOR\_AND\_MAX\_FUNCTIONS
- HYPERV\_CPUID\_SYNDBG\_INTERFACE
- HYPERV\_CPUID\_SYNDBG\_PLATFORM\_CAPABILITIES

HYPERV\_CPUID\_NESTED\_FEATURES leaf is only exposed when Enlightened VMCS was enabled on the corresponding vCPU (KVM\_CAP\_HYPERV\_ENLIGHTENED\_VMCS).

Userspace invokes KVM\_GET\_SUPPORTED\_HV\_CPUID by passing a `kvm_cpuid2` structure with the 'nent' field indicating the number of entries in the variable-size array 'entries'. If the number of entries is too low to describe all Hyper-V feature leaves, an error (E2BIG) is returned. If the number is more or equal to the number of Hyper-V feature leaves, the 'nent' field is adjusted to the number of valid entries in the 'entries' array, which is then filled.

'index' and 'flags' fields in 'struct kvm\_cpuid\_entry2' are currently reserved, userspace should not expect to get any particular value there.

## 4.119 KVM\_ARM\_VCPU\_FINALIZE

### Architectures

arm, arm64

### Type

vcpu ioctl

### Parameters

int feature (in)

### Returns

0 on success, -1 on error

Errors:

EPERM	feature not enabled, needs configuration, or already finalized
EIN- VAL	feature unknown or not present

Recognised values for feature:

arm64	KVM_ARM_VCPU_SVE (requires KVM_CAP_ARM_SVE)
-------	---

Finalizes the configuration of the specified vcpu feature.

The vcpu must already have been initialised, enabling the affected feature, by means of a successful KVM\_ARM\_VCPU\_INIT call with the appropriate flag set in features[].

For affected vcpu features, this is a mandatory step that must be performed before the vcpu is fully usable.

Between KVM\_ARM\_VCPU\_INIT and KVM\_ARM\_VCPU\_FINALIZE, the feature may be configured by use of ioctls such as KVM\_SET\_ONE\_REG. The exact configuration that should be performed and how to do it are feature-dependent.

Other calls that depend on a particular feature being finalized, such as KVM\_RUN, KVM\_GET\_REG\_LIST, KVM\_GET\_ONE\_REG and KVM\_SET\_ONE\_REG, will fail with -EPERM unless the feature has already been finalized by means of a KVM\_ARM\_VCPU\_FINALIZE call.

See KVM\_ARM\_VCPU\_INIT for details of vcpu features that require finalization using this ioctl.

## 4.120 KVM\_SET\_PMU\_EVENT\_FILTER

**Capability**

KVM\_CAP\_PMU\_EVENT\_FILTER

**Architectures**

x86

**Type**

vm ioctl

**Parameters**

struct kvm\_pmu\_event\_filter (in)

**Returns**

0 on success, -1 on error

```
struct kvm_pmu_event_filter {
    __u32 action;
    __u32 nevents;
    __u32 fixed_counter_bitmap;
    __u32 flags;
    __u32 pad[4];
    __u64 events[0];
};
```

This ioctl restricts the set of PMU events that the guest can program. The argument holds a list of events which will be allowed or denied. The eventsel+umask of each event the guest attempts to program is compared against the events field to determine whether the guest should have access. The events field only controls general purpose counters; fixed purpose counters are controlled by the fixed\_counter\_bitmap.

No flags are defined yet, the field must be zero.

Valid values for 'action' :

```
#define KVM_PMU_EVENT_ALLOW 0
#define KVM_PMU_EVENT_DENY 1
```

## 4.121 KVM\_PPC\_SVM\_OFF

**Capability**

basic

**Architectures**

powerpc

**Type**

vm ioctl

**Parameters**

none

**Returns**

0 on successful completion,



Errors:

EINVAL	if ultravisor failed to terminate the secure guest
ENOMEM	if hypervisor failed to allocate new radix page tables for guest

This ioctl is used to turn off the secure mode of the guest or transition the guest from secure mode to normal mode. This is invoked when the guest is reset. This has no effect if called for a normal guest.

This ioctl issues an ultravisor call to terminate the secure guest, unpins the VPA pages and releases all the device pages that are used to track the secure pages by hypervisor.

#### 4.122 KVM\_S390\_NORMAL\_RESET

**Capability**

KVM\_CAP\_S390\_VCPU\_RESETS

**Architectures**

s390

**Type**

vcpu ioctl

**Parameters**

none

**Returns**

0

This ioctl resets VCPU registers and control structures according to the cpu reset definition in the POP (Principles Of Operation).

#### 4.123 KVM\_S390\_INITIAL\_RESET

**Capability**

none

**Architectures**

s390

**Type**

vcpu ioctl

**Parameters**

none

**Returns**

0

This ioctl resets VCPU registers and control structures according to the initial cpu reset definition in the POP. However, the cpu is not put into ESA mode. This reset is a superset of the normal reset.

#### 4.124 KVM\_S390\_CLEAR\_RESET

**Capability**

KVM\_CAP\_S390\_VCPU\_RESETS

**Architectures**

s390

**Type**

vcpu ioctl

**Parameters**

none

**Returns**

0

This ioctl resets VCPU registers and control structures according to the clear cpu reset definition in the POP. However, the cpu is not put into ESA mode. This reset is a superset of the initial reset.

#### 4.125 KVM\_S390\_PV\_COMMAND

**Capability**

KVM\_CAP\_S390\_PROTECTED

**Architectures**

s390

**Type**

vm ioctl

**Parameters**

struct kvm\_pv\_cmd

**Returns**

0 on success, &lt; 0 on error

```
struct kvm_pv_cmd {
    __u32 cmd;           /* Command to be executed */
    __u16 rc;            /* Ultravisor return code */
    __u16 rrc;           /* Ultravisor return reason code */
    __u64 data;          /* Data or address */
    __u32 flags;         /* flags for future extensions. Must be 0 for
↳ now */
    __u32 reserved[3];
};
```

cmd values:

**KVM\_PV\_ENABLE**

Allocate memory and register the VM with the Ultravisor, thereby donating memory to the Ultravisor that will become inaccessible to KVM. All existing CPUs are converted to protected ones. After this command has succeeded, any CPU added via hotplug will become protected during its creation as well.

Errors:

EINTR an unmasked signal is pending
-------------------------------------

## KVM\_PV\_DISABLE

Deregister the VM from the Ultravisor and reclaim the memory that had been donated to the Ultravisor, making it usable by the kernel again. All registered VCPUs are converted back to non-protected ones.

## KVM\_PV\_VM\_SET\_SEC\_PARMs

Pass the image header from VM memory to the Ultravisor in preparation of image unpacking and verification.

## KVM\_PV\_VM\_UNPACK

Unpack (protect and decrypt) a page of the encrypted boot image.

## KVM\_PV\_VM\_VERIFY

Verify the integrity of the unpacked image. Only if this succeeds, KVM is allowed to start protected VCPUs.

## 4.126 KVM\_X86\_SET\_MSR\_FILTER

### Capability

KVM\_X86\_SET\_MSR\_FILTER

### Architectures

x86

### Type

vm ioctl

### Parameters

struct kvm\_msr\_filter

### Returns

0 on success, < 0 on error

```
struct kvm_msr_filter_range {
#define KVM_MSR_FILTER_READ  (1 << 0)
#define KVM_MSR_FILTER_WRITE (1 << 1)
    __u32 flags;
    __u32 nmsrs; /* number of msrs in bitmap */
    __u32 base; /* MSR index the bitmap starts at */
    __u8 *bitmap; /* a 1 bit allows the operations in flags, 0
↪denies */
};

#define KVM_MSR_FILTER_MAX_RANGES 16
struct kvm_msr_filter {
#define KVM_MSR_FILTER_DEFAULT_ALLOW (0 << 0)
#define KVM_MSR_FILTER_DEFAULT_DENY (1 << 0)
    __u32 flags;
```

(continues on next page)

(continued from previous page)

```
struct kvm_msr_filter_range ranges[KVM_MSR_FILTER_MAX_RANGES];  
};
```

flags values for struct `kvm_msr_filter_range`:

#### `KVM_MSR_FILTER_READ`

Filter read accesses to MSRs using the given bitmap. A 0 in the bitmap indicates that a read should immediately fail, while a 1 indicates that a read for a particular MSR should be handled regardless of the default filter action.

#### `KVM_MSR_FILTER_WRITE`

Filter write accesses to MSRs using the given bitmap. A 0 in the bitmap indicates that a write should immediately fail, while a 1 indicates that a write for a particular MSR should be handled regardless of the default filter action.

#### `KVM_MSR_FILTER_READ | KVM_MSR_FILTER_WRITE`

Filter both read and write accesses to MSRs using the given bitmap. A 0 in the bitmap indicates that both reads and writes should immediately fail, while a 1 indicates that reads and writes for a particular MSR are not filtered by this range.

flags values for struct `kvm_msr_filter`:

#### `KVM_MSR_FILTER_DEFAULT_ALLOW`

If no filter range matches an MSR index that is getting accessed, KVM will fall back to allowing access to the MSR.

#### `KVM_MSR_FILTER_DEFAULT_DENY`

If no filter range matches an MSR index that is getting accessed, KVM will fall back to rejecting access to the MSR. In this mode, all MSRs that should be processed by KVM need to explicitly be marked as allowed in the bitmaps.

This ioctl allows user space to define up to 16 bitmaps of MSR ranges to specify whether a certain MSR access should be explicitly filtered for or not.

If this ioctl has never been invoked, MSR accesses are not guarded and the default KVM in-kernel emulation behavior is fully preserved.

Calling this ioctl with an empty set of ranges (all `nmsrs == 0`) disables MSR filtering. In that mode, `KVM_MSR_FILTER_DEFAULT_DENY` is invalid and causes an error.

As soon as the filtering is in place, every MSR access is processed through the filtering except for accesses to the x2APIC MSRs (from 0x800 to 0x8ff); x2APIC MSRs are always allowed, independent of the `default_allow` setting, and their behavior depends on the `X2APIC_ENABLE` bit of the APIC base register.

If a bit is within one of the defined ranges, read and write accesses are guarded by the bitmap's value for the MSR index if the kind of access is included in

the struct `kvm_msr_filter_range` flags. If no range cover this particular access, the behavior is determined by the flags field in the `kvm_msr_filter` struct: `KVM_MSR_FILTER_DEFAULT_ALLOW` and `KVM_MSR_FILTER_DEFAULT_DENY`.

Each bitmap range specifies a range of MSRs to potentially allow access on. The range goes from MSR index `[base .. base+nmsrs]`. The flags field indicates whether reads, writes or both reads and writes are filtered by setting a 1 bit in the bitmap for the corresponding MSR index.

If an MSR access is not permitted through the filtering, it generates a #GP inside the guest. When combined with `KVM_CAP_X86_USER_SPACE_MSR`, that allows user space to deflect and potentially handle various MSR accesses into user space.

Note, invoking this ioctl with a vCPU is running is inherently racy. However, KVM does guarantee that vCPUs will see either the previous filter or the new filter, e.g. MSRs with identical settings in both the old and new filter will have deterministic behavior.

### 1.1.5 5. The `kvm_run` structure

Application code obtains a pointer to the `kvm_run` structure by `mmap()`ing a `vcpu` fd. From that point, application code can control execution by changing fields in `kvm_run` prior to calling the `KVM_RUN` ioctl, and obtain information about the reason `KVM_RUN` returned by looking up structure members.

```
struct kvm_run {
    /* in */
    __u8 request_interrupt_window;
```

Request that `KVM_RUN` return when it becomes possible to inject external interrupts into the guest. Useful in conjunction with `KVM_INTERRUPT`.

```
__u8 immediate_exit;
```

This field is polled once when `KVM_RUN` starts; if non-zero, `KVM_RUN` exits immediately, returning `-EINTR`. In the common scenario where a signal is used to “kick” a VCPU out of `KVM_RUN`, this field can be used to avoid usage of `KVM_SET_SIGNAL_MASK`, which has worse scalability. Rather than blocking the signal outside `KVM_RUN`, userspace can set up a signal handler that sets `run->immediate_exit` to a non-zero value.

This field is ignored if `KVM_CAP_IMMEDIATE_EXIT` is not available.

```
__u8 padding1[6];

/* out */
__u32 exit_reason;
```

When `KVM_RUN` has returned successfully (return value 0), this informs application code why `KVM_RUN` has returned. Allowable values for this field are detailed below.

```
__u8 ready_for_interrupt_injection;
```

If request\_interrupt\_window has been specified, this field indicates an interrupt can be injected now with KVM\_INTERRUPT.

```
__u8 if_flag;
```

The value of the current interrupt flag. Only valid if in-kernel local APIC is not used.

```
__u16 flags;
```

More architecture-specific flags detailing state of the VCPU that may affect the device's behavior. The only currently defined flag is KVM\_RUN\_X86\_SMM, which is valid on x86 machines and is set if the VCPU is in system management mode.

```
/* in (pre_kvm_run), out (post_kvm_run) */
__u64 cr8;
```

The value of the cr8 register. Only valid if in-kernel local APIC is not used. Both input and output.

```
__u64 apic_base;
```

The value of the APIC BASE msr. Only valid if in-kernel local APIC is not used. Both input and output.

```
union {
    /* KVM_EXIT_UNKNOWN */
    struct {
        __u64 hardware_exit_reason;
    } hw;
};
```

If exit\_reason is KVM\_EXIT\_UNKNOWN, the vcpu has exited due to unknown reasons. Further architecture-specific information is available in hardware\_exit\_reason.

```
/* KVM_EXIT_FAIL_ENTRY */
struct {
    __u64 hardware_entry_failure_reason;
    __u32 cpu; /* if KVM_LAST_CPU */
} fail_entry;
```

If exit\_reason is KVM\_EXIT\_FAIL\_ENTRY, the vcpu could not be run due to unknown reasons. Further architecture-specific information is available in hardware\_entry\_failure\_reason.

```
/* KVM_EXIT_EXCEPTION */
struct {
    __u32 exception;
    __u32 error_code;
} ex;
```

Unused.

```

        /* KVM_EXIT_IO */
        struct {
#define KVM_EXIT_IO_IN 0
#define KVM_EXIT_IO_OUT 1
                __u8 direction;
                __u8 size; /* bytes */
                __u16 port;
                __u32 count;
                __u64 data_offset; /* relative to kvm_run_
↪start */
        } io;

```

If `exit_reason` is `KVM_EXIT_IO`, then the `vcpu` has executed a port I/O instruction which could not be satisfied by `kvm`. `data_offset` describes where the data is located (`KVM_EXIT_IO_OUT`) or where `kvm` expects application code to place the data for the next `KVM_RUN` invocation (`KVM_EXIT_IO_IN`). Data format is a packed array.

```

/* KVM_EXIT_DEBUG */
struct {
    struct kvm_debug_exit_arch arch;
} debug;

```

If the `exit_reason` is `KVM_EXIT_DEBUG`, then a `vcpu` is processing a debug event for which architecture specific information is returned.

```

/* KVM_EXIT_MMIO */
struct {
    __u64 phys_addr;
    __u8 data[8];
    __u32 len;
    __u8 is_write;
} mmio;

```

If `exit_reason` is `KVM_EXIT_MMIO`, then the `vcpu` has executed a memory-mapped I/O instruction which could not be satisfied by `kvm`. The ‘data’ member contains the written data if ‘is\_write’ is true, and should be filled by application code otherwise.

The ‘data’ member contains, in its first ‘len’ bytes, the value as it would appear if the VCPU performed a load or store of the appropriate width directly to the byte array.

**Note:** For `KVM_EXIT_IO`, `KVM_EXIT_MMIO`, `KVM_EXIT_OSI`, `KVM_EXIT_PAPR`, `KVM_EXIT_EPR`, `KVM_EXIT_X86_RDMSR` and `KVM_EXIT_X86_WRMSR` the corresponding operations are complete (and guest state is consistent) only after userspace has re-entered the kernel with `KVM_RUN`. The kernel side will first finish incomplete operations and then check for pending signals. Userspace can re-enter the guest with an unmasked signal pending to complete pending opera-

tions.

---

```
/* KVM_EXIT_HYPERCALL */
struct {
    __u64 nr;
    __u64 args[6];
    __u64 ret;
    __u32 longmode;
    __u32 pad;
} hypercall;
```

Unused. This was once used for ‘hypercall to userspace’ . To implement such functionality, use KVM\_EXIT\_IO (x86) or KVM\_EXIT\_MMIO (all except s390).

---

**Note:** KVM\_EXIT\_IO is significantly faster than KVM\_EXIT\_MMIO.

---

```
/* KVM_EXIT_TPR_ACCESS */
struct {
    __u64 rip;
    __u32 is_write;
    __u32 pad;
} tpr_access;
```

To be documented (KVM\_TPR\_ACCESS\_REPORTING).

```
/* KVM_EXIT_S390_SIEIC */
struct {
    __u8 icptcode;
    __u64 mask; /* psw upper half */
    __u64 addr; /* psw lower half */
    __u16 ipa;
    __u32 ipb;
} s390_sieic;
```

s390 specific.

```
/* KVM_EXIT_S390_RESET */
#define KVM_S390_RESET_POR      1
#define KVM_S390_RESET_CLEAR    2
#define KVM_S390_RESET_SUBSYSTEM 4
#define KVM_S390_RESET_CPU_INIT 8
#define KVM_S390_RESET_IPL      16
    __u64 s390_reset_flags;
```

s390 specific.

```
/* KVM_EXIT_S390_UCONTROL */
struct {
    __u64 trans_exc_code;
```

(continues on next page)



(continued from previous page)

```

    __u32 pgm_code;
} s390_ucontrol;

```

s390 specific. A page fault has occurred for a user controlled virtual machine (KVM\_VM\_S390\_UNCONTROL) on it's host page table that cannot be resolved by the kernel. The program code and the translation exception code that were placed in the cpu's lowcore are presented here as defined by the z Architecture Principles of Operation Book in the Chapter for Dynamic Address Translation (DAT)

```

/* KVM_EXIT_DCR */
struct {
    __u32 dcrn;
    __u32 data;
    __u8  is_write;
} dcr;

```

Deprecated - was used for 440 KVM.

```

/* KVM_EXIT_OSI */
struct {
    __u64 gprs[32];
} osi;

```

MOL uses a special hypercall interface it calls 'OSI' . To enable it, we catch hypercalls and exit with this exit struct that contains all the guest gprs.

If exit\_reason is KVM\_EXIT\_OSI, then the vcpu has triggered such a hypercall. Userspace can now handle the hypercall and when it's done modify the gprs as necessary. Upon guest entry all guest GPRs will then be replaced by the values in this struct.

```

/* KVM_EXIT_PAPR_HCALL */
struct {
    __u64 nr;
    __u64 ret;
    __u64 args[9];
} papr_hcall;

```

This is used on 64-bit PowerPC when emulating a pSeries partition, e.g. with the 'pseries' machine type in qemu. It occurs when the guest does a hypercall using the 'sc 1' instruction. The 'nr' field contains the hypercall number (from the guest R3), and 'args' contains the arguments (from the guest R4 - R12). Userspace should put the return code in 'ret' and any extra returned values in args[]. The possible hypercalls are defined in the Power Architecture Platform Requirements (PAPR) document available from [www.power.org](http://www.power.org) (free developer registration required to access it).

```

/* KVM_EXIT_S390_TSCH */
struct {
    __u16 subchannel_id;
    __u16 subchannel_nr;
}

```

(continues on next page)

(continued from previous page)

```
    __u32 io_int_parm;
    __u32 io_int_word;
    __u32 ipb;
    __u8 dequeued;
} s390_tsch;
```

s390 specific. This exit occurs when KVM\_CAP\_S390\_CSS\_SUPPORT has been enabled and TEST SUBCHANNEL was intercepted. If dequeued is set, a pending I/O interrupt for the target subchannel has been dequeued and subchannel\_id, subchannel\_nr, io\_int\_parm and io\_int\_word contain the parameters for that interrupt. ipb is needed for instruction parameter decoding.

```
/* KVM_EXIT_EPR */
struct {
    __u32 epr;
} epr;
```

On FSL BookE PowerPC chips, the interrupt controller has a fast patch interrupt acknowledge path to the core. When the core successfully delivers an interrupt, it automatically populates the EPR register with the interrupt vector number and acknowledges the interrupt inside the interrupt controller.

In case the interrupt controller lives in user space, we need to do the interrupt acknowledge cycle through it to fetch the next to be delivered interrupt vector using this exit.

It gets triggered whenever both KVM\_CAP\_PPC\_EPR are enabled and an external interrupt has just been delivered into the guest. User space should put the acknowledged interrupt vector into the 'epr' field.

```
/* KVM_EXIT_SYSTEM_EVENT */
struct {
#define KVM_SYSTEM_EVENT_SHUTDOWN      1
#define KVM_SYSTEM_EVENT_RESET        2
#define KVM_SYSTEM_EVENT_CRASH        3
    __u32 type;
    __u64 flags;
} system_event;
```

If exit\_reason is KVM\_EXIT\_SYSTEM\_EVENT then the vcpu has triggered a system-level event using some architecture specific mechanism (hypercall or some special instruction). In case of ARM/ARM64, this is triggered using HVC instruction based PSCI call from the vcpu. The 'type' field describes the system-level event type. The 'flags' field describes architecture specific flags for the system-level event.

Valid values for 'type' are:

- KVM\_SYSTEM\_EVENT\_SHUTDOWN - the guest has requested a shutdown of the VM. Userspace is not obliged to honour this, and if it does honour this does not need to destroy the VM synchronously (ie it may call KVM\_RUN again before shutdown finally occurs).

- `KVM_SYSTEM_EVENT_RESET` - the guest has requested a reset of the VM. As with `SHUTDOWN`, userspace can choose to ignore the request, or to schedule the reset to occur in the future and may call `KVM_RUN` again.
- `KVM_SYSTEM_EVENT_CRASH` - the guest crash occurred and the guest has requested a crash condition maintenance. Userspace can choose to ignore the request, or to gather VM memory core dump and/or reset/shutdown of the VM.

```
/* KVM_EXIT_IOAPIC_EOI */
struct {
    __u8 vector;
} eoi;
```

Indicates that the VCPU's in-kernel local APIC received an EOI for a level-triggered IOAPIC interrupt. This exit only triggers when the IOAPIC is implemented in userspace (i.e. `KVM_CAP_SPLIT_IRQCHIP` is enabled); the userspace IOAPIC should process the EOI and retrigger the interrupt if it is still asserted. Vector is the LAPIC interrupt vector for which the EOI was received.

```
    struct kvm_hyperv_exit {
#define KVM_EXIT_HYPERV_SYNIC          1
#define KVM_EXIT_HYPERV_HCALL          2
#define KVM_EXIT_HYPERV_SYNDBG          3
        __u32 type;
        __u32 pad1;
        union {
            struct {
                __u32 msr;
                __u32 pad2;
                __u64 control;
                __u64 evt_page;
                __u64 msg_page;
            } synic;
            struct {
                __u64 input;
                __u64 result;
                __u64 params[2];
            } hcall;
            struct {
                __u32 msr;
                __u32 pad2;
                __u64 control;
                __u64 status;
                __u64 send_page;
                __u64 recv_page;
                __u64 pending_page;
            } syndbg;
        } u;
    };
/* KVM_EXIT_HYPERV */
struct kvm_hyperv_exit hyperv;
```

Indicates that the VCPU exits into userspace to process some tasks related to Hyper-V emulation.

Valid values for ‘type’ are:

- KVM\_EXIT\_HYPERV\_SYNIC - synchronously notify user-space about

Hyper-V SynIC state change. Notification is used to remap SynIC event/message pages and to enable/disable SynIC messages/events processing in userspace.

- KVM\_EXIT\_HYPERV\_SYNDBG - synchronously notify user-space about

Hyper-V Synthetic debugger state change. Notification is used to either update the pending\_page location or to send a control command (send the buffer located in send\_page or recv a buffer to recv\_page).

```
/* KVM_EXIT_ARM_NISV */
struct {
    __u64 esr_iss;
    __u64 fault_ipa;
} arm_nisv;
```

Used on arm and arm64 systems. If a guest accesses memory not in a memslot, KVM will typically return to userspace and ask it to do MMIO emulation on its behalf. However, for certain classes of instructions, no instruction decode (direction, length of memory access) is provided, and fetching and decoding the instruction from the VM is overly complicated to live in the kernel.

Historically, when this situation occurred, KVM would print a warning and kill the VM. KVM assumed that if the guest accessed non-memslot memory, it was trying to do I/O, which just couldn't be emulated, and the warning message was phrased accordingly. However, what happened more often was that a guest bug caused access outside the guest memory areas which should lead to a more meaningful warning message and an external abort in the guest, if the access did not fall within an I/O window.

Userspace implementations can query for KVM\_CAP\_ARM\_NISV\_TO\_USER, and enable this capability at VM creation. Once this is done, these types of errors will instead return to userspace with KVM\_EXIT\_ARM\_NISV, with the valid bits from the HSR (arm) and ESR\_EL2 (arm64) in the esr\_iss field, and the faulting IPA in the fault\_ipa field. Userspace can either fix up the access if it's actually an I/O access by decoding the instruction from guest memory (if it's very brave) and continue executing the guest, or it can decide to suspend, dump, or restart the guest.

Note that KVM does not skip the faulting instruction as it does for KVM\_EXIT\_MMIO, but userspace has to emulate any change to the processing state if it decides to decode and emulate the instruction.

```
/* KVM_EXIT_X86_RDMSR / KVM_EXIT_X86_WRMSR */
struct {
    __u8 error; /* user -> kernel */
    __u8 pad[7];
    __u32 reason; /* kernel -> user */
    __u32 index; /* kernel -> user */
}
```

(continues on next page)

(continued from previous page)

```

    __u64 data; /* kernel <-> user */
} msr;

```

Used on x86 systems. When the VM capability `KVM_CAP_X86_USER_SPACE_MSR` is enabled, MSR accesses to registers that would invoke a `#GP` by KVM kernel code will instead trigger a `KVM_EXIT_X86_RDMSR` exit for reads and `KVM_EXIT_X86_WRMSR` exit for writes.

The “reason” field specifies why the MSR trap occurred. User space will only receive MSR exit traps when a particular reason was requested during `ENABLE_CAP`. Currently valid exit reasons are:

`KVM_MSR_EXIT_REASON_UNKNOWN` - access to MSR that is unknown to KVM  
`KVM_MSR_EXIT_REASON_INVALID` - access to invalid MSRs or reserved bits  
`KVM_MSR_EXIT_REASON_FILTER` - access blocked by `KVM_X86_SET_MSR_FILTER`

For `KVM_EXIT_X86_RDMSR`, the “index” field tells user space which MSR the guest wants to read. To respond to this request with a successful read, user space writes the respective data into the “data” field and must continue guest execution to ensure the read data is transferred into guest register state.

If the RDMSR request was unsuccessful, user space indicates that with a “1” in the “error” field. This will inject a `#GP` into the guest when the VCPU is executed again.

For `KVM_EXIT_X86_WRMSR`, the “index” field tells user space which MSR the guest wants to write. Once finished processing the event, user space must continue vCPU execution. If the MSR write was unsuccessful, user space also sets the “error” field to “1”.

```

    /* Fix the size of the union. */
    char padding[256];
};

/*
 * shared registers between kvm and userspace.
 * kvm_valid_regs specifies the register classes set by the host
 * kvm_dirty_regs specified the register classes dirtied by
 * ↪userspace
 * struct kvm_sync_regs is architecture specific, as well as the
 * bits for kvm_valid_regs and kvm_dirty_regs
 */
__u64 kvm_valid_regs;
__u64 kvm_dirty_regs;
union {
    struct kvm_sync_regs regs;
    char padding[SYNC_REGS_SIZE_BYTES];
} s;

```

If `KVM_CAP_SYNC_REGS` is defined, these fields allow userspace to access certain guest registers without having to call `SET/GET_*REGS`. Thus we can avoid some system call overhead if userspace has to handle the exit. Userspace can query the

validity of the structure by checking `kvm_valid_regs` for specific bits. These bits are architecture specific and usually define the validity of a groups of registers. (e.g. one bit for general purpose registers)

Please note that the kernel is allowed to use the `kvm_run` structure as the primary storage for certain register types. Therefore, the kernel may use the values in `kvm_run` even if the corresponding bit in `kvm_dirty_regs` is not set.

```
};
```

### 1.1.6 6. Capabilities that can be enabled on vCPUs

There are certain capabilities that change the behavior of the virtual CPU or the virtual machine when enabled. To enable them, please see section 4.37. Below you can find a list of capabilities and what their effect on the vCPU or the virtual machine is when enabling them.

The following information is provided along with the description:

**Architectures:**

which instruction set architectures provide this ioctl. x86 includes both i386 and x86\_64.

**Target:**

whether this is a per-vcpu or per-vm capability.

**Parameters:**

what parameters are accepted by the capability.

**Returns:**

the return value. General error numbers (EBADF, ENOMEM, EINVAL) are not detailed, but errors with specific meanings are.

#### 6.1 KVM\_CAP\_PPC\_OSI

**Architectures**

ppc

**Target**

vcpu

**Parameters**

none

**Returns**

0 on success; -1 on error

This capability enables interception of OSI hypercalls that otherwise would be treated as normal system calls to be injected into the guest. OSI hypercalls were invented by Mac-on-Linux to have a standardized communication mechanism between the guest and the host.

When this capability is enabled, `KVM_EXIT_OSI` can occur.

## 6.2 KVM\_CAP\_PPC\_PAPR

### Architectures

ppc

### Target

vcpu

### Parameters

none

### Returns

0 on success; -1 on error

This capability enables interception of PAPR hypercalls. PAPR hypercalls are done using the hypercall instruction “sc 1” .

It also sets the guest privilege level to “supervisor” mode. Usually the guest runs in “hypervisor” privilege mode with a few missing features.

In addition to the above, it changes the semantics of SDR1. In this mode, the HTAB address part of SDR1 contains an HVA instead of a GPA, as PAPR keeps the HTAB invisible to the guest.

When this capability is enabled, KVM\_EXIT\_PAPR\_HCALL can occur.

## 6.3 KVM\_CAP\_SW\_TLB

### Architectures

ppc

### Target

vcpu

### Parameters

args[0] is the address of a struct `kvm_config_tlb`

### Returns

0 on success; -1 on error

```
struct kvm_config_tlb {
    __u64 params;
    __u64 array;
    __u32 mmu_type;
    __u32 array_len;
};
```

Configures the virtual CPU’ s TLB array, establishing a shared memory area between userspace and KVM. The “params” and “array” fields are userspace addresses of mmu-type-specific data structures. The “array\_len” field is an safety mechanism, and should be set to the size in bytes of the memory that userspace has reserved for the array. It must be at least the size dictated by “mmu\_type” and “params” .

While KVM\_RUN is active, the shared region is under control of KVM. Its contents are undefined, and any modification by userspace results in boundedly undefined

behavior.

On return from KVM\_RUN, the shared region will reflect the current state of the guest's TLB. If userspace makes any changes, it must call KVM\_DIRTY\_TLB to tell KVM which entries have been changed, prior to calling KVM\_RUN again on this vcpu.

For mmu types KVM\_MMU\_FSL\_BOOKE\_NOHV and KVM\_MMU\_FSL\_BOOKE\_HV:

- The “params” field is of type “struct kvm\_book3e\_206\_tlb\_params” .
- The “array” field points to an array of type “struct kvm\_book3e\_206\_tlb\_entry” .
- The array consists of all entries in the first TLB, followed by all entries in the second TLB.
- Within a TLB, entries are ordered first by increasing set number. Within a set, entries are ordered by way (increasing ESEL).
- The hash for determining set number in TLB0 is:  $(MAS2 \gg 12) \& (num\_sets - 1)$  where “num\_sets” is the tlb\_sizes[] value divided by the tlb\_ways[] value.
- The tsize field of mas1 shall be set to 4K on TLB0, even though the hardware ignores this value for TLB0.

## 6.4 KVM\_CAP\_S390\_CSS\_SUPPORT

### Architectures

s390

### Target

vcpu

### Parameters

none

### Returns

0 on success; -1 on error

This capability enables support for handling of channel I/O instructions.

TEST PENDING INTERRUPTION and the interrupt portion of TEST SUBCHANNEL are handled in-kernel, while the other I/O instructions are passed to userspace.

When this capability is enabled, KVM\_EXIT\_S390\_TSCH will occur on TEST SUBCHANNEL intercepts.

Note that even though this capability is enabled per-vcpu, the complete virtual machine is affected.



## 6.5 KVM\_CAP\_PPC\_EPR

### Architectures

ppc

### Target

vcpu

### Parameters

args[0] defines whether the proxy facility is active

### Returns

0 on success; -1 on error

This capability enables or disables the delivery of interrupts through the external proxy facility.

When enabled (args[0] != 0), every time the guest gets an external interrupt delivered, it automatically exits into user space with a KVM\_EXIT\_EPR exit to receive the topmost interrupt vector.

When disabled (args[0] == 0), behavior is as if this facility is unsupported.

When this capability is enabled, KVM\_EXIT\_EPR can occur.

## 6.6 KVM\_CAP\_IRQ\_MPIC

### Architectures

ppc

### Parameters

args[0] is the MPIC device fd; args[1] is the MPIC CPU number for this vcpu

This capability connects the vcpu to an in-kernel MPIC device.

## 6.7 KVM\_CAP\_IRQ\_XICS

### Architectures

ppc

### Target

vcpu

### Parameters

args[0] is the XICS device fd; args[1] is the XICS CPU number (server ID) for this vcpu

This capability connects the vcpu to an in-kernel XICS device.

## **6.8 KVM\_CAP\_S390\_IRQCHIP**

### **Architectures**

s390

### **Target**

vm

### **Parameters**

none

This capability enables the in-kernel irqchip for s390. Please refer to “4.24 KVM\_CREATE\_IRQCHIP” for details.

## **6.9 KVM\_CAP\_MIPS\_FPU**

### **Architectures**

mips

### **Target**

vcpu

### **Parameters**

args[0] is reserved for future use (should be 0).

This capability allows the use of the host Floating Point Unit by the guest. It allows the Config1.FP bit to be set to enable the FPU in the guest. Once this is done the KVM\_REG\_MIPS\_FPR\_\* and KVM\_REG\_MIPS\_FCR\_\* registers can be accessed (depending on the current guest FPU register mode), and the Status.FR, Config5.FRE bits are accessible via the KVM API and also from the guest, depending on them being supported by the FPU.

## **6.10 KVM\_CAP\_MIPS\_MSA**

### **Architectures**

mips

### **Target**

vcpu

### **Parameters**

args[0] is reserved for future use (should be 0).

This capability allows the use of the MIPS SIMD Architecture (MSA) by the guest. It allows the Config3.MSAP bit to be set to enable the use of MSA by the guest. Once this is done the KVM\_REG\_MIPS\_VEC\_\* and KVM\_REG\_MIPS\_MSA\_\* registers can be accessed, and the Config5.MSAEn bit is accessible via the KVM API and also from the guest.

## 6.74 KVM\_CAP\_SYNC\_REGS

### Architectures

s390, x86

### Target

s390: always enabled, x86: vcpu

### Parameters

none

### Returns

x86: KVM\_CHECK\_EXTENSION returns a bit-array indicating which register sets are supported (bitfields defined in arch/x86/include/uapi/asm/kvm.h).

As described above in the `kvm_sync_regs` struct info in section 5 (`kvm_run`): KVM\_CAP\_SYNC\_REGS “allow[s] userspace to access certain guest registers without having to call SET/GET\_\*REGS”. This reduces overhead by eliminating repeated `ioctl` calls for setting and/or getting register values. This is particularly important when userspace is making synchronous guest state modifications, e.g. when emulating and/or intercepting instructions in userspace.

For s390 specifics, please refer to the source code.

For x86:

- the register sets to be copied out to `kvm_run` are selectable by userspace (rather than all sets being copied out for every exit).
- `vcpu_events` are available in addition to `regs` and `sregs`.

For x86, the ‘`kvm_valid_regs`’ field of struct `kvm_run` is overloaded to function as an input bit-array field set by userspace to indicate the specific register sets to be copied out on the next exit.

To indicate when userspace has modified values that should be copied into the vCPU, the all architecture bitarray field, ‘`kvm_dirty_regs`’ must be set. This is done using the same bitflags as for the ‘`kvm_valid_regs`’ field. If the dirty bit is not set, then the register set values will not be copied into the vCPU even if they’ve been modified.

Unused bitfields in the bitarrays must be set to zero.

```
struct kvm_sync_regs {
    struct kvm_regs regs;
    struct kvm_sregs sregs;
    struct kvm_vcpu_events events;
};
```

## 6.75 KVM\_CAP\_PPC\_IRQ\_XIVE

**Architectures**

ppc

**Target**

vcpu

**Parameters**

args[0] is the XIVE device fd; args[1] is the XIVE CPU number (server ID) for this vcpu

This capability connects the vcpu to an in-kernel XIVE device.

### 1.1.7 7. Capabilities that can be enabled on VMs

There are certain capabilities that change the behavior of the virtual machine when enabled. To enable them, please see section 4.37. Below you can find a list of capabilities and what their effect on the VM is when enabling them.

The following information is provided along with the description:

**Architectures:**

which instruction set architectures provide this ioctl. x86 includes both i386 and x86\_64.

**Parameters:**

what parameters are accepted by the capability.

**Returns:**

the return value. General error numbers (EBADF, ENOMEM, EINVAL) are not detailed, but errors with specific meanings are.

## 7.1 KVM\_CAP\_PPC\_ENABLE\_HCALL

**Architectures**

ppc

**Parameters**

args[0] is the sPAPR hcall number; args[1] is 0 to disable, 1 to enable in-kernel handling

This capability controls whether individual sPAPR hypercalls (hcalls) get handled by the kernel or not. Enabling or disabling in-kernel handling of an hcall is effective across the VM. On creation, an initial set of hcalls are enabled for in-kernel handling, which consists of those hcalls for which in-kernel handlers were implemented before this capability was implemented. If disabled, the kernel will not attempt to handle the hcall, but will always exit to userspace to handle it. Note that it may not make sense to enable some and disable others of a group of related hcalls, but KVM does not prevent userspace from doing that.

If the hcall number specified is not one that has an in-kernel implementation, the KVM\_ENABLE\_CAP ioctl will fail with an EINVAL error.

## 7.2 KVM\_CAP\_S390\_USER\_SIGP

### Architectures

s390

### Parameters

none

This capability controls which SIGP orders will be handled completely in user space. With this capability enabled, all fast orders will be handled completely in the kernel:

- SENSE
- SENSE RUNNING
- EXTERNAL CALL
- EMERGENCY SIGNAL
- CONDITIONAL EMERGENCY SIGNAL

All other orders will be handled completely in user space.

Only privileged operation exceptions will be checked for in the kernel (or even in the hardware prior to interception). If this capability is not enabled, the old way of handling SIGP orders is used (partially in kernel and user space).

## 7.3 KVM\_CAP\_S390\_VECTOR\_REGISTERS

### Architectures

s390

### Parameters

none

### Returns

0 on success, negative value on error

Allows use of the vector registers introduced with z13 processor, and provides for the synchronization between host and user space. Will return -EINVAL if the machine does not support vectors.

## 7.4 KVM\_CAP\_S390\_USER\_STSI

### Architectures

s390

### Parameters

none

This capability allows post-handlers for the STSI instruction. After initial handling in the kernel, KVM exits to user space with KVM\_EXIT\_S390\_STSI to allow user space to insert further data.

Before exiting to userspace, kvm handlers should fill in s390\_stsi field of vcpu->run:

```
struct {
    __u64 addr;
    __u8 ar;
    __u8 reserved;
    __u8 fc;
    __u8 sel1;
    __u16 sel2;
} s390_stsi;

@addr - guest address of STSI SYSIB
@fc    - function code
@sel1  - selector 1
@sel2  - selector 2
@ar    - access register number
```

KVM handlers should exit to userspace with rc = -EREMOTE.

## 7.5 KVM\_CAP\_SPLIT\_IRQCHIP

### Architectures

x86

### Parameters

args[0] - number of routes reserved for userspace IOAPICs

### Returns

0 on success, -1 on error

Create a local apic for each processor in the kernel. This can be used instead of KVM\_CREATE\_IRQCHIP if the userspace VMM wishes to emulate the IOAPIC and PIC (and also the PIT, even though this has to be enabled separately).

This capability also enables in kernel routing of interrupt requests; when KVM\_CAP\_SPLIT\_IRQCHIP only routes of KVM\_IRQ\_ROUTING\_MSI type are used in the IRQ routing table. The first args[0] MSI routes are reserved for the IOAPIC pins. Whenever the LAPIC receives an EOI for these routes, a KVM\_EXIT\_IOAPIC\_EOI vmexit will be reported to userspace.

Fails if VCPU has already been created, or if the irqchip is already in the kernel (i.e. KVM\_CREATE\_IRQCHIP has already been called).

## 7.6 KVM\_CAP\_S390\_RI

### Architectures

s390

### Parameters

none

Allows use of runtime-instrumentation introduced with zEC12 processor. Will return -EINVAL if the machine does not support runtime-instrumentation. Will return -EBUSY if a VCPU has already been created.

## 7.7 KVM\_CAP\_X2APIC\_API

### Architectures

x86

### Parameters

args[0] - features that should be enabled

### Returns

0 on success, -EINVAL when args[0] contains invalid features

Valid feature flags in args[0] are:

```
#define KVM_X2APIC_API_USE_32BIT_IDS          (1ULL << 0)
#define KVM_X2APIC_API_DISABLE_BROADCAST QUIRK (1ULL << 1)
```

Enabling `KVM_X2APIC_API_USE_32BIT_IDS` changes the behavior of `KVM_SET_GSI_ROUTING`, `KVM_SIGNAL_MSI`, `KVM_SET_LAPIC`, and `KVM_GET_LAPIC`, allowing the use of 32-bit APIC IDs. See `KVM_CAP_X2APIC_API` in their respective sections.

`KVM_X2APIC_API_DISABLE_BROADCAST QUIRK` must be enabled for x2APIC to work in logical mode or with more than 255 VCPUs. Otherwise, KVM treats 0xff as a broadcast even in x2APIC mode in order to support physical x2APIC without interrupt remapping. This is undesirable in logical mode, where 0xff represents CPUs 0-7 in cluster 0.

## 7.8 KVM\_CAP\_S390\_USER\_INSTR0

### Architectures

s390

### Parameters

none

With this capability enabled, all illegal instructions 0x0000 (2 bytes) will be intercepted and forwarded to user space. User space can use this mechanism e.g. to realize 2-byte software breakpoints. The kernel will not inject an operating exception for these instructions, user space has to take care of that.

This capability can be enabled dynamically even if VCPUs were already created and are running.

## 7.9 KVM\_CAP\_S390\_GS

### Architectures

s390

### Parameters

none

### Returns

0 on success; -EINVAL if the machine does not support guarded storage; -EBUSY if a VCPU has already been created.

Allows use of guarded storage for the KVM guest.

### 7.10 KVM\_CAP\_S390\_AIS

#### Architectures

s390

#### Parameters

none

Allow use of adapter-interruption suppression. :Returns: 0 on success; -EBUSY if a VCPU has already been created.

### 7.11 KVM\_CAP\_PPC\_SMT

#### Architectures

ppc

#### Parameters

vsmt\_mode, flags

Enabling this capability on a VM provides userspace with a way to set the desired virtual SMT mode (i.e. the number of virtual CPUs per virtual core). The virtual SMT mode, `vsmt_mode`, must be a power of 2 between 1 and 8. On POWER8, `vsmt_mode` must also be no greater than the number of threads per subcore for the host. Currently flags must be 0. A successful call to enable this capability will result in `vsmt_mode` being returned when the `KVM_CAP_PPC_SMT` capability is subsequently queried for the VM. This capability is only supported by HV KVM, and can only be set before any VCPUs have been created. The `KVM_CAP_PPC_SMT_POSSIBLE` capability indicates which virtual SMT modes are available.

### 7.12 KVM\_CAP\_PPC\_FWNMI

#### Architectures

ppc

#### Parameters

none

With this capability a machine check exception in the guest address space will cause KVM to exit the guest with NMI exit reason. This enables QEMU to build error log and branch to guest kernel registered machine check handling routine. Without this capability KVM will branch to guests' 0x200 interrupt vector.



## 7.13 KVM\_CAP\_X86\_DISABLE\_EXITS

### Architectures

x86

### Parameters

args[0] defines which exits are disabled

### Returns

0 on success, -EINVAL when args[0] contains invalid exits

Valid bits in args[0] are:

```
#define KVM_X86_DISABLE_EXITS_MWAIT      (1 << 0)
#define KVM_X86_DISABLE_EXITS_HLT        (1 << 1)
#define KVM_X86_DISABLE_EXITS_PAUSE      (1 << 2)
#define KVM_X86_DISABLE_EXITS_CSTATE     (1 << 3)
```

Enabling this capability on a VM provides userspace with a way to no longer intercept some instructions for improved latency in some workloads, and is suggested when vCPUs are associated to dedicated physical CPUs. More bits can be added in the future; userspace can just pass the KVM\_CHECK\_EXTENSION result to KVM\_ENABLE\_CAP to disable all such vmexits.

Do not enable KVM\_FEATURE\_PV\_UNHALT if you disable HLT exits.

## 7.14 KVM\_CAP\_S390\_HPAGE\_1M

### Architectures

s390

### Parameters

none

### Returns

0 on success, -EINVAL if hpage module parameter was not set or cmma is enabled, or the VM has the KVM\_VM\_S390\_UCONTROL flag set

With this capability the KVM support for memory backing with 1m pages through hugetlbfs can be enabled for a VM. After the capability is enabled, cmma can't be enabled anymore and pfmf and the storage key interpretation are disabled. If cmma has already been enabled or the hpage module parameter is not set to 1, -EINVAL is returned.

While it is generally possible to create a huge page backed VM without this capability, the VM will not be able to run.

## 7.15 KVM\_CAP\_MSR\_PLATFORM\_INFO

### Architectures

x86

### Parameters

args[0] whether feature should be enabled or not

With this capability, a guest may read the MSR\_PLATFORM\_INFO MSR. Otherwise, a #GP would be raised when the guest tries to access. Currently, this capability does not enable write permissions of this MSR for the guest.

## 7.16 KVM\_CAP\_PPC\_NESTED\_HV

### Architectures

ppc

### Parameters

none

### Returns

0 on success, -EINVAL when the implementation doesn't support nested-HV virtualization.

HV-KVM on POWER9 and later systems allows for “nested-HV” virtualization, which provides a way for a guest VM to run guests that can run using the CPU's supervisor mode (privileged non-hypervisor state). Enabling this capability on a VM depends on the CPU having the necessary functionality and on the facility being enabled with a `kvm-hv` module parameter.

## 7.17 KVM\_CAP\_EXCEPTION\_PAYLOAD

### Architectures

x86

### Parameters

args[0] whether feature should be enabled or not

With this capability enabled, CR2 will not be modified prior to the emulated VM-exit when L1 intercepts a #PF exception that occurs in L2. Similarly, for `kvm-intel` only, DR6 will not be modified prior to the emulated VM-exit when L1 intercepts a #DB exception that occurs in L2. As a result, when `KVM_GET_VCPU_EVENTS` reports a pending #PF (or #DB) exception for L2, `exception.has_payload` will be set and the faulting address (or the new DR6 bits\*) will be reported in the `exception_payload` field. Similarly, when userspace injects a #PF (or #DB) into L2 using `KVM_SET_VCPU_EVENTS`, it is expected to set `exception.has_payload` and to put the faulting address - or the new DR6 bits<sup>3</sup> - in the `exception_payload` field.

This capability also enables `exception.pending` in `struct kvm_vcpu_events`, which allows userspace to distinguish between pending and injected exceptions.

## 7.18 KVM\_CAP\_MANUAL\_DIRTY\_LOG\_PROTECT2

---

<sup>3</sup> For the new DR6 bits, note that bit 16 is set iff the #DB exception will clear DR6.RTM.

**Architectures**

x86, arm, arm64, mips

**Parameters**

args[0] whether feature should be enabled or not

Valid flags are:

```
#define KVM_DIRTY_LOG_MANUAL_PROTECT_ENABLE (1 << 0)
#define KVM_DIRTY_LOG_INITIALLY_SET        (1 << 1)
```

With `KVM_DIRTY_LOG_MANUAL_PROTECT_ENABLE` is set, `KVM_GET_DIRTY_LOG` will not automatically clear and write-protect all pages that are returned as dirty. Rather, userspace will have to do this operation separately using `KVM_CLEAR_DIRTY_LOG`.

At the cost of a slightly more complicated operation, this provides better scalability and responsiveness for two reasons. First, `KVM_CLEAR_DIRTY_LOG` ioctl can operate on a 64-page granularity rather than requiring to sync a full memslot; this ensures that KVM does not take spinlocks for an extended period of time. Second, in some cases a large amount of time can pass between a call to `KVM_GET_DIRTY_LOG` and userspace actually using the data in the page. Pages can be modified during this time, which is inefficient for both the guest and userspace: the guest will incur a higher penalty due to write protection faults, while userspace can see false reports of dirty pages. Manual reprotection helps reducing this time, improving guest performance and reducing the number of dirty log false positives.

With `KVM_DIRTY_LOG_INITIALLY_SET` set, all the bits of the dirty bitmap will be initialized to 1 when created. This also improves performance because dirty logging can be enabled gradually in small chunks on the first call to `KVM_CLEAR_DIRTY_LOG`. `KVM_DIRTY_LOG_INITIALLY_SET` depends on `KVM_DIRTY_LOG_MANUAL_PROTECT_ENABLE` (it is also only available on x86 and arm64 for now).

`KVM_CAP_MANUAL_DIRTY_LOG_PROTECT2` was previously available under the name `KVM_CAP_MANUAL_DIRTY_LOG_PROTECT`, but the implementation had bugs that make it hard or impossible to use it correctly. The availability of `KVM_CAP_MANUAL_DIRTY_LOG_PROTECT2` signals that those bugs are fixed. Userspace should not try to use `KVM_CAP_MANUAL_DIRTY_LOG_PROTECT`.

**7.19 KVM\_CAP\_PPC\_SECURE\_GUEST****Architectures**

ppc

This capability indicates that KVM is running on a host that has ultravisor firmware and thus can support a secure guest. On such a system, a guest can ask the ultravisor to make it a secure guest, one whose memory is inaccessible to the host except for pages which are explicitly requested to be shared with the host. The ultravisor notifies KVM when a guest requests to become a secure guest, and KVM has the opportunity to veto the transition.

If present, this capability can be enabled for a VM, meaning that KVM will allow the transition to secure guest mode. Otherwise KVM will veto the transition.

## **7.20 KVM\_CAP\_HALT\_POLL**

### **Architectures**

all

### **Target**

VM

### **Parameters**

args[0] is the maximum poll time in nanoseconds

### **Returns**

0 on success; -1 on error

This capability overrides the kvm module parameter `halt_poll_ns` for the target VM.

VCPU polling allows a VCPU to poll for wakeup events instead of immediately scheduling during guest halts. The maximum time a VCPU can spend polling is controlled by the kvm module parameter `halt_poll_ns`. This capability allows the maximum halt time to specified on a per-VM basis, effectively overriding the module parameter for the target VM.

## **7.21 KVM\_CAP\_X86\_USER\_SPACE\_MSR**

### **Architectures**

x86

### **Target**

VM

### **Parameters**

args[0] contains the mask of `KVM_MSR_EXIT_REASON_*` events to report

### **Returns**

0 on success; -1 on error

This capability enables trapping of #GP invoking RDMSR and WRMSR instructions into user space.

When a guest requests to read or write an MSR, KVM may not implement all MSRs that are relevant to a respective system. It also does not differentiate by CPU type.

To allow more fine grained control over MSR handling, user space may enable this capability. With it enabled, MSR accesses that match the mask specified in args[0] and trigger a #GP event inside the guest by KVM will instead trigger `KVM_EXIT_X86_RDMSR` and `KVM_EXIT_X86_WRMSR` exit notifications which user space can then handle to implement model specific MSR handling and/or user notifications to inform a user that an MSR was not handled.

### 1.1.8 8. Other capabilities.

This section lists capabilities that give information about other features of the KVM implementation.

#### 8.1 KVM\_CAP\_PPC\_HWRNG

##### Architectures

ppc

This capability, if KVM\_CHECK\_EXTENSION indicates that it is available, means that the kernel has an implementation of the H\_RANDOM hypercall backed by a hardware random-number generator. If present, the kernel H\_RANDOM handler can be enabled for guest use with the KVM\_CAP\_PPC\_ENABLE\_HCALL capability.

#### 8.2 KVM\_CAP\_HYPERV\_SYNIC

##### Architectures

x86

This capability, if KVM\_CHECK\_EXTENSION indicates that it is available, means that the kernel has an implementation of the Hyper-V Synthetic interrupt controller(SynIC). Hyper-V SynIC is used to support Windows Hyper-V based guest paravirt drivers(VMBus).

In order to use SynIC, it has to be activated by setting this capability via KVM\_ENABLE\_CAP ioctl on the vcpu fd. Note that this will disable the use of APIC hardware virtualization even if supported by the CPU, as it's incompatible with SynIC auto-EOI behavior.

#### 8.3 KVM\_CAP\_PPC\_RADIX\_MMU

##### Architectures

ppc

This capability, if KVM\_CHECK\_EXTENSION indicates that it is available, means that the kernel can support guests using the radix MMU defined in Power ISA V3.00 (as implemented in the POWER9 processor).

#### 8.4 KVM\_CAP\_PPC\_HASH\_MMU\_V3

##### Architectures

ppc

This capability, if KVM\_CHECK\_EXTENSION indicates that it is available, means that the kernel can support guests using the hashed page table MMU defined in Power ISA V3.00 (as implemented in the POWER9 processor), including in-memory segment tables.

## 8.5 KVM\_CAP\_MIPS\_VZ

### Architectures

mips

This capability, if KVM\_CHECK\_EXTENSION on the main kvm handle indicates that it is available, means that full hardware assisted virtualization capabilities of the hardware are available for use through KVM. An appropriate KVM\_VM\_MIPS\_\* type must be passed to KVM\_CREATE\_VM to create a VM which utilises it.

If KVM\_CHECK\_EXTENSION on a kvm VM handle indicates that this capability is available, it means that the VM is using full hardware assisted virtualization capabilities of the hardware. This is useful to check after creating a VM with KVM\_VM\_MIPS\_DEFAULT.

The value returned by KVM\_CHECK\_EXTENSION should be compared against known values (see below). All other values are reserved. This is to allow for the possibility of other hardware assisted virtualization implementations which may be incompatible with the MIPS VZ ASE.

- |   |   |
|---|---|
| 0 | The trap & emulate implementation is in use to run guest code in user mode. Guest virtual memory segments are rearranged to fit the guest in the user mode address space. |
| 1 | The MIPS VZ ASE is in use, providing full hardware assisted virtualization, including standard guest virtual memory segments.   |

## 8.6 KVM\_CAP\_MIPS\_TE

### Architectures

mips

This capability, if KVM\_CHECK\_EXTENSION on the main kvm handle indicates that it is available, means that the trap & emulate implementation is available to run guest code in user mode, even if KVM\_CAP\_MIPS\_VZ indicates that hardware assisted virtualisation is also available. KVM\_VM\_MIPS\_TE (0) must be passed to KVM\_CREATE\_VM to create a VM which utilises it.

If KVM\_CHECK\_EXTENSION on a kvm VM handle indicates that this capability is available, it means that the VM is using trap & emulate.

## 8.7 KVM\_CAP\_MIPS\_64BIT

### Architectures

mips

This capability indicates the supported architecture type of the guest, i.e. the supported register and address width.

The values returned when this capability is checked by KVM\_CHECK\_EXTENSION on a kvm VM handle correspond roughly to the CP0\_Config.AT register field, and should be checked specifically against known values (see below). All other values are reserved.

- |   |   |
|---|---|
| 0 | MIPS32 or microMIPS32. Both registers and addresses are 32-bits wide. It will only be possible to run 32-bit guest code.  |
| 1 | MIPS64 or microMIPS64 with access only to 32-bit compatibility segments. Registers are 64-bits wide, but addresses are 32-bits wide. 64-bit guest code may run but cannot access MIPS64 memory segments. It will also be possible to run 32-bit guest code. |
| 2 | MIPS64 or microMIPS64 with access to all address segments. Both registers and addresses are 64-bits wide. It will be possible to run 64-bit or 32-bit guest code.   |

## 8.9 KVM\_CAP\_ARM\_USER\_IRQ

### Architectures

arm, arm64

This capability, if KVM\_CHECK\_EXTENSION indicates that it is available, means that if userspace creates a VM without an in-kernel interrupt controller, it will be notified of changes to the output level of in-kernel emulated devices, which can generate virtual interrupts, presented to the VM. For such VMs, on every return to userspace, the kernel updates the vcpu's `run->s.regs.device_irq_level` field to represent the actual output level of the device.

Whenever kvm detects a change in the device output level, kvm guarantees at least one return to userspace before running the VM. This exit could either be a KVM\_EXIT\_INTR or any other exit event, like KVM\_EXIT\_MMIO. This way, userspace can always sample the device output level and re-compute the state of the userspace interrupt controller. Userspace should always check the state of `run->s.regs.device_irq_level` on every kvm exit. The value in `run->s.regs.device_irq_level` can represent both level and edge triggered interrupt signals, depending on the device. Edge triggered interrupt signals will exit to userspace with the bit in `run->s.regs.device_irq_level` set exactly once per edge signal.

The field `run->s.regs.device_irq_level` is available independent of `run->kvm_valid_regs` or `run->kvm_dirty_regs` bits.

If KVM\_CAP\_ARM\_USER\_IRQ is supported, the KVM\_CHECK\_EXTENSION ioctl returns a number larger than 0 indicating the version of this capability is implemented and thereby which bits in `run->s.regs.device_irq_level` can signal values.

Currently the following bits are defined for the `device_irq_level` bitmap:

KVM\_CAP\_ARM\_USER\_IRQ >= 1:

KVM_ARM_DEV_EL1_VTIMER	-	EL1 virtual timer
KVM_ARM_DEV_EL1_PTIMER	-	EL1 physical timer
KVM_ARM_DEV_PMU	-	ARM PMU overflow interrupt signal

Future versions of kvm may implement additional events. These will get indicated by returning a higher number from KVM\_CHECK\_EXTENSION and will be listed above.

## **8.10 KVM\_CAP\_PPC\_SMT\_POSSIBLE**

### **Architectures**

ppc

Querying this capability returns a bitmap indicating the possible virtual SMT modes that can be set using KVM\_CAP\_PPC\_SMT. If bit N (counting from the right) is set, then a virtual SMT mode of  $2^N$  is available.

## **8.11 KVM\_CAP\_HYPERV\_SYNIC2**

### **Architectures**

x86

This capability enables a newer version of Hyper-V Synthetic interrupt controller (SynIC). The only difference with KVM\_CAP\_HYPERV\_SYNIC is that KVM doesn't clear SynIC message and event flags pages when they are enabled by writing to the respective MSRs.

## **8.12 KVM\_CAP\_HYPERV\_VP\_INDEX**

### **Architectures**

x86

This capability indicates that userspace can load HV\_X64\_MSR\_VP\_INDEX msr. Its value is used to denote the target vcpu for a SynIC interrupt. For compatibility, KVM initializes this msr to KVM's internal vcpu index. When this capability is absent, userspace can still query this msr's value.

## **8.13 KVM\_CAP\_S390\_AIS\_MIGRATION**

### **Architectures**

s390

### **Parameters**

none

This capability indicates if the flic device will be able to get/set the AIS states for migration via the KVM\_DEV\_FLIC\_AISM\_ALL attribute and allows to discover this without having to create a flic device.

## **8.14 KVM\_CAP\_S390\_PSW**

### **Architectures**

s390

This capability indicates that the PSW is exposed via the kvm\_run structure.



## 8.15 KVM\_CAP\_S390\_GMAP

### Architectures

s390

This capability indicates that the user space memory used as guest mapping can be anywhere in the user memory address space, as long as the memory slots are aligned and sized to a segment (1MB) boundary.

## 8.16 KVM\_CAP\_S390\_COW

### Architectures

s390

This capability indicates that the user space memory used as guest mapping can use copy-on-write semantics as well as dirty pages tracking via read-only page tables.

## 8.17 KVM\_CAP\_S390\_BPB

### Architectures

s390

This capability indicates that kvm will implement the interfaces to handle reset, migration and nested KVM for branch prediction blocking. The stfle facility 82 should not be provided to the guest without this capability.

## 8.18 KVM\_CAP\_HYPERV\_TLBFLUSH

### Architectures

x86

This capability indicates that KVM supports paravirtualized Hyper-V TLB Flush hypercalls: HvFlushVirtualAddressSpace, HvFlushVirtualAddressSpaceEx, HvFlushVirtualAddressList, HvFlushVirtualAddressListEx.

## 8.19 KVM\_CAP\_ARM\_INJECT\_ERROR\_ESR

### Architectures

arm, arm64

This capability indicates that userspace can specify (via the KVM\_SET\_VCPU\_EVENTS ioctl) the syndrome value reported to the guest when it takes a virtual SError interrupt exception. If KVM advertises this capability, userspace can only specify the ISS field for the ESR syndrome. Other parts of the ESR, such as the EC are generated by the CPU when the exception is taken. If this virtual SError is taken to EL1 using AArch64, this value will be reported in the ISS field of ESR\_ELx.

See KVM\_CAP\_VCPU\_EVENTS for more details.

## **8.20 KVM\_CAP\_HYPERV\_SEND\_IPI**

### **Architectures**

x86

This capability indicates that KVM supports paravirtualized Hyper-V IPI send hypercalls: `HvCallSendSyntheticClusterIpi`, `HvCallSendSyntheticClusterIpiEx`.

## **8.21 KVM\_CAP\_HYPERV\_DIRECT\_TLBFLUSH**

### **Architectures**

x86

This capability indicates that KVM running on top of Hyper-V hypervisor enables Direct TLB flush for its guests meaning that TLB flush hypercalls are handled by Level 0 hypervisor (Hyper-V) bypassing KVM. Due to the different ABI for hypercall parameters between Hyper-V and KVM, enabling this capability effectively disables all hypercall handling by KVM (as some KVM hypercall may be mistakenly treated as TLB flush hypercalls by Hyper-V) so userspace should disable KVM identification in `CPUID` and only exposes Hyper-V identification. In this case, guest thinks it's running on Hyper-V and only use Hyper-V hypercalls.

## **8.22 KVM\_CAP\_S390\_VCPU\_RESETS**

### **Architectures**

s390

This capability indicates that the `KVM_S390_NORMAL_RESET` and `KVM_S390_CLEAR_RESET` ioctls are available.

## **8.23 KVM\_CAP\_S390\_PROTECTED**

### **Architectures**

s390

This capability indicates that the Ultravisor has been initialized and KVM can therefore start protected VMs. This capability governs the `KVM_S390_PV_COMMAND` ioctl and the `KVM_MP_STATE_LOAD` `MP_STATE`. `KVM_SET_MP_STATE` can fail for protected guests when the state change is invalid.

## **8.24 KVM\_CAP\_STEAL\_TIME**

### **Architectures**

arm64, x86

This capability indicates that KVM supports steal time accounting. When steal time accounting is supported it may be enabled with architecture-specific interfaces. This capability and the architecture-specific interfaces must be consistent, i.e. if one says the feature is supported, then the other should as well and vice

versa. For arm64 see *Generic vcpu interface* “KVM\_ARM\_VCPU\_PVTIME\_CTRL” . For x86 see *KVM-specific MSRs* “MSR\_KVM\_STEAL\_TIME” .

## 8.25 KVM\_CAP\_S390\_DIAG318

### Architectures

s390

This capability enables a guest to set information about its control program (i.e. guest kernel type and version). The information is helpful during system/firmware service events, providing additional data about the guest environments running on the machine.

The information is associated with the DIAGNOSE 0x318 instruction, which sets an 8-byte value consisting of a one-byte Control Program Name Code (CPNC) and a 7-byte Control Program Version Code (CPVC). The CPNC determines what environment the control program is running in (e.g. Linux, z/VM...), and the CPVC is used for information specific to OS (e.g. Linux version, Linux distribution...)

If this capability is available, then the CPNC and CPVC can be synchronized between KVM and userspace via the sync regs mechanism (KVM\_SYNC\_DIAG318).

## 8.26 KVM\_CAP\_X86\_USER\_SPACE\_MSR

### Architectures

x86

This capability indicates that KVM supports deflection of MSR reads and writes to user space. It can be enabled on a VM level. If enabled, MSR accesses that would usually trigger a #GP by KVM into the guest will instead get bounced to user space through the KVM\_EXIT\_X86\_RDMSR and KVM\_EXIT\_X86\_WRMSR exit notifications.

## 8.27 KVM\_X86\_SET\_MSR\_FILTER

### Architectures

x86

This capability indicates that KVM supports that accesses to user defined MSRs may be rejected. With this capability exposed, KVM exports new VM ioctl KVM\_X86\_SET\_MSR\_FILTER which user space can call to specify bitmaps of MSR ranges that KVM should reject access to.

In combination with KVM\_CAP\_X86\_USER\_SPACE\_MSR, this allows user space to trap and emulate MSRs that are outside of the scope of KVM as well as limit the attack surface on KVM's MSR emulation code.

## **8.28 KVM\_CAP\_ENFORCE\_PV\_CPUID**

Architectures: x86

When enabled, KVM will disable paravirtual features provided to the guest according to the bits in the KVM\_CPUID\_FEATURES CPUID leaf (0x40000001). Otherwise, a guest may use the paravirtual features regardless of what has actually been exposed through the CPUID leaf.

### **1.1.9 9. Known KVM API problems**

In some cases, KVM's API has some inconsistencies or common pitfalls that userspace need to be aware of. This section details some of these issues.

Most of them are architecture specific, so the section is split by architecture.

#### **9.1. x86**

##### **KVM\_GET\_SUPPORTED\_CPUID issues**

In general, KVM\_GET\_SUPPORTED\_CPUID is designed so that it is possible to take its result and pass it directly to KVM\_SET\_CPUID2. This section documents some cases in which that requires some care.

##### **Local APIC features**

CPU[EAX=1]:ECX[21] (X2APIC) is reported by KVM\_GET\_SUPPORTED\_CPUID, but it can only be enabled if KVM\_CREATE\_IRQCHIP or KVM\_ENABLE\_CAP(KVM\_CAP\_IRQCHIP\_SPLIT) are used to enable in-kernel emulation of the local APIC.

The same is true for the KVM\_FEATURE\_PV\_UNHALT paravirtualized feature.

CPU[EAX=1]:ECX[24] (TSC\_DEADLINE) is not reported by KVM\_GET\_SUPPORTED\_CPUID. It can be enabled if KVM\_CAP\_TSC\_DEADLINE\_TIMER is present and the kernel has enabled in-kernel emulation of the local APIC.

##### **CPU topology**

Several CPUID values include topology information for the host CPU: 0x0b and 0x1f for Intel systems, 0x8000001e for AMD systems. Different versions of KVM return different values for this information and userspace should not rely on it. Currently they return all zeroes.

If userspace wishes to set up a guest topology, it should be careful that the values of these three leaves differ for each CPU. In particular, the APIC ID is found in EDX for all subleaves of 0x0b and 0x1f, and in EAX for 0x8000001e; the latter also encodes the core id and node id in bits 7:0 of EBX and ECX respectively.

## Obsolete ioctls and capabilities

KVM\_CAP\_DISABLE\_QUIRKS does not let userspace know which quirks are actually available. Use KVM\_CHECK\_EXTENSION(KVM\_CAP\_DISABLE\_QUIRKS2) instead if available.

## Ordering of KVM\_GET\_\*/KVM\_SET\_\* ioctls

TBD

## 1.2 Secure Encrypted Virtualization (SEV)

### 1.2.1 Overview

Secure Encrypted Virtualization (SEV) is a feature found on AMD processors.

SEV is an extension to the AMD-V architecture which supports running virtual machines (VMs) under the control of a hypervisor. When enabled, the memory contents of a VM will be transparently encrypted with a key unique to that VM.

The hypervisor can determine the SEV support through the CPUID instruction. The CPUID function 0x8000001f reports information related to SEV:

```
0x8000001f[eax]:
    Bit[1] indicates support for SEV
    ...
    [ecx]:
        Bits[31:0] Number of encrypted guests supported,
        ↳simultaneously
```

If support for SEV is present, MSR 0xc001\_0010 (MSR\_K8\_SYSCFG) and MSR 0xc001\_0015 (MSR\_K7\_HWCR) can be used to determine if it can be enabled:

```
0xc001_0010:
    Bit[23]    1 = memory encryption can be enabled
               0 = memory encryption can not be enabled

0xc001_0015:
    Bit[0]     1 = memory encryption can be enabled
               0 = memory encryption can not be enabled
```

When SEV support is available, it can be enabled in a specific VM by setting the SEV bit before executing VMRUN.:

```
VMCB[0x90]:
    Bit[1]     1 = SEV is enabled
               0 = SEV is disabled
```

SEV hardware uses ASIDs to associate a memory encryption key with a VM. Hence, the ASID for the SEV-enabled guests must be from 1 to a maximum value defined in the CPUID 0x8000001f[ecx] field.

### 1.2.2 SEV Key Management

The SEV guest key management is handled by a separate processor called the AMD Secure Processor (AMD-SP). Firmware running inside the AMD-SP provides a secure key management interface to perform common hypervisor activities such as encrypting bootstrap code, snapshot, migrating and debugging the guest. For more information, see the SEV Key Management spec [[api-spec](#)]

The main ioctl to access SEV is KVM\_MEMORY\_ENCRYPT\_OP. If the argument to KVM\_MEMORY\_ENCRYPT\_OP is NULL, the ioctl returns 0 if SEV is enabled and ENOTTY if it is disabled (on some older versions of Linux, the ioctl runs normally even with a NULL argument, and therefore will likely return -EFAULT). If non-NULL, the argument to KVM\_MEMORY\_ENCRYPT\_OP must be a struct kvm\_sev\_cmd:

```
struct kvm_sev_cmd {
    __u32 id;
    __u64 data;
    __u32 error;
    __u32 sev_fd;
};
```

The id field contains the subcommand, and the data field points to another struct containing arguments specific to command. The sev\_fd should point to a file descriptor that is opened on the /dev/sev device, if needed (see individual commands).

On output, error is zero on success, or an error code. Error codes are defined in <linux/psp-dev.h>.

KVM implements the following commands to support common lifecycle events of SEV guests, such as launching, running, snapshotting, migrating and decommissioning.

#### 1. KVM\_SEV\_INIT

The KVM\_SEV\_INIT command is used by the hypervisor to initialize the SEV platform context. In a typical workflow, this command should be the first command issued.

Returns: 0 on success, -negative on error

## 2. KVM\_SEV\_LAUNCH\_START

The KVM\_SEV\_LAUNCH\_START command is used for creating the memory encryption context. To create the encryption context, user must provide a guest policy, the owner's public Diffie-Hellman (PDH) key and session information.

Parameters: struct kvm\_sev\_launch\_start (in/out)

Returns: 0 on success, -negative on error

```
struct kvm_sev_launch_start {
    __u32 handle;           /* if zero then firmware creates a ↵
    ↪new handle */
    __u32 policy;           /* guest's policy */
    __u64 dh_uaddr;         /* userspace address pointing to ↵
    ↪the guest owner's PDH key */
    __u32 dh_len;
    __u64 session_addr;     /* userspace address which points ↵
    ↪to the guest session information */
    __u32 session_len;
};
```

On success, the 'handle' field contains a new handle and on error, a negative value.

KVM\_SEV\_LAUNCH\_START requires the sev\_fd field to be valid.

For more details, see SEV spec Section 6.2.

## 3. KVM\_SEV\_LAUNCH\_UPDATE\_DATA

The KVM\_SEV\_LAUNCH\_UPDATE\_DATA is used for encrypting a memory region. It also calculates a measurement of the memory contents. The measurement is a signature of the memory contents that can be sent to the guest owner as an attestation that the memory was encrypted correctly by the firmware.

Parameters (in): struct kvm\_sev\_launch\_update\_data

Returns: 0 on success, -negative on error

```
struct kvm_sev_launch_update {
    __u64 uaddr;           /* userspace address to be encrypted (must ↵
    ↪be 16-byte aligned) */
    __u32 len;             /* length of the data to be encrypted (must ↵
    ↪be 16-byte aligned) */
};
```

For more details, see SEV spec Section 6.3.

## 4. KVM\_SEV\_LAUNCH\_MEASURE

The KVM\_SEV\_LAUNCH\_MEASURE command is used to retrieve the measurement of the data encrypted by the KVM\_SEV\_LAUNCH\_UPDATE\_DATA command. The guest owner may wait to provide the guest with confidential information until it can verify the measurement. Since the guest owner knows the initial contents of the guest at boot, the measurement can be verified by comparing it to what the guest owner expects.

Parameters (in): struct kvm\_sev\_launch\_measure

Returns: 0 on success, -negative on error

```
struct kvm_sev_launch_measure {
    __u64 uaddr;    /* where to copy the measurement */
    __u32 len;      /* length of measurement blob */
};
```

For more details on the measurement verification flow, see SEV spec Section 6.4.

## 5. KVM\_SEV\_LAUNCH\_FINISH

After completion of the launch flow, the KVM\_SEV\_LAUNCH\_FINISH command can be issued to make the guest ready for the execution.

Returns: 0 on success, -negative on error

## 6. KVM\_SEV\_GUEST\_STATUS

The KVM\_SEV\_GUEST\_STATUS command is used to retrieve status information about a SEV-enabled guest.

Parameters (out): struct kvm\_sev\_guest\_status

Returns: 0 on success, -negative on error

```
struct kvm_sev_guest_status {
    __u32 handle;    /* guest handle */
    __u32 policy;    /* guest policy */
    __u8 state;      /* guest state (see enum below) */
};
```

SEV guest state:

```
enum {
    SEV_STATE_INVALID = 0;
    SEV_STATE_LAUNCHING,    /* guest is currently being launched */
    SEV_STATE_SECRET,       /* guest is being launched and ready to
    ↪ accept the ciphertext data */
    SEV_STATE_RUNNING,      /* guest is fully launched and running */
    SEV_STATE_RECEIVING,    /* guest is being migrated in from another
    ↪ SEV machine */
};
```

(continues on next page)



(continued from previous page)

```
SEV_STATE_SENDING      /* guest is getting migrated out to another
↳ SEV machine */
};
```

## 7. KVM\_SEV\_DBG\_DECRYPT

The KVM\_SEV\_DEBUG\_DECRYPT command can be used by the hypervisor to request the firmware to decrypt the data at the given memory region.

Parameters (in): struct kvm\_sev\_dbg

Returns: 0 on success, -negative on error

```
struct kvm_sev_dbg {
    __u64 src_uaddr;      /* userspace address of data to
↳ decrypt */
    __u64 dst_uaddr;      /* userspace address of destination
↳ */
    __u32 len;           /* length of memory region to
↳ decrypt */
};
```

The command returns an error if the guest policy does not allow debugging.

## 8. KVM\_SEV\_DBG\_ENCRYPT

The KVM\_SEV\_DEBUG\_ENCRYPT command can be used by the hypervisor to request the firmware to encrypt the data at the given memory region.

Parameters (in): struct kvm\_sev\_dbg

Returns: 0 on success, -negative on error

```
struct kvm_sev_dbg {
    __u64 src_uaddr;      /* userspace address of data to
↳ encrypt */
    __u64 dst_uaddr;      /* userspace address of destination
↳ */
    __u32 len;           /* length of memory region to
↳ encrypt */
};
```

The command returns an error if the guest policy does not allow debugging.

## 9. KVM\_SEV\_LAUNCH\_SECRET

The KVM\_SEV\_LAUNCH\_SECRET command can be used by the hypervisor to inject secret data after the measurement has been validated by the guest owner.

Parameters (in): struct kvm\_sev\_launch\_secret

Returns: 0 on success, -negative on error

```
struct kvm_sev_launch_secret {
    __u64 hdr_uaddr;          /* userspace address containing the
    ↪ packet header */
    __u32 hdr_len;

    __u64 guest_uaddr;        /* the guest memory region where
    ↪ the secret should be injected */
    __u32 guest_len;

    __u64 trans_uaddr;        /* the hypervisor memory region
    ↪ which contains the secret */
    __u32 trans_len;
};
```

### 1.2.3 References

See [white-paper], [api-spec], [amd-apm] and [kvm-forum] for more info.

## 1.3 KVM CPUID bits

### Author

Glauber Costa <glommer@gmail.com>

A guest running on a kvm host, can check some of its features using cpuid. This is not always guaranteed to work, since userspace can mask-out some, or even all KVM-related cpuid features before launching a guest.

KVM cpuid functions are:

function: KVM\_CPUID\_SIGNATURE (0x40000000)

returns:

```
eax = 0x40000001
ebx = 0x4b4d564b
ecx = 0x564b4d56
edx = 0x4d
```

Note that this value in ebx, ecx and edx corresponds to the string “KVMKVMKVM”. The value in eax corresponds to the maximum cpuid function present in this leaf, and will be updated if more functions are added in the future. Note also that old hosts set eax value to 0x0. This should be interpreted as if the value was 0x40000001. This function queries the presence of KVM cpuid leafs.

function: define KVM\_CPUID\_FEATURES (0x40000001)

returns:

ebx, ecx  
 eax = an OR'ed group of (1 << flag)

where flag is defined as below:

flag	valu	meaning
KVM_FEATURE_CL	0	kvmclock available at msrs 0x11 and 0x12
KVM_FEATURE_NC	1	not necessary to perform delays on PIO operations
KVM_FEATURE_M	2	deprecated
KVM_FEATURE_CL	3	kvmclock available at msrs 0x4b564d00 and 0x4b564d01
KVM_FEATURE_AS	4	async pf can be enabled by writing to msr 0x4b564d02
KVM_FEATURE_ST	5	steal time can be enabled by writing to msr 0x4b564d03
KVM_FEATURE_PV	6	paravirtualized end of interrupt handler can be enabled by writing to msr 0x4b564d04
KVM_FEATURE_PV	7	guest checks this feature bit before enabling paravirtualized spinlock support
KVM_FEATURE_PV	9	guest checks this feature bit before enabling paravirtualized tlb flush
KVM_FEATURE_AS	10	paravirtualized async PF VM EXIT can be enabled by setting bit 2 when writing to msr 0x4b564d02
KVM_FEATURE_PV	11	guest checks this feature bit before enabling paravirtualized send IPIs
KVM_FEATURE_PO	12	host-side polling on HLT can be disabled by writing to msr 0x4b564d05.
KVM_FEATURE_PV	13	guest checks this feature bit before using paravirtualized sched yield.
KVM_FEATURE_AS	14	guest checks this feature bit before using the second async pf control msr 0x4b564d06 and async pf acknowledgment msr 0x4b564d07.
KVM_FEATURE_M	15	guest checks this feature bit before using extended destination ID bits in MSI address bits 11-5.
KVM_FEATURE_CL	24	host will warn if no guest-side per-cpu warps are expected in kvmclock

edx = an OR'ed group of (1 << flag)

Where flag here is defined as below:

flag	valu	meaning
KVM_HINTS_	0	guest checks this feature bit to determine that vCPUs are never preempted for an unlimited time allowing optimizations

## 1.4 The KVM halt polling system

The KVM halt polling system provides a feature within KVM whereby the latency of a guest can, under some circumstances, be reduced by polling in the host for some time period after the guest has elected to no longer run by cedeing. That is, when a guest vcpu has ceded, or in the case of powerpc when all of the vcpus of a single vcore have ceded, the host kernel polls for wakeup conditions before giving up the cpu to the scheduler in order to let something else run.

Polling provides a latency advantage in cases where the guest can be run again very quickly by at least saving us a trip through the scheduler, normally on the order of a few micro-seconds, although performance benefits are workload dependant. In the event that no wakeup source arrives during the polling interval or some other task on the runqueue is runnable the scheduler is invoked. Thus halt polling is especially useful on workloads with very short wakeup periods where the time spent halt polling is minimised and the time savings of not invoking the scheduler are distinguishable.

The generic halt polling code is implemented in:

virt/kvm/kvm\_main.c: `kvm_vcpu_block()`

The powerpc kvm-hv specific case is implemented in:

arch/powerpc/kvm/book3s\_hv.c: `kvmppc_vcore_blocked()`

### 1.4.1 Halt Polling Interval

The maximum time for which to poll before invoking the scheduler, referred to as the halt polling interval, is increased and decreased based on the perceived effectiveness of the polling in an attempt to limit pointless polling. This value is stored in either the vcpu struct:

`kvm_vcpu->halt_poll_ns`

or in the case of powerpc kvm-hv, in the vcore struct:

`kvmppc_vcore->halt_poll_ns`

Thus this is a per vcpu (or vcore) value.

During polling if a wakeup source is received within the halt polling interval, the interval is left unchanged. In the event that a wakeup source isn't received during the polling interval (and thus schedule is invoked) there are two options, either the polling interval and total block time[0] were less than the global max polling interval (see module params below), or the total block time was greater than the global max polling interval.

In the event that both the polling interval and total block time were less than the global max polling interval then the polling interval can be increased in the hope that next time during the longer polling interval the wake up source will be received while the host is polling and the latency benefits will be received. The polling interval is grown in the function `grow_halt_poll_ns()` and is multiplied by the module parameters `halt_poll_ns_grow` and `halt_poll_ns_grow_start`.

In the event that the total block time was greater than the global max polling interval then the host will never poll for long enough (limited by the global max) to wakeup during the polling interval so it may as well be shrunk in order to avoid pointless polling. The polling interval is shrunk in the function `shrink_halt_poll_ns()` and is divided by the module parameter `halt_poll_ns_shrink`, or set to 0 iff `halt_poll_ns_shrink == 0`.

It is worth noting that this adjustment process attempts to hone in on some steady state polling interval but will only really do a good job for wakeups which come at an approximately constant rate, otherwise there will be constant adjustment of the polling interval.

**[0] total block time:**

the time between when the halt polling function is invoked and a wakeup source received (irrespective of whether the scheduler is invoked within that function).

### 1.4.2 Module Parameters

The `kvm` module has 3 tuneable module parameters to adjust the global max polling interval as well as the rate at which the polling interval is grown and shrunk. These variables are defined in `include/linux/kvm_host.h` and as module parameters in `virt/kvm/kvm_main.c`, or `arch/powerpc/kvm/book3s_hv.c` in the powerpc `kvm-hv` case.

Module Parameter	Description	Default Value
<code>halt_poll_ns</code>	The global max polling interval which defines the ceiling value of the polling interval for each vcpu.	<code>KVM_HALT_POLL_NS_DEFAULT</code> (per arch value)
<code>halt_poll_ns</code>	The value by which the halt polling interval is multiplied in the <code>grow_halt_poll_ns()</code> function.	2
<code>halt_poll_ns</code>	The initial value to grow to from zero in the <code>grow_halt_poll_ns()</code> function.	10000
<code>halt_poll_ns</code>	The value by which the halt polling interval is divided in the <code>shrink_halt_poll_ns()</code> function.	0

These module parameters can be set from the debugfs files in:

`/sys/module/kvm/parameters/`

**Note:** that these module parameters are system wide values and are not able to be tuned on a per vm basis.

### 1.4.3 Further Notes

- Care should be taken when setting the `halt_poll_ns` module parameter as a large value has the potential to drive the cpu usage to 100% on a machine which would be almost entirely idle otherwise. This is because even if a guest has wakeups during which very little work is done and which are quite far apart, if the period is shorter than the global max polling interval (`halt_poll_ns`) then the host will always poll for the entire block time and thus cpu utilisation will go to 100%.
- Halt polling essentially presents a trade off between power usage and latency and the module parameters should be used to tune the affinity for this. Idle cpu time is essentially converted to host kernel time with the aim of decreasing latency when entering the guest.
- Halt polling will only be conducted by the host when no other tasks are runnable on that cpu, otherwise the polling will cease immediately and schedule will be invoked to allow that other task to run. Thus this doesn't allow a guest to deny of service the cpu.

## 1.5 Linux KVM Hypercall

### X86:

KVM Hypercalls have a three-byte sequence of either the `vmcall` or the `vmmcall` instruction. The hypervisor can replace it with instructions that are guaranteed to be supported.

Up to four arguments may be passed in `rbx`, `rcx`, `rdx`, and `rsi` respectively. The hypercall number should be placed in `rax` and the return value will be placed in `rax`. No other registers will be clobbered unless explicitly stated by the particular hypercall.

### S390:

R2-R7 are used for parameters 1-6. In addition, R1 is used for hypercall number. The return value is written to R2.

S390 uses `diagnose` instruction as hypercall (0x500) along with hypercall number in R1.

For further information on the S390 `diagnose` call as supported by KVM, refer to *The s390 DIAGNOSE call on KVM*.

### PowerPC:

It uses R3-R10 and hypercall number in R11. R4-R11 are used as output registers. Return value is placed in R3.

KVM hypercalls uses 4 byte opcode, that are patched with 'hypercall-instructions' property inside the device tree's `/hypervisor` node. For more information refer to *The PPC KVM paravirtual interface*

### MIPS:

KVM hypercalls use the `HYPCALL` instruction with code 0 and the hypercall number in `$2` (`v0`). Up to four arguments may be placed in `$4-$7` (`a0-a3`) and the return value is placed in `$2` (`v0`).

### 1.5.1 KVM Hypercalls Documentation

The template for each hypercall is: 1. Hypercall name. 2. Architecture(s) 3. Status (deprecated, obsolete, active) 4. Purpose

#### 1. KVM\_HC\_VAPIC\_POLL\_IRQ

**Architecture**

x86

**Status**

active

**Purpose**

Trigger guest exit so that the host can check for pending interrupts on reentry.

#### 2. KVM\_HC\_MMU\_OP

**Architecture**

x86

**Status**

deprecated.

**Purpose**

Support MMU operations such as writing to PTE, flushing TLB, release PT.

#### 3. KVM\_HC\_FEATURES

**Architecture**

PPC

**Status**

active

**Purpose**

Expose hypercall availability to the guest. On x86 platforms, cpuid used to enumerate which hypercalls are available. On PPC, either device tree based lookup ( which is also what EPAPR dictates) OR KVM specific enumeration mechanism (which is this hypercall) can be used.

## **4. KVM\_HC\_PPC\_MAP\_MAGIC\_PAGE**

**Architecture**

PPC

**Status**

active

**Purpose**

To enable communication between the hypervisor and guest there is a shared page that contains parts of supervisor visible register state. The guest can map this shared page to access its supervisor register through memory using this hypercall.

## **5. KVM\_HC\_KICK\_CPU**

**Architecture**

x86

**Status**

active

**Purpose**

Hypercall used to wakeup a vcpu from HLT state

**Usage example**

A vcpu of a paravirtualized guest that is busywaiting in guest kernel mode for an event to occur (ex: a spinlock to become available) can execute HLT instruction once it has busy-waited for more than a threshold time-interval. Execution of HLT instruction would cause the hypervisor to put the vcpu to sleep until occurrence of an appropriate event. Another vcpu of the same guest can wakeup the sleeping vcpu by issuing KVM\_HC\_KICK\_CPU hypercall, specifying APIC ID (a1) of the vcpu to be woken up. An additional argument (a0) is used in the hypercall for future use.

## **6. KVM\_HC\_CLOCK\_PAIRING**

**Architecture**

x86

**Status**

active

**Purpose**

Hypercall used to synchronize host and guest clocks.

Usage:

a0: guest physical address where host copies “struct kvm\_clock\_offset” structure.

a1: clock\_type, ATM only KVM\_CLOCK\_PAIRING\_WALLCLOCK (0) is supported (corresponding to the host’s CLOCK\_REALTIME clock).



```

struct kvm_clock_pairing {
    __s64 sec;
    __s64 nsec;
    __u64 tsc;
    __u32 flags;
    __u32 pad[9];
};

```

**Where:**

- sec: seconds from clock\_type clock.
- nsec: nanoseconds from clock\_type clock.
- tsc: guest TSC value used to calculate sec/nsec pair
- flags: flags, unused (0) at the moment.

The hypercall lets a guest compute a precise timestamp across host and guest. The guest can use the returned TSC value to compute the CLOCK\_REALTIME for its clock, at the same instant.

Returns KVM\_EOPNOTSUPP if the host does not use TSC clocksource, or if clock type is different than KVM\_CLOCK\_PAIRING\_WALLCLOCK.

**6. KVM\_HC\_SEND\_IPI****Architecture**

x86

**Status**

active

**Purpose**

Send IPIs to multiple vCPUs.

- a0: lower part of the bitmap of destination APIC IDs
- a1: higher part of the bitmap of destination APIC IDs
- a2: the lowest APIC ID in bitmap
- a3: APIC ICR

The hypercall lets a guest send multicast IPIs, with at most 128 128 destinations per hypercall in 64-bit mode and 64 vCPUs per hypercall in 32-bit mode. The destinations are represented by a bitmap contained in the first two arguments (a0 and a1). Bit 0 of a0 corresponds to the APIC ID in the third argument (a2), bit 1 corresponds to the APIC ID a2+1, and so on.

Returns the number of CPUs to which the IPIs were delivered successfully.

## 7. KVM\_HC\_SCHED\_YIELD

### Architecture

x86

### Status

active

### Purpose

Hypercall used to yield if the IPI target vCPU is preempted

a0: destination APIC ID

### Usage example

When sending a call-function IPI-many to vCPUs, yield if any of the IPI target vCPUs was preempted.

## 1.6 KVM Lock Overview

### 1.6.1 1. Acquisition Orders

The acquisition orders for mutexes are as follows:

- `kvm->lock` is taken outside `vcpu->mutex`
- `kvm->lock` is taken outside `kvm->slots_lock` and `kvm->irq_lock`
- `kvm->slots_lock` is taken outside `kvm->irq_lock`, though acquiring them together is quite rare.

On x86, `vcpu->mutex` is taken outside `kvm->arch.hyperv.hv_lock`.

Everything else is a leaf: no other lock is taken inside the critical sections.

### 1.6.2 2. Exception

Fast page fault:

Fast page fault is the fast path which fixes the guest page fault out of the mmu-lock on x86. Currently, the page fault can be fast in one of the following two cases:

1. Access Tracking: The SPTE is not present, but it is marked for access tracking i.e. the `SPTE_SPECIAL_MASK` is set. That means we need to restore the saved R/X bits. This is described in more detail later below.
2. Write-Protection: The SPTE is present and the fault is caused by write-protect. That means we just need to change the W bit of the spte.

What we use to avoid all the race is the `SPTE_HOST_WRITEABLE` bit and `SPTE_MMU_WRITEABLE` bit on the spte:

- `SPTE_HOST_WRITEABLE` means the gfn is writable on host.
- `SPTE_MMU_WRITEABLE` means the gfn is writable on mmu. The bit is set when the gfn is writable on guest mmu and it is not write-protected by shadow page write-protection.

On fast page fault path, we will use `cmpxchg` to atomically set the `spte W` bit if `spte.SPTE_HOST_WRITEABLE = 1` and `spte.SPTE_WRITE_PROTECT = 1`, or restore the saved R/X bits if `VMX_EPT_TRACK_ACCESS` mask is set, or both. This is safe because whenever changing these bits can be detected by `cmpxchg`.

But we need carefully check these cases:

#### 1) The mapping from gfn to pfn

The mapping from gfn to pfn may be changed since we can only ensure the pfn is not changed during `cmpxchg`. This is a ABA problem, for example, below case will happen:

At the beginning:

```
gpte = gfn1
gfn1 is mapped to pfn1 on host
spte is the shadow page table entry corresponding with gpte and
spte = pfn1
```

On fast page fault path:

CPU 0:	CPU 1:
--------	--------

```
old_spte = *spte;
```

pfn1 is swapped out:

```
spte = 0;
```

pfn1 is re-allocated for gfn2.  
gpte is changed to point to gfn2 by the guest:

```
spte = pfn1;
```

```
if (cmpxchg(spte, old_spte, old_spte+W)
    mark_page_dirty(vcpu->kvm, gfn1)
    OOPS!!!
```

We dirty-log for gfn1, that means gfn2 is lost in dirty-bitmap.

For direct sp, we can easily avoid it since the `spte` of direct sp is fixed to gfn. For indirect sp, we disabled fast page fault for simplicity.

A solution for indirect sp could be to pin the gfn, for example via `kvm_vcpu_gfn_to_pfn_atomic`, before the `cmpxchg`. After the pinning:

- We have held the refcount of pfn that means the pfn can not be freed and be reused for another gfn.
- The pfn is writable and therefore it cannot be shared between different gfn by KSM.

Then, we can ensure the dirty bitmaps is correctly set for a gfn.

#### 2) Dirty bit tracking

In the origin code, the spte can be fast updated (non-atomically) if the spte is read-only and the Accessed bit has already been set since the Accessed bit and Dirty bit can not be lost.

But it is not true after fast page fault since the spte can be marked writable between reading spte and updating spte. Like below case:

At the beginning:

```
spte.W = 0
spte.Accessed = 1
```

CPU 0:

In mmu\_spte\_clear\_track\_bits():

```
old_spte = *spte;

/* 'if' condition is satisfied.
↳ */
if (old_spte.Accessed == 1 &&
    old_spte.W == 0)
    spte = 0ull;
```

CPU 1:

on fast page fault path:

```
spte.W = 1
```

memory write on the spte:

```
spte.Dirty = 1
```

```
else
    old_spte = xchg(spte, 0ull)
if (old_spte.Accessed == 1)
    kvm_set_pfn_accessed(spte.
↳ pfn);
if (old_spte.Dirty == 1)
    kvm_set_pfn_dirty(spte.pfn);
OOPS!!!
```

The Dirty bit is lost in this case.

In order to avoid this kind of issue, we always treat the spte as “volatile” if it can be updated out of mmu-lock, see `spte_has_volatile_bits()`, it means, the spte is always atomically updated in this case.

### 3) flush tlbs due to spte updated

If the spte is updated from writable to readonly, we should flush all TLBs, otherwise `rmap_write_protect` will find a read-only spte, even though the writable spte might be cached on a CPU’ s TLB.

As mentioned before, the spte can be updated to writable out of mmu-lock on fast page fault path, in order to easily audit the path, we see if TLBs need be flushed

caused by this reason in `mmu_spte_update()` since this is a common function to update spte (present -> present).

Since the spte is “volatile” if it can be updated out of mmu-lock, we always atomically update the spte, the race caused by fast page fault can be avoided, See the comments in `spte_has_volatile_bits()` and `mmu_spte_update()`.

Lockless Access Tracking:

This is used for Intel CPUs that are using EPT but do not support the EPT A/D bits. In this case, when the KVM MMU notifier is called to track accesses to a page (via `kvm_mmu_notifier_clear_flush_young()`), it marks the PTE as not-present by clearing the RWX bits in the PTE and storing the original R & X bits in some unused/ignored bits. In addition, the `SPTE_SPECIAL_MASK` is also set on the PTE (using the ignored bit 62). When the VM tries to access the page later on, a fault is generated and the fast page fault mechanism described above is used to atomically restore the PTE to a Present state. The W bit is not saved when the PTE is marked for access tracking and during restoration to the Present state, the W bit is set depending on whether or not it was a write access. If it wasn't, then the W bit will remain clear until a write access happens, at which time it will be set using the Dirty tracking mechanism described above.

### 1.6.3 3. Reference

**Name**

`kvm_lock`

**Type**

`mutex`

**Arch**

any

**Protects**

- `vm_list`

**Name**

`kvm_count_lock`

**Type**

`raw_spinlock_t`

**Arch**

any

**Protects**

- hardware virtualization enable/disable

**Comment**

‘raw’ because hardware enabling/disabling must be atomic /wrt migration.

**Name**

`kvm_arch::tsc_write_lock`

**Type**

raw\_spinlock

**Arch**

x86

**Protects**

- kvm\_arch::{last\_tsc\_write,last\_tsc\_nsec,last\_tsc\_offset}
- tsc offset in vmcb

**Comment**

‘raw’ because updating the tsc offsets must not be preempted.

**Name**

kvm->mmu\_lock

**Type**

spinlock\_t

**Arch**

any

**Protects**

-shadow page/shadow tlb entry

**Comment**

it is a spinlock since it is used in mmu notifier.

**Name**

kvm->srcu

**Type**

srcu lock

**Arch**

any

**Protects**

- kvm->memslots
- kvm->buses

**Comment**

The srcu read lock must be held while accessing memslots (e.g. when using gfn\_to\_\* functions) and while accessing in-kernel MMIO/PIO address->device structure mapping (kvm->buses). The srcu index can be stored in kvm\_vcpu->srcu\_idx per vcpu if it is needed by multiple functions.

**Name**

blocked\_vcpu\_on\_cpu\_lock

**Type**

spinlock\_t

**Arch**

x86

**Protects**

`blocked_vcpu_on_cpu`

**Comment**

This is a per-CPU lock and it is used for VT-d posted-interrupts. When VT-d posted-interrupts is supported and the VM has assigned devices, we put the blocked vCPU on the list `blocked_vcpu_on_cpu` protected by `blocked_vcpu_on_cpu_lock`, when VT-d hardware issues wakeup notification event since external interrupts from the assigned devices happens, we will find the vCPU on the list to wakeup.

## 1.7 The x86 kvm shadow mmu

The mmu (in `arch/x86/kvm`, files `mmu.[ch]` and `paging_tmpl.h`) is responsible for presenting a standard x86 mmu to the guest, while translating guest physical addresses to host physical addresses.

The mmu code attempts to satisfy the following requirements:

- **correctness:**  
the guest should not be able to determine that it is running on an emulated mmu except for timing (we attempt to comply with the specification, not emulate the characteristics of a particular implementation such as tlb size)
- **security:**  
the guest must not be able to touch host memory not assigned to it
- **performance:**  
minimize the performance penalty imposed by the mmu
- **scaling:**  
need to scale to large memory and large vcpu guests
- **hardware:**  
support the full range of x86 virtualization hardware
- **integration:**  
Linux memory management code must be in control of guest memory so that swapping, page migration, page merging, transparent hugepages, and similar features work without change
- **dirty tracking:**  
report writes to guest memory to enable live migration and framebuffer-based displays
- **footprint:**  
keep the amount of pinned kernel memory low (most memory should be shrinkable)
- **reliability:**  
avoid multipage or GFP\_ATOMIC allocations

### 1.7.1 Acronyms

pfn	host page frame number
hpa	host physical address
hva	host virtual address
gfn	guest frame number
gpa	guest physical address
gva	guest virtual address
ngpa	nested guest physical address
ngva	nested guest virtual address
pte	page table entry (used also to refer generically to paging structure entries)
gpte	guest pte (referring to gfns)
spte	shadow pte (referring to pfns)
tdp	two dimensional paging (vendor neutral term for NPT and EPT)

### 1.7.2 Virtual and real hardware supported

The mmu supports first-generation mmu hardware, which allows an atomic switch of the current paging mode and cr3 during guest entry, as well as two-dimensional paging (AMD' s NPT and Intel' s EPT). The emulated hardware it exposes is the traditional 2/3/4 level x86 mmu, with support for global pages, pae, pse, pse36, cr0.wp, and 1GB pages. Emulated hardware also able to expose NPT capable hardware on NPT capable hosts.

### 1.7.3 Translation

The primary job of the mmu is to program the processor' s mmu to translate addresses for the guest. Different translations are required at different times:

- when guest paging is disabled, we translate guest physical addresses to host physical addresses (gpa->hpa)
- when guest paging is enabled, we translate guest virtual addresses, to guest physical addresses, to host physical addresses (gva->gpa->hpa)
- when the guest launches a guest of its own, we translate nested guest virtual addresses, to nested guest physical addresses, to guest physical addresses, to host physical addresses (ngva->ngpa->gpa->hpa)

The primary challenge is to encode between 1 and 3 translations into hardware that support only 1 (traditional) and 2 (tdp) translations. When the number of required translations matches the hardware, the mmu operates in direct mode; otherwise it operates in shadow mode (see below).



### 1.7.4 Memory

Guest memory (gpa) is part of the user address space of the process that is using kvm. Userspace defines the translation between guest addresses and user addresses (gpa->hva); note that two gpas may alias to the same hva, but not vice versa.

These hvas may be backed using any method available to the host: anonymous memory, file backed memory, and device memory. Memory might be paged by the host at any time.

### 1.7.5 Events

The mmu is driven by events, some from the guest, some from the host.

Guest generated events:

- writes to control registers (especially cr3)
- invlpg/invlpga instruction execution
- access to missing or protected translations

Host generated events:

- changes in the gpa->hpa translation (either through gpa->hva changes or through hva->hpa changes)
- memory pressure (the shrinker)

### 1.7.6 Shadow pages

The principal data structure is the shadow page, 'struct kvm\_mmu\_page'. A shadow page contains 512 sptes, which can be either leaf or nonleaf sptes. A shadow page may contain a mix of leaf and nonleaf sptes.

A nonleaf spte allows the hardware mmu to reach the leaf pages and is not related to a translation directly. It points to other shadow pages.

A leaf spte corresponds to either one or two translations encoded into one paging structure entry. These are always the lowest level of the translation stack, with optional higher level translations left to NPT/EPT. Leaf ptes point at guest pages.

The following table shows translations encoded by leaf ptes, with higher-level translations in parentheses:

Non-nested guests:

nonpaging:	gpa->hpa
paging:	gva->gpa->hpa
paging, tdp:	(gva->)gpa->hpa

Nested guests:

```
non-tdp:      ngva->gpa->hpa  (*)
tdp:          (ngva->)ngpa->gpa->hpa
```

(\*) the guest hypervisor will encode the ngva->gpa\_↵  
↵translation into its page  
tables if npt is not present

### Shadow pages contain the following information:

#### **role.level:**

The level in the shadow paging hierarchy that this shadow page belongs to. 1=4k sptes, 2=2M sptes, 3=1G sptes, etc.

#### **role.direct:**

If set, leaf sptes reachable from this page are for a linear range. Examples include real mode translation, large guest pages backed by small host pages, and gpa->hpa translations when NPT or EPT is active. The linear range starts at (gfn << PAGE\_SHIFT) and its size is determined by role.level (2MB for first level, 1GB for second level, 0.5TB for third level, 256TB for fourth level) If clear, this page corresponds to a guest page table denoted by the gfn field.

#### **role.quadrant:**

When role.gpte\_is\_8\_bytes=0, the guest uses 32-bit gptes while the host uses 64-bit sptes. That means a guest page table contains more ptes than the host, so multiple shadow pages are needed to shadow one guest page. For first-level shadow pages, role.quadrant can be 0 or 1 and denotes the first or second 512-gpte block in the guest page table. For second-level page tables, each 32-bit gpte is converted to two 64-bit sptes (since each first-level guest page is shadowed by two first-level shadow pages) so role.quadrant takes values in the range 0..3. Each quadrant maps 1GB virtual address space.

#### **role.access:**

Inherited guest access permissions from the parent ptes in the form uwx. Note execute permission is positive, not negative.

#### **role.invalid:**

The page is invalid and should not be used. It is a root page that is currently pinned (by a cpu hardware register pointing to it); once it is unpinned it will be destroyed.

#### **role.gpte\_is\_8\_bytes:**

Reflects the size of the guest PTE for which the page is valid, i.e. '1' if 64-bit gptes are in use, '0' if 32-bit gptes are in use.

#### **role.nxe:**

Contains the value of efer.nxe for which the page is valid.

#### **role.cr0\_wp:**

Contains the value of cr0.wp for which the page is valid.

#### **role.smep\_andnot\_wp:**

Contains the value of cr4.smep && !cr0.wp for which the page is valid

(pages for which this is true are different from other pages; see the treatment of `cr0.wp=0` below).

**role.smap\_andnot\_wp:**

Contains the value of `cr4.smap && !cr0.wp` for which the page is valid (pages for which this is true are different from other pages; see the treatment of `cr0.wp=0` below).

**role.ept\_sp:**

This is a virtual flag to denote a shadowed nested EPT page. `ept_sp` is true if “`cr0_wp && smap_andnot_wp`”, an otherwise invalid combination.

**role.smm:**

Is 1 if the page is valid in system management mode. This field determines which of the `kvm_memslots` array was used to build this shadow page; it is also used to go back from a `struct kvm_mmu_page` to a memslot, through the `kvm_memslots_for_spte_role` macro and `__gfn_to_memslot`.

**role.ad\_disabled:**

Is 1 if the MMU instance cannot use A/D bits. EPT did not have A/D bits before Haswell; shadow EPT page tables also cannot use A/D bits if the L1 hypervisor does not enable them.

**gfn:**

Either the guest page table containing the translations shadowed by this page, or the base page frame for linear translations. See `role.direct`.

**spt:**

A pageful of 64-bit sptes containing the translations for this page. Accessed by both `kvm` and hardware. The page pointed to by `spt` will have its `page->private` pointing back at the shadow page structure. `sptes` in `spt` point either at guest pages, or at lower-level shadow pages. Specifically, if `sp1` and `sp2` are shadow pages, then `sp1->spt[n]` may point at `__pa(sp2->spt)`. `sp2` will point back at `sp1` through `parent_pte`. The `spt` array forms a DAG structure with the shadow page as a node, and guest pages as leaves.

**gfns:**

An array of 512 guest frame numbers, one for each present pte. Used to perform a reverse map from a pte to a gfn. When `role.direct` is set, any element of this array can be calculated from the `gfn` field when used, in this case, the array of `gfns` is not allocated. See `role.direct` and `gfn`.

**root\_count:**

A counter keeping track of how many hardware registers (guest `cr3` or `pdptrs`) are now pointing at the page. While this counter is nonzero, the page cannot be destroyed. See `role.invalid`.

**parent\_ptes:**

The reverse mapping for the pte/ptes pointing at this page's `spt`. If `parent_ptes` bit 0 is zero, only one spte points at this page and `parent_ptes` points at this single spte, otherwise, there exists multiple sptes pointing at this page and `(parent_ptes & ~0x1)` points at a data structure with a list of parent sptes.

**unsync:**

If true, then the translations in this page may not match the guest's translation. This is equivalent to the state of the tlb when a pte is changed but before the tlb entry is flushed. Accordingly, unsync ptes are synchronized when the guest executes `inlpg` or flushes its tlb by other means. Valid for leaf pages.

**unsync\_children:**

How many sptes in the page point at pages that are unsync (or have unsynchronized children).

**unsync\_child\_bitmap:**

A bitmap indicating which sptes in spt point (directly or indirectly) at pages that may be unsynchronized. Used to quickly locate all unsynchronized pages reachable from a given page.

**clear\_spte\_count:**

Only present on 32-bit hosts, where a 64-bit spte cannot be written atomically. The reader uses this while running out of the MMU lock to detect in-progress updates and retry them until the writer has finished the write.

**write\_flooding\_count:**

A guest may write to a page table many times, causing a lot of emulations if the page needs to be write-protected (see “Synchronized and unsynchronized pages” below). Leaf pages can be unsynchronized so that they do not trigger frequent emulation, but this is not possible for non-leaves. This field counts the number of emulations since the last time the page table was actually used; if emulation is triggered too frequently on this page, KVM will unmap the page to avoid emulation in the future.

### 1.7.7 Reverse map

The mmu maintains a reverse mapping whereby all ptes mapping a page can be reached given its gfn. This is used, for example, when swapping out a page.

### 1.7.8 Synchronized and unsynchronized pages

The guest uses two events to synchronize its tlb and page tables: tlb flushes and page invalidations (`inlpg`).

A tlb flush means that we need to synchronize all sptes reachable from the guest's `cr3`. This is expensive, so we keep all guest page tables write protected, and synchronize sptes to gptes when a gpte is written.

A special case is when a guest page table is reachable from the current guest `cr3`. In this case, the guest is obliged to issue an `inlpg` instruction before using the translation. We take advantage of that by removing write protection from the guest page, and allowing the guest to modify it freely. We synchronize modified gptes when the guest invokes `inlpg`. This reduces the amount of emulation we have to do when the guest modifies multiple gptes, or when the a guest page is no longer used as a page table and is used for random guest data.

As a side effect we have to resynchronize all reachable unsynchronized shadow pages on a tlb flush.

### 1.7.9 Reaction to events

- guest page fault (or npt page fault, or ept violation)

This is the most complicated event. The cause of a page fault can be:

- a true guest fault (the guest translation won't allow the access) (\*)
- access to a missing translation
- access to a protected translation - when logging dirty pages, memory is write protected - synchronized shadow pages are write protected (\*)
- access to untranslatable memory (mmio)

(\*) not applicable in direct mode

Handling a page fault is performed as follows:

- if the RSV bit of the error code is set, the page fault is caused by guest accessing MMIO and cached MMIO information is available.
  - walk shadow page table
  - check for valid generation number in the spte (see “Fast invalidation of MMIO sptes” below)
  - cache the information to `vcpu->arch.mmio_gva`, `vcpu->arch.mmio_access` and `vcpu->arch.mmio_gfn`, and call the emulator
- If both P bit and R/W bit of error code are set, this could possibly be handled as a “fast page fault” (fixed without taking the MMU lock). See the description in [KVM Lock Overview](#).
- if needed, walk the guest page tables to determine the guest translation (`gva->gpa` or `ngpa->gpa`)
  - if permissions are insufficient, reflect the fault back to the guest
- determine the host page
  - if this is an mmio request, there is no host page; cache the info to `vcpu->arch.mmio_gva`, `vcpu->arch.mmio_access` and `vcpu->arch.mmio_gfn`
- walk the shadow page table to find the spte for the translation, instantiating missing intermediate page tables as necessary
  - If this is an mmio request, cache the mmio info to the spte and set some reserved bit on the spte (see callers of `kvm_mmu_set_mmio_spte_mask`)
- try to unsynchronize the page
  - if successful, we can let the guest continue and modify the gpte
- emulate the instruction
  - if failed, unshadow the page and let the guest continue

- update any translations that were modified by the instruction

inlpg handling:

- walk the shadow page hierarchy and drop affected translations
- try to reinstantiate the indicated translation in the hope that the guest will use it in the near future

Guest control register updates:

- mov to cr3
  - look up new shadow roots
  - synchronize newly reachable shadow pages
- mov to cr0/cr4/efer
  - set up mmu context for new paging mode
  - look up new shadow roots
  - synchronize newly reachable shadow pages

Host translation updates:

- mmu notifier called with updated hva
- look up affected sptes through reverse map
- drop (or update) translations

### 1.7.10 Emulating cr0.wp

If tdp is not enabled, the host must keep cr0.wp=1 so page write protection works for the guest kernel, not guest user space. When the guest cr0.wp=1, this does not present a problem. However when the guest cr0.wp=0, we cannot map the permissions for gpte.u=1, gpte.w=0 to any spte (the semantics require allowing any guest kernel access plus user read access).

We handle this by mapping the permissions to two possible sptes, depending on fault type:

- kernel write fault: spte.u=0, spte.w=1 (allows full kernel access, disallows user access)
- read fault: spte.u=1, spte.w=0 (allows full read access, disallows kernel write access)

(user write faults generate a #PF)

In the first case there are two additional complications:

- if CR4.SMEP is enabled: since we've turned the page into a kernel page, the kernel may now execute it. We handle this by also setting spte.nx. If we get a user fetch or read fault, we'll change spte.u=1 and spte.nx=gpte.nx back. For this to work, KVM forces EFER.NX to 1 when shadow paging is in use.
- if CR4.SMAP is disabled: since the page has been changed to a kernel page, it can not be reused when CR4.SMAP is enabled. We set CR4.SMAP && !CR0.WP into shadow page's role to avoid this case. Note, here we do not

care the case that CR4.SMAP is enabled since KVM will directly inject #PF to guest due to failed permission check.

To prevent an spte that was converted into a kernel page with `cr0.wp=0` from being written by the kernel after `cr0.wp` has changed to 1, we make the value of `cr0.wp` part of the page role. This means that an spte created with one value of `cr0.wp` cannot be used when `cr0.wp` has a different value - it will simply be missed by the shadow page lookup code. A similar issue exists when an spte created with `cr0.wp=0` and `cr4.smep=0` is used after changing `cr4.smep` to 1. To avoid this, the value of `!cr0.wp && cr4.smep` is also made a part of the page role.

### 1.7.11 Large pages

The mmu supports all combinations of large and small guest and host pages. Supported page sizes include 4k, 2M, 4M, and 1G. 4M pages are treated as two separate 2M pages, on both guest and host, since the mmu always uses PAE paging.

To instantiate a large spte, four constraints must be satisfied:

- the spte must point to a large host page
- the guest pte must be a large pte of at least equivalent size (if tdp is enabled, there is no guest pte and this condition is satisfied)
- if the spte will be writeable, the large page frame may not overlap any write-protected pages
- the guest page must be wholly contained by a single memory slot

To check the last two conditions, the mmu maintains a `->disallow_lpage` set of arrays for each memory slot and large page size. Every write protected page causes its `disallow_lpage` to be incremented, thus preventing instantiation of a large spte. The frames at the end of an unaligned memory slot have artificially inflated `->disallow_lpages` so they can never be instantiated.

### 1.7.12 Fast invalidation of MMIO sptes

As mentioned in “Reaction to events” above, kvm will cache MMIO information in leaf sptes. When a new memslot is added or an existing memslot is changed, this information may become stale and needs to be invalidated. This also needs to hold the MMU lock while walking all shadow pages, and is made more scalable with a similar technique.

MMIO sptes have a few spare bits, which are used to store a generation number. The global generation number is stored in `kvm_memslots(kvm)->generation`, and increased whenever guest memory info changes.

When KVM finds an MMIO spte, it checks the generation number of the spte. If the generation number of the spte does not equal the global generation number, it will ignore the cached MMIO information and handle the page fault through the slow path.

Since only 18 bits are used to store generation-number on mmio spte, all pages are zapped when there is an overflow.

Unfortunately, a single memory access might access `kvm_memslots(kvm)` multiple times, the last one happening when the generation number is retrieved and stored into the MMIO spte. Thus, the MMIO spte might be created based on out-of-date information, but with an up-to-date generation number.

To avoid this, the generation number is incremented again after `synchronize_srcu` returns; thus, bit 63 of `kvm_memslots(kvm)->generation` set to 1 only during a memslot update, while some SRCU readers might be using the old copy. We do not want to use an MMIO sptes created with an odd generation number, and we can do this without losing a bit in the MMIO spte. The “update in-progress” bit of the generation is not stored in MMIO spte, and is so is implicitly zero when the generation is extracted out of the spte. If KVM is unlucky and creates an MMIO spte while an update is in-progress, the next access to the spte will always be a cache miss. For example, a subsequent access during the update window will miss due to the in-progress flag diverging, while an access after the update window closes will have a higher generation number (as compared to the spte).

### 1.7.13 Further reading

- NPT presentation from KVM Forum 2008 [https://www.linux-kvm.org/images/c/c8/KvmForum2008%24kdf2008\\_21.pdf](https://www.linux-kvm.org/images/c/c8/KvmForum2008%24kdf2008_21.pdf)

## 1.8 KVM-specific MSRs

### Author

Glauber Costa <[glommer@redhat.com](mailto:glommer@redhat.com)>, Red Hat Inc, 2010

KVM makes use of some custom MSRs to service some requests.

Custom MSRs have a range reserved for them, that goes from `0x4b564d00` to `0x4b564dff`. There are MSRs outside this area, but they are deprecated and their use is discouraged.

### 1.8.1 Custom MSR list

The current supported Custom MSR list is:

#### **MSR\_KVM\_WALL\_CLOCK\_NEW:**

`0x4b564d00`

#### **data:**

4-byte alignment physical address of a memory area which must be in guest RAM. This memory is expected to hold a copy of the following structure:

```
struct pvclock_wall_clock {
    u32    version;
    u32    sec;
    u32    nsec;
} __attribute__((__packed__));
```



whose data will be filled in by the hypervisor. The hypervisor is only guaranteed to update this data at the moment of MSR write. Users that want to reliably query this information more than once have to write more than once to this MSR. Fields have the following meanings:

**version:**

guest has to check version before and after grabbing time information and check that they are both equal and even. An odd version indicates an in-progress update.

**sec:**

number of seconds for wallclock at time of boot.

**nsec:**

number of nanoseconds for wallclock at time of boot.

In order to get the current wallclock time, the `system_time` from `MSR_KVM_SYSTEM_TIME_NEW` needs to be added.

Note that although MSRs are per-CPU entities, the effect of this particular MSR is global.

Availability of this MSR must be checked via bit 3 in `0x4000001` cpuid leaf prior to usage.

**MSR\_KVM\_SYSTEM\_TIME\_NEW:**

`0x4b564d01`

**data:**

4-byte aligned physical address of a memory area which must be in guest RAM, plus an enable bit in bit 0. This memory is expected to hold a copy of the following structure:

```
struct pvclock_vcpu_time_info {
    u32    version;
    u32    pad0;
    u64    tsc_timestamp;
    u64    system_time;
    u32    tsc_to_system_mul;
    s8     tsc_shift;
    u8     flags;
    u8     pad[2];
} __attribute__((__packed__)); /* 32 bytes */
```

whose data will be filled in by the hypervisor periodically. Only one write, or registration, is needed for each VCPU. The interval between updates of this structure is arbitrary and implementation-dependent. The hypervisor may update this structure at any time it sees fit until anything with `bit0 == 0` is written to it.

Fields have the following meanings:

**version:**

guest has to check version before and after grabbing time information and check that they are both equal and even. An odd version indicates an in-progress update.

**tsc\_timestamp:**

the tsc value at the current VCPU at the time of the update of this structure. Guests can subtract this value from current tsc to derive a notion of elapsed time since the structure update.

**system\_time:**

a host notion of monotonic time, including sleep time at the time this structure was last updated. Unit is nanoseconds.

**tsc\_to\_system\_mul:**

multiplier to be used when converting tsc-related quantity to nanoseconds

**tsc\_shift:**

shift to be used when converting tsc-related quantity to nanoseconds. This shift will ensure that multiplication with tsc\_to\_system\_mul does not overflow. A positive value denotes a left shift, a negative value a right shift.

The conversion from tsc to nanoseconds involves an additional right shift by 32 bits. With this information, guests can derive per-CPU time by doing:

```
time = (current_tsc - tsc_timestamp)
if (tsc_shift >= 0)
    time <=< tsc_shift;
else
    time >>= -tsc_shift;
time = (time * tsc_to_system_mul) >> 32
time = time + system_time
```

**flags:**

bits in this field indicate extended capabilities coordinated between the guest and the hypervisor. Availability of specific flags has to be checked in 0x40000001 cpuid leaf. Current flags are:

flag bit	cpuid bit	meaning
0	24	time measures taken across multiple cpus are guaranteed to be monotonic
1	N/A	guest vcpu has been paused by the host See 4.70 in api.txt

Availability of this MSR must be checked via bit 3 in 0x40000001 cpuid leaf prior to usage.

**MSR\_KVM\_WALL\_CLOCK:**

0x11

**data and functioning:**

same as MSR\_KVM\_WALL\_CLOCK\_NEW. Use that instead.

This MSR falls outside the reserved KVM range and may be removed in the future. Its usage is deprecated.

Availability of this MSR must be checked via bit 0 in 0x4000001 cpuid leaf prior to usage.

### **MSR\_KVM\_SYSTEM\_TIME:**

0x12

#### **data and functioning:**

same as MSR\_KVM\_SYSTEM\_TIME\_NEW. Use that instead.

This MSR falls outside the reserved KVM range and may be removed in the future. Its usage is deprecated.

Availability of this MSR must be checked via bit 0 in 0x4000001 cpuid leaf prior to usage.

The suggested algorithm for detecting kvmclock presence is then:

```
if (!kvm_para_available())    /* refer to cpuid.txt */
    return NON_PRESENT;

flags = cpuid_eax(0x40000001);
if (flags & 3) {
    msr_kvm_system_time = MSR_KVM_SYSTEM_TIME_NEW;
    msr_kvm_wall_clock = MSR_KVM_WALL_CLOCK_NEW;
    return PRESENT;
} else if (flags & 0) {
    msr_kvm_system_time = MSR_KVM_SYSTEM_TIME;
    msr_kvm_wall_clock = MSR_KVM_WALL_CLOCK;
    return PRESENT;
} else
    return NON_PRESENT;
```

### **MSR\_KVM\_ASYNC\_PF\_EN:**

0x4b564d02

#### **data:**

Asynchronous page fault (APF) control MSR.

Bits 63-6 hold 64-byte aligned physical address of a 64 byte memory area which must be in guest RAM and must be zeroed. This memory is expected to hold a copy of the following structure:

```
struct kvm_vcpu_pv_apf_data {
    /* Used for 'page not present' events delivered via #PF */
    __u32 flags;

    /* Used for 'page ready' events delivered via interrupt_
    ↪notification */
    __u32 token;

    __u8 pad[56];
    __u32 enabled;
};
```

Bits 5-4 of the MSR are reserved and should be zero. Bit 0 is set to 1 when

asynchronous page faults are enabled on the vcpu, 0 when disabled. Bit 1 is 1 if asynchronous page faults can be injected when vcpu is in cpl == 0. Bit 2 is 1 if asynchronous page faults are delivered to L1 as #PF vmexits. Bit 2 can be set only if KVM\_FEATURE\_ASYNC\_PF\_VMEXIT is present in CPUID. Bit 3 enables interrupt based delivery of 'page ready' events. Bit 3 can only be set if KVM\_FEATURE\_ASYNC\_PF\_INT is present in CPUID.

'Page not present' events are currently always delivered as synthetic #PF exception. During delivery of these events APF CR2 register contains a token that will be used to notify the guest when missing page becomes available. Also, to make it possible to distinguish between real #PF and APF, first 4 bytes of 64 byte memory location ( 'flags' ) will be written to by the hypervisor at the time of injection. Only first bit of 'flags' is currently supported, when set, it indicates that the guest is dealing with asynchronous 'page not present' event. If during a page fault APF 'flags' is '0' it means that this is regular page fault. Guest is supposed to clear 'flags' when it is done handling #PF exception so the next event can be delivered.

Note, since APF 'page not present' events use the same exception vector as regular page fault, guest must reset 'flags' to '0' before it does something that can generate normal page fault.

Bytes 5-7 of 64 byte memory location ( 'token' ) will be written to by the hypervisor at the time of APF 'page ready' event injection. The content of these bytes is a token which was previously delivered as 'page not present' event. The event indicates the page is now available. Guest is supposed to write '0' to 'token' when it is done handling 'page ready' event and to write 1 to MSR\_KVM\_ASYNC\_PF\_ACK after clearing the location; writing to the MSR forces KVM to re-scan its queue and deliver the next pending notification.

Note, MSR\_KVM\_ASYNC\_PF\_INT MSR specifying the interrupt vector for 'page ready' APF delivery needs to be written to before enabling APF mechanism in MSR\_KVM\_ASYNC\_PF\_EN or interrupt #0 can get injected. The MSR is available if KVM\_FEATURE\_ASYNC\_PF\_INT is present in CPUID.

Note, previously, 'page ready' events were delivered via the same #PF exception as 'page not present' events but this is now deprecated. If bit 3 (interrupt based delivery) is not set APF events are not delivered.

If APF is disabled while there are outstanding APFs, they will not be delivered.

Currently 'page ready' APF events will be always delivered on the same vcpu as 'page not present' event was, but guest should not rely on that.

### **MSR\_KVM\_STEAL\_TIME:**

0x4b564d03

#### **data:**

64-byte alignment physical address of a memory area which must be in guest RAM, plus an enable bit in bit 0. This memory is expected to hold a copy of the following structure:

```
struct kvm_steal_time {
    __u64 steal;
    __u32 version;
```

(continues on next page)

(continued from previous page)

```

    __u32 flags;
    __u8  preempted;
    __u8  u8_pad[3];
    __u32 pad[11];
}

```

whose data will be filled in by the hypervisor periodically. Only one write, or registration, is needed for each VCPU. The interval between updates of this structure is arbitrary and implementation-dependent. The hypervisor may update this structure at any time it sees fit until anything with bit0 == 0 is written to it. Guest is required to make sure this structure is initialized to zero.

Fields have the following meanings:

**version:**

a sequence counter. In other words, guest has to check this field before and after grabbing time information and make sure they are both equal and even. An odd version indicates an in-progress update.

**flags:**

At this point, always zero. May be used to indicate changes in this structure in the future.

**steal:**

the amount of time in which this vCPU did not run, in nanoseconds. Time during which the vcpu is idle, will not be reported as steal time.

**preempted:**

indicate the vCPU who owns this struct is running or not. Non-zero values mean the vCPU has been preempted. Zero means the vCPU is not preempted. NOTE, it is always zero if the the hypervisor doesn't support this field.

**MSR\_KVM\_EOI\_EN:**

0x4b564d04

**data:**

Bit 0 is 1 when PV end of interrupt is enabled on the vcpu; 0 when disabled. Bit 1 is reserved and must be zero. When PV end of interrupt is enabled (bit 0 set), bits 63-2 hold a 4-byte aligned physical address of a 4 byte memory area which must be in guest RAM and must be zeroed.

The first, least significant bit of 4 byte memory location will be written to by the hypervisor, typically at the time of interrupt injection. Value of 1 means that guest can skip writing EOI to the apic (using MSR or MMIO write); instead, it is sufficient to signal EOI by clearing the bit in guest memory - this location will later be polled by the hypervisor. Value of 0 means that the EOI write is required.

It is always safe for the guest to ignore the optimization and perform the APIC EOI write anyway.

Hypervisor is guaranteed to only modify this least significant bit while in the current VCPU context, this means that guest does not need to use either lock

prefix or memory ordering primitives to synchronise with the hypervisor.

However, hypervisor can set and clear this memory bit at any time: therefore to make sure hypervisor does not interrupt the guest and clear the least significant bit in the memory area in the window between guest testing it to detect whether it can skip EOI apic write and between guest clearing it to signal EOI to the hypervisor, guest must both read the least significant bit in the memory area and clear it using a single CPU instruction, such as test and clear, or compare and exchange.

**MSR\_KVM\_POLL\_CONTROL:**

0x4b564d05

Control host-side polling.

**data:**

Bit 0 enables (1) or disables (0) host-side HLT polling logic.

KVM guests can request the host not to poll on HLT, for example if they are performing polling themselves.

**MSR\_KVM\_ASYNC\_PF\_INT:**

0x4b564d06

**data:**

Second asynchronous page fault (APF) control MSR.

Bits 0-7: APIC vector for delivery of ‘page ready’ APF events. Bits 8-63: Reserved

Interrupt vector for asynchronous ‘page ready’ notifications delivery. The vector has to be set up before asynchronous page fault mechanism is enabled in MSR\_KVM\_ASYNC\_PF\_EN. The MSR is only available if KVM\_FEATURE\_ASYNC\_PF\_INT is present in CPUID.

**MSR\_KVM\_ASYNC\_PF\_ACK:**

0x4b564d07

**data:**

Asynchronous page fault (APF) acknowledgment.

When the guest is done processing ‘page ready’ APF event and ‘token’ field in ‘struct kvm\_vcpu\_pv\_apf\_data’ is cleared it is supposed to write ‘1’ to bit 0 of the MSR, this causes the host to re-scan its queue and check if there are more notifications pending. The MSR is available if KVM\_FEATURE\_ASYNC\_PF\_INT is present in CPUID.

## 1.9 Nested VMX

### 1.9.1 Overview

On Intel processors, KVM uses Intel’s VMX (Virtual-Machine eXtensions) to easily and efficiently run guest operating systems. Normally, these guests *cannot* themselves be hypervisors running their own guests, because in VMX, guests cannot use VMX instructions.

The “Nested VMX” feature adds this missing capability - of running guest hypervisors (which use VMX) with their own nested guests. It does so by allowing a guest to use VMX instructions, and correctly and efficiently emulating them using the single level of VMX available in the hardware.

We describe in much greater detail the theory behind the nested VMX feature, its implementation and its performance characteristics, in the OSDI 2010 paper “The Turtles Project: Design and Implementation of Nested Virtualization” , available at:

[https://www.usenix.org/events/osdi10/tech/full\\_papers/Ben-Yehuda.pdf](https://www.usenix.org/events/osdi10/tech/full_papers/Ben-Yehuda.pdf)

### 1.9.2 Terminology

Single-level virtualization has two levels - the host (KVM) and the guests. In nested virtualization, we have three levels: The host (KVM), which we call L0, the guest hypervisor, which we call L1, and its nested guest, which we call L2.

### 1.9.3 Running nested VMX

The nested VMX feature is disabled by default. It can be enabled by giving the “nested=1” option to the kvm-intel module.

No modifications are required to user space (qemu). However, qemu’s default emulated CPU type (qemu64) does not list the “VMX” CPU feature, so it must be explicitly enabled, by giving qemu one of the following options:

- `cpu host` (emulated CPU has all features of the real CPU)
- `cpu qemu64,+vmx` (add just the vmx feature to a named CPU type)

### 1.9.4 ABIs

Nested VMX aims to present a standard and (eventually) fully-functional VMX implementation for the a guest hypervisor to use. As such, the official specification of the ABI that it provides is Intel’s VMX specification, namely volume 3B of their “Intel 64 and IA-32 Architectures Software Developer’s Manual” . Not all of VMX’s features are currently fully supported, but the goal is to eventually support them all, starting with the VMX features which are used in practice by popular hypervisors (KVM and others).

As a VMX implementation, nested VMX presents a VMCS structure to L1. As mandated by the spec, other than the two fields `revision_id` and `abort`, this structure is *opaque* to its user, who is not supposed to know or care about its internal structure. Rather, the structure is accessed through the `VMREAD` and `VMWRITE` instructions. Still, for debugging purposes, KVM developers might be interested to know the internals of this structure; This is `struct vmcs12` from `arch/x86/kvm/vmx.c`.

The name “`vmcs12`” refers to the VMCS that L1 builds for L2. In the code we also have “`vmcs01`” , the VMCS that L0 built for L1, and “`vmcs02`” is the VMCS which L0 builds to actually run L2 - how this is done is explained in the aforementioned paper.

For convenience, we repeat the content of struct vmcs12 here. If the internals of this structure changes, this can break live migration across KVM versions. VMCS12\_REVISION (from vmx.c) should be changed if struct vmcs12 or its inner struct shadow\_vmcs is ever changed.

```
typedef u64 natural_width;
struct __packed vmcs12 {
    /* According to the Intel spec, a VMCS region must start
    ↪with
        * these two user-visible fields */
    u32 revision_id;
    u32 abort;

    u32 launch_state; /* set to 0 by VMCLEAR, to 1 by VMLAUNCH,
    ↪*/
    u32 padding[7]; /* room for future expansion */

    u64 io_bitmap_a;
    u64 io_bitmap_b;
    u64 msr_bitmap;
    u64 vm_exit_msr_store_addr;
    u64 vm_exit_msr_load_addr;
    u64 vm_entry_msr_load_addr;
    u64 tsc_offset;
    u64 virtual_apic_page_addr;
    u64 apic_access_addr;
    u64 ept_pointer;
    u64 guest_physical_address;
    u64 vmcs_link_pointer;
    u64 guest_ia32_debugctl;
    u64 guest_ia32_pat;
    u64 guest_ia32_efer;
    u64 guest_pdptr0;
    u64 guest_pdptr1;
    u64 guest_pdptr2;
    u64 guest_pdptr3;
    u64 host_ia32_pat;
    u64 host_ia32_efer;
    u64 padding64[8]; /* room for future expansion */
    natural_width cr0_guest_host_mask;
    natural_width cr4_guest_host_mask;
    natural_width cr0_read_shadow;
    natural_width cr4_read_shadow;
    natural_width dead_space[4]; /* Last remnants of cr3_target_
    ↪value[0-3]. */
    natural_width exit_qualification;
    natural_width guest_linear_address;
    natural_width guest_cr0;
    natural_width guest_cr3;
    natural_width guest_cr4;
    natural_width guest_es_base;
```

(continues on next page)



(continued from previous page)

```
natural_width guest_cs_base;
natural_width guest_ss_base;
natural_width guest_ds_base;
natural_width guest_fs_base;
natural_width guest_gs_base;
natural_width guest_ldtr_base;
natural_width guest_tr_base;
natural_width guest_gdtr_base;
natural_width guest_idtr_base;
natural_width guest_dr7;
natural_width guest_rsp;
natural_width guest_rip;
natural_width guest_rflags;
natural_width guest_pending_dbg_exceptions;
natural_width guest_sysenter_esp;
natural_width guest_sysenter_eip;
natural_width host_cr0;
natural_width host_cr3;
natural_width host_cr4;
natural_width host_fs_base;
natural_width host_gs_base;
natural_width host_tr_base;
natural_width host_gdtr_base;
natural_width host_idtr_base;
natural_width host_ia32_sysenter_esp;
natural_width host_ia32_sysenter_eip;
natural_width host_rsp;
natural_width host_rip;
natural_width paddingl[8]; /* room for future expansion */
u32 pin_based_vm_exec_control;
u32 cpu_based_vm_exec_control;
u32 exception_bitmap;
u32 page_fault_error_code_mask;
u32 page_fault_error_code_match;
u32 cr3_target_count;
u32 vm_exit_controls;
u32 vm_exit_msr_store_count;
u32 vm_exit_msr_load_count;
u32 vm_entry_controls;
u32 vm_entry_msr_load_count;
u32 vm_entry_intr_info_field;
u32 vm_entry_exception_error_code;
u32 vm_entry_instruction_len;
u32 tpr_threshold;
u32 secondary_vm_exec_control;
u32 vm_instruction_error;
u32 vm_exit_reason;
u32 vm_exit_intr_info;
u32 vm_exit_intr_error_code;
```

(continues on next page)

(continued from previous page)

```
u32 idt_vectoring_info_field;
u32 idt_vectoring_error_code;
u32 vm_exit_instruction_len;
u32 vmx_instruction_info;
u32 guest_es_limit;
u32 guest_cs_limit;
u32 guest_ss_limit;
u32 guest_ds_limit;
u32 guest_fs_limit;
u32 guest_gs_limit;
u32 guest_ldtr_limit;
u32 guest_tr_limit;
u32 guest_gdtr_limit;
u32 guest_idtr_limit;
u32 guest_es_ar_bytes;
u32 guest_cs_ar_bytes;
u32 guest_ss_ar_bytes;
u32 guest_ds_ar_bytes;
u32 guest_fs_ar_bytes;
u32 guest_gs_ar_bytes;
u32 guest_ldtr_ar_bytes;
u32 guest_tr_ar_bytes;
u32 guest_interruptibility_info;
u32 guest_activity_state;
u32 guest_sysenter_cs;
u32 host_ia32_sysenter_cs;
u32 padding32[8]; /* room for future expansion */
u16 virtual_processor_id;
u16 guest_es_selector;
u16 guest_cs_selector;
u16 guest_ss_selector;
u16 guest_ds_selector;
u16 guest_fs_selector;
u16 guest_gs_selector;
u16 guest_ldtr_selector;
u16 guest_tr_selector;
u16 host_es_selector;
u16 host_cs_selector;
u16 host_ss_selector;
u16 host_ds_selector;
u16 host_fs_selector;
u16 host_gs_selector;
u16 host_tr_selector;
};
```

### 1.9.5 Authors

**These patches were written by:**

- Abel Gordon, abelg <at> il.ibm.com
- Nadav Har' El, nyh <at> il.ibm.com
- Orit Wasserman, oritw <at> il.ibm.com
- Ben-Ami Yassor, benami <at> il.ibm.com
- Muli Ben-Yehuda, muli <at> il.ibm.com

**With contributions by:**

- Anthony Liguori, aliguori <at> us.ibm.com
- Mike Day, mdday <at> us.ibm.com
- Michael Factor, factor <at> il.ibm.com
- Zvi Dubitzky, dubi <at> il.ibm.com

**And valuable reviews by:**

- Avi Kivity, avi <at> redhat.com
- Gleb Natapov, gleb <at> redhat.com
- Marcelo Tosatti, mtosatti <at> redhat.com
- Kevin Tian, kevin.tian <at> intel.com
- and others.

## 1.10 The PPC KVM paravirtual interface

The basic execution principle by which KVM on PowerPC works is to run all kernel space code in PR=1 which is user space. This way we trap all privileged instructions and can emulate them accordingly.

Unfortunately that is also the downfall. There are quite some privileged instructions that needlessly return us to the hypervisor even though they could be handled differently.

This is what the PPC PV interface helps with. It takes privileged instructions and transforms them into unprivileged ones with some help from the hypervisor. This cuts down virtualization costs by about 50% on some of my benchmarks.

The code for that interface can be found in arch/powerpc/kernel/kvm\*

### 1.10.1 Querying for existence

To find out if we're running on KVM or not, we leverage the device tree. When Linux is running on KVM, a node /hypervisor exists. That node contains a compatible property with the value "linux,kvm".

Once you determined you're running under a PV capable KVM, you can now use hypercalls as described below.

### 1.10.2 KVM hypercalls

Inside the device tree's /hypervisor node there's a property called 'hypercall-instructions'. This property contains at most 4 opcodes that make up the hypercall. To call a hypercall, just call these instructions.

The parameters are as follows:

Register	IN	OUT
r0	.	volatile
r3	1st parameter	Return code
r4	2nd parameter	1st output value
r5	3rd parameter	2nd output value
r6	4th parameter	3rd output value
r7	5th parameter	4th output value
r8	6th parameter	5th output value
r9	7th parameter	6th output value
r10	8th parameter	7th output value
r11	hypercall number	8th output value
r12	.	volatile

Hypercall definitions are shared in generic code, so the same hypercall numbers apply for x86 and powerpc alike with the exception that each KVM hypercall also needs to be ORed with the KVM vendor code which is  $(42 \ll 16)$ .

Return codes can be as follows:

Code	Meaning
0	Success
12	Hypercall not implemented
<0	Error

### 1.10.3 The magic page

To enable communication between the hypervisor and guest there is a new shared page that contains parts of supervisor visible register state. The guest can map this shared page using the KVM hypercall `KVM_HC_PPC_MAP_MAGIC_PAGE`.

With this hypercall issued the guest always gets the magic page mapped at the desired location. The first parameter indicates the effective address when the MMU is enabled. The second parameter indicates the address in real mode, if applicable to the target. For now, we always map the page to -4096. This way we can access it using absolute load and store functions. The following instruction reads the first field of the magic page:

```
ld      rX, -4096(0)
```

The interface is designed to be extensible should there be need later to add additional registers to the magic page. If you add fields to the magic page, also define a new hypercall feature to indicate that the host can give you more registers. Only if the host supports the additional features, make use of them.

The magic page layout is described by struct `kvm_vcpu_arch_shared` in `arch/powerpc/include/asm/kvm_para.h`.

### 1.10.4 Magic page features

When mapping the magic page using the KVM hypercall `KVM_HC_PPC_MAP_MAGIC_PAGE`, a second return value is passed to the guest. This second return value contains a bitmap of available features inside the magic page.

The following enhancements to the magic page are currently available:

<code>KVM_MAGIC_FEAT_SR</code>	Maps SR registers r/w in the magic page
<code>KVM_MAGIC_FEAT_MAS0_TO_SF</code>	Maps MASn, ESR, PIR and high SPRGs

For enhanced features in the magic page, please check for the existence of the feature before using them!

### 1.10.5 Magic page flags

In addition to features that indicate whether a host is capable of a particular feature we also have a channel for a guest to tell the guest whether it's capable of something. This is what we call "flags".

Flags are passed to the host in the low 12 bits of the Effective Address.

The following flags are currently available for a guest to expose:

MAGIC\_PAGE\_FLAG\_NOT\_MAPPED\_NX Guest handles NX bits correctly wrt magic page

### 1.10.6 MSR bits

The MSR contains bits that require hypervisor intervention and bits that do not require direct hypervisor intervention because they only get interpreted when entering the guest or don't have any impact on the hypervisor's behavior.

The following bits are safe to be set inside the guest:

- MSR\_EE
- MSR\_RI

If any other bit changes in the MSR, please still use `mtmsr(d)`.

### 1.10.7 Patched instructions

The “ld” and “std” instructions are transformed to “lwz” and “stw” instructions respectively on 32 bit systems with an added offset of 4 to accommodate for big endianness.

The following is a list of mapping the Linux kernel performs when running as guest. Implementing any of those mappings is optional, as the instruction traps also act on the shared page. So calling privileged instructions still works as before.

From	To
mfmsr rX	ld rX, magic_page->msr
mfsprg rX, 0	ld rX, magic_page->sprg0
mfsprg rX, 1	ld rX, magic_page->sprg1
mfsprg rX, 2	ld rX, magic_page->sprg2
mfsprg rX, 3	ld rX, magic_page->sprg3
mfrr0 rX	ld rX, magic_page->srr0
mfrr1 rX	ld rX, magic_page->srr1
mfdr rX	ld rX, magic_page->dar
mfdsisr rX	lwz rX, magic_page->dsisr
mtmsr rX	std rX, magic_page->msr
mtsprg 0, rX	std rX, magic_page->sprg0
mtsprg 1, rX	std rX, magic_page->sprg1
mtsprg 2, rX	std rX, magic_page->sprg2
mtsprg 3, rX	std rX, magic_page->sprg3
mtsrr0 rX	std rX, magic_page->srr0
mtsrr1 rX	std rX, magic_page->srr1
mtdr rX	std rX, magic_page->dar
mtdsisr rX	stw rX, magic_page->dsisr
tlbsync	nop
mtmsrd rX, 0	b <special mtmsr section>
mtmsr rX	b <special mtmsr section>
mtmsrd rX, 1	b <special mtmsrd section>
[Book3S only]	
mtsrin rX, rY	b <special mtsrin section>
[BookE only]	
wrteei [0 1]	b <special wrteei section>

Some instructions require more logic to determine what's going on than a load or store instruction can deliver. To enable patching of those, we keep some RAM around where we can live translate instructions to. What happens is the following:

- 1) copy emulation code to memory
- 2) patch that code to fit the emulated instruction
- 3) patch that code to return to the original pc + 4
- 4) patch the original instruction to branch to the new code

That way we can inject an arbitrary amount of code as replacement for a single instruction. This allows us to check for pending interrupts when setting EE=1 for example.

### 1.10.8 Hypercall ABIs in KVM on PowerPC

#### 1) KVM hypercalls (ePAPR)

These are ePAPR compliant hypercall implementation (mentioned above). Even generic hypercalls are implemented here, like the ePAPR idle hcall. These are available on all targets.

#### 2) PAPR hypercalls

PAPR hypercalls are needed to run server PowerPC PAPR guests (-M pseries in QEMU). These are the same hypercalls that pHyp, the POWER hypervisor implements. Some of them are handled in the kernel, some are handled in user space. This is only available on book3s\_64.

#### 3) OSI hypercalls

Mac-on-Linux is another user of KVM on PowerPC, which has its own hypercall (long before KVM). This is supported to maintain compatibility. All these hypercalls get forwarded to user space. This is only useful on book3s\_32, but can be used with book3s\_64 as well.

## 1.11 The s390 DIAGNOSE call on KVM

KVM on s390 supports the DIAGNOSE call for making hypercalls, both for native hypercalls and for selected hypercalls found on other s390 hypervisors.

Note that bits are numbered as by the usual s390 convention (most significant bit on the left).

### 1.11.1 General remarks

DIAGNOSE calls by the guest cause a mandatory intercept. This implies all supported DIAGNOSE calls need to be handled by either KVM or its userspace.

All DIAGNOSE calls supported by KVM use the RS-a format:

-----						
	'83'		R1		R3	
			B2		D2	
-----						
0	8	12	16	20		31

The second-operand address (obtained by the base/displacement calculation) is not used to address data. Instead, bits 48-63 of this address specify the function code, and bits 0-47 are ignored.

The supported DIAGNOSE function codes vary by the userspace used. For DIAGNOSE function codes not specific to KVM, please refer to the documentation for the s390 hypervisors defining them.



### 1.11.2 DIAGNOSE function code ‘X’ 500’ - KVM virtio functions

If the function code specifies 0x500, various virtio-related functions are performed.

General register 1 contains the virtio subfunction code. Supported virtio subfunctions depend on KVM’s userspace. Generally, userspace provides either s390-virtio (subcodes 0-2) or virtio-ccw (subcode 3).

Upon completion of the DIAGNOSE instruction, general register 2 contains the function’s return code, which is either a return code or a subcode specific value.

#### **Subcode 0 - s390-virtio notification and early console printk**

Handled by userspace.

#### **Subcode 1 - s390-virtio reset**

Handled by userspace.

#### **Subcode 2 - s390-virtio set status**

Handled by userspace.

#### **Subcode 3 - virtio-ccw notification**

Handled by either userspace or KVM (ioeventfd case).

General register 2 contains a subchannel-identification word denoting the subchannel of the virtio-ccw proxy device to be notified.

General register 3 contains the number of the virtqueue to be notified.

General register 4 contains a 64bit identifier for KVM usage (the `kvm_io_bus` cookie). If general register 4 does not contain a valid identifier, it is ignored.

After completion of the DIAGNOSE call, general register 2 may contain a 64bit identifier (in the `kvm_io_bus` cookie case), or a negative error value, if an internal error occurred.

See also the virtio standard for a discussion of this hypercall.

### 1.11.3 DIAGNOSE function code ‘X’ 501 - KVM breakpoint

If the function code specifies 0x501, breakpoint functions may be performed. This function code is handled by userspace.

This diagnose function code has no subfunctions and uses no parameters.

## 1.12 s390 (IBM Z) Ultravisor and Protected VMs

### 1.12.1 Summary

Protected virtual machines (PVM) are KVM VMs that do not allow KVM to access VM state like guest memory or guest registers. Instead, the PVMs are mostly managed by a new entity called Ultravisor (UV). The UV provides an API that can be used by PVMs and KVM to request management actions.

Each guest starts in non-protected mode and then may make a request to transition into protected mode. On transition, KVM registers the guest and its VCPUs with the Ultravisor and prepares everything for running it.

The Ultravisor will secure and decrypt the guest's boot memory (i.e. kernel/initrd). It will safeguard state changes like VCPU starts/stops and injected interrupts while the guest is running.

As access to the guest's state, such as the SIE state description, is normally needed to be able to run a VM, some changes have been made in the behavior of the SIE instruction. A new format 4 state description has been introduced, where some fields have different meanings for a PVM. SIE exits are minimized as much as possible to improve speed and reduce exposed guest state.

### 1.12.2 Interrupt injection

Interrupt injection is safeguarded by the Ultravisor. As KVM doesn't have access to the VCPUs' lowcores, injection is handled via the format 4 state description.

Machine check, external, IO and restart interruptions each can be injected on SIE entry via a bit in the interrupt injection control field (offset 0x54). If the guest cpu is not enabled for the interrupt at the time of injection, a validity interception is recognized. The format 4 state description contains fields in the interception data block where data associated with the interrupt can be transported.

Program and Service Call exceptions have another layer of safeguarding; they can only be injected for instructions that have been intercepted into KVM. The exceptions need to be a valid outcome of an instruction emulation by KVM, e.g. we can never inject a addressing exception as they are reported by SIE since KVM has no access to the guest memory.

### 1.12.3 Mask notification interceptions

KVM cannot intercept lctl(g) and lpsw(e) anymore in order to be notified when a PVM enables a certain class of interrupt. As a replacement, two new interception codes have been introduced: One indicating that the contents of CRs 0, 6, or 14 have been changed, indicating different interruption subclasses; and one indicating that PSW bit 13 has been changed, indicating that a machine check intervention was requested and those are now enabled.

### 1.12.4 Instruction emulation

With the format 4 state description for PVMs, the SIE instruction already interprets more instructions than it does with format 2. It is not able to interpret every instruction, but needs to hand some tasks to KVM; therefore, the SIE and the ultravisor safeguard emulation inputs and outputs.

The control structures associated with SIE provide the Secure Instruction Data Area (SIDA), the Interception Parameters (IP) and the Secure Interception General Register Save Area. Guest GRs and most of the instruction data, such as I/O data structures, are filtered. Instruction data is copied to and from the SIDA when needed. Guest GRs are put into / retrieved from the Secure Interception General Register Save Area.

Only GR values needed to emulate an instruction will be copied into this save area and the real register numbers will be hidden.

The Interception Parameters state description field still contains the bytes of the instruction text, but with pre-set register values instead of the actual ones. I.e. each instruction always uses the same instruction text, in order not to leak guest instruction text. This also implies that the register content that a guest had in `r<n>` may be in `r<m>` from the hypervisor's point of view.

The Secure Instruction Data Area contains instruction storage data. Instruction data, i.e. data being referenced by an instruction like the SCCB for `scbp`, is moved via the SIDA. When an instruction is intercepted, the SIE will only allow data and program interrupts for this instruction to be moved to the guest via the two data areas discussed before. Other data is either ignored or results in validity interceptions.

### 1.12.5 Instruction emulation interceptions

There are two types of SIE secure instruction intercepts: the normal and the notification type. Normal secure instruction intercepts will make the guest pending for instruction completion of the intercepted instruction type, i.e. on SIE entry it is attempted to complete emulation of the instruction with the data provided by KVM. That might be a program exception or instruction completion.

The notification type intercepts inform KVM about guest environment changes due to guest instruction interpretation. Such an interception is recognized, for example, for the store prefix instruction to provide the new lowcore location. On SIE reentry, any KVM data in the data areas is ignored and execution continues as if the guest instruction had completed. For that reason KVM is not allowed to inject a program interrupt.

### 1.12.6 Links

[KVM Forum 2019 presentation](#)

## 1.13 s390 (IBM Z) Boot/IPL of Protected VMs

### 1.13.1 Summary

The memory of Protected Virtual Machines (PVMs) is not accessible to I/O or the hypervisor. In those cases where the hypervisor needs to access the memory of a PVM, that memory must be made accessible. Memory made accessible to the hypervisor will be encrypted. See *s390 (IBM Z) Ultravisor and Protected VMs* for details.”

On IPL (boot) a small plaintext bootloader is started, which provides information about the encrypted components and necessary metadata to KVM to decrypt the protected virtual machine.

Based on this data, KVM will make the protected virtual machine known to the Ultravisor (UV) and instruct it to secure the memory of the PVM, decrypt the components and verify the data and address list hashes, to ensure integrity. Afterwards KVM can run the PVM via the SIE instruction which the UV will intercept and execute on KVM's behalf.

As the guest image is just like an opaque kernel image that does the switch into PV mode itself, the user can load encrypted guest executables and data via every available method (network, dasd, scsi, direct kernel, ...) without the need to change the boot process.

### 1.13.2 Diag308

This diagnose instruction is the basic mechanism to handle IPL and related operations for virtual machines. The VM can set and retrieve IPL information blocks, that specify the IPL method/devices and request VM memory and subsystem resets, as well as IPLs.

For PVMs this concept has been extended with new subcodes:

Subcode 8: Set an IPL Information Block of type 5 (information block for PVMs)

Subcode 9: Store the saved block in guest memory  
Subcode 10: Move into Protected Virtualization mode

The new PV load-device-specific-parameters field specifies all data that is necessary to move into PV mode.

- PV Header origin
- PV Header length
- **List of Components composed of**
  - AES-XTS Tweak prefix
  - Origin
  - Size

The PV header contains the keys and hashes, which the UV will use to decrypt and verify the PV, as well as control flags and a start PSW.

The components are for instance an encrypted kernel, kernel parameters and initrd. The components are decrypted by the UV.

After the initial import of the encrypted data, all defined pages will contain the guest content. All non-specified pages will start out as zero pages on first access.

When running in protected virtualization mode, some subcodes will result in exceptions or return error codes.

Subcodes 4 and 7, which specify operations that do not clear the guest memory, will result in specification exceptions. This is because the UV will clear all memory when a secure VM is removed, and therefore non-clearing IPL subcodes are not allowed.

Subcodes 8, 9, 10 will result in specification exceptions. Re-IPL into a protected mode is only possible via a detour into non protected mode.

### 1.13.3 Keys

Every CEC will have a unique public key to enable tooling to build encrypted images. See [s390-tools](#) for the tooling.

## 1.14 Timekeeping Virtualization for X86-Based Architectures

### Author

Zachary Amsden <[zamsden@redhat.com](mailto:zamsden@redhat.com)>

### Copyright

(c) 2010, Red Hat. All rights reserved.

### 1.14.1 1. Overview

One of the most complicated parts of the X86 platform, and specifically, the virtualization of this platform is the plethora of timing devices available and the complexity of emulating those devices. In addition, virtualization of time introduces a new set of challenges because it introduces a multiplexed division of time beyond the control of the guest CPU.

First, we will describe the various timekeeping hardware available, then present some of the problems which arise and solutions available, giving specific recommendations for certain classes of KVM guests.

The purpose of this document is to collect data and information relevant to timekeeping which may be difficult to find elsewhere, specifically, information relevant to KVM and hardware-based virtualization.

### 1.14.2 2. Timing Devices

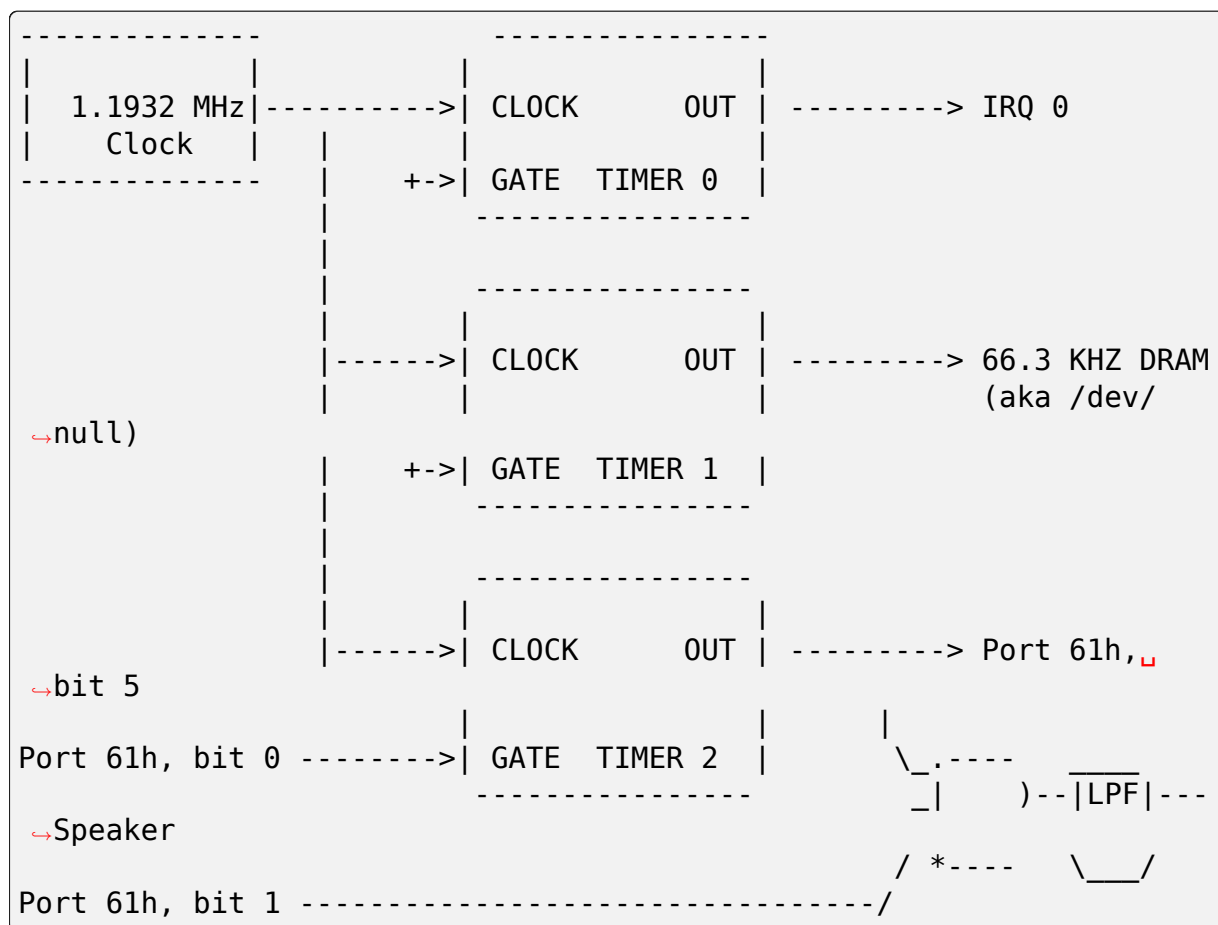
First we discuss the basic hardware devices available. TSC and the related KVM clock are special enough to warrant a full exposition and are described in the following section.

#### 2.1. i8254 - PIT

One of the first timer devices available is the programmable interrupt timer, or PIT. The PIT has a fixed frequency 1.193182 MHz base clock and three channels which can be programmed to deliver periodic or one-shot interrupts. These three channels can be configured in different modes and have individual counters. Channel 1 and 2 were not available for general use in the original IBM PC, and historically were connected to control RAM refresh and the PC speaker. Now the PIT is typically integrated as part of an emulated chipset and a separate physical PIT is not used.

The PIT uses I/O ports 0x40 - 0x43. Access to the 16-bit counters is done using single or multiple byte access to the I/O ports. There are 6 modes available, but

not all modes are available to all timers, as only timer 2 has a connected gate input, required for modes 1 and 5. The gate line is controlled by port 61h, bit 0, as illustrated in the following diagram:



The timer modes are now described.

### Mode 0: Single Timeout.

This is a one-shot software timeout that counts down when the gate is high (always true for timers 0 and 1). When the count reaches zero, the output goes high.

### Mode 1: Triggered One-shot.

The output is initially set high. When the gate line is set high, a countdown is initiated (which does not stop if the gate is lowered), during which the output is set low. When the count reaches zero, the output goes high.

### Mode 2: Rate Generator.

The output is initially set high. When the countdown reaches 1, the output goes low for one count and then returns high. The value is reloaded and the countdown automatically resumes. If the gate line goes low, the count is halted. If the output is low when the gate is lowered, the output automatically goes high (this only affects timer 2).

### Mode 3: Square Wave.

This generates a high / low square wave. The count determines the length of the pulse, which alternates between high and low when zero is reached. The count only proceeds when gate is high and is automatically reloaded on

reaching zero. The count is decremented twice at each clock to generate a full high / low cycle at the full periodic rate. If the count is even, the clock remains high for  $N/2$  counts and low for  $N/2$  counts; if the clock is odd, the clock is high for  $(N+1)/2$  counts and low for  $(N-1)/2$  counts. Only even values are latched by the counter, so odd values are not observed when reading. This is the intended mode for timer 2, which generates sine-like tones by low-pass filtering the square wave output.

#### Mode 4: Software Strobe.

After programming this mode and loading the counter, the output remains high until the counter reaches zero. Then the output goes low for 1 clock cycle and returns high. The counter is not reloaded. Counting only occurs when gate is high.

#### Mode 5: Hardware Strobe.

After programming and loading the counter, the output remains high. When the gate is raised, a countdown is initiated (which does not stop if the gate is lowered). When the counter reaches zero, the output goes low for 1 clock cycle and then returns high. The counter is not reloaded.

In addition to normal binary counting, the PIT supports BCD counting. The command port, 0x43 is used to set the counter and mode for each of the three timers.

PIT commands, issued to port 0x43, using the following bit encoding:

Bit 7-4: Command (See table below)  
 Bit 3-1: Mode (000 = Mode 0, 101 = Mode 5, 11X = undefined)  
 Bit 0 : Binary (0) / BCD (1)

Command table:

```
0000 - Latch Timer 0 count for port 0x40
       sample and hold the count to be read in port 0x40;
       additional commands ignored until counter is read;
       mode bits ignored.

0001 - Set Timer 0 LSB mode for port 0x40
       set timer to read LSB only and force MSB to zero;
       mode bits set timer mode

0010 - Set Timer 0 MSB mode for port 0x40
       set timer to read MSB only and force LSB to zero;
       mode bits set timer mode

0011 - Set Timer 0 16-bit mode for port 0x40
       set timer to read / write LSB first, then MSB;
       mode bits set timer mode

0100 - Latch Timer 1 count for port 0x41 - as described above
0101 - Set Timer 1 LSB mode for port 0x41 - as described above
0110 - Set Timer 1 MSB mode for port 0x41 - as described above
0111 - Set Timer 1 16-bit mode for port 0x41 - as described above
```

(continues on next page)

(continued from previous page)

```
1000 - Latch Timer 2 count for port 0x42 - as described above
1001 - Set Timer 2 LSB mode for port 0x42 - as described above
1010 - Set Timer 2 MSB mode for port 0x42 - as described above
1011 - Set Timer 2 16-bit mode for port 0x42 as described above

1101 - General counter latch
      Latch combination of counters into corresponding ports
      Bit 3 = Counter 2
      Bit 2 = Counter 1
      Bit 1 = Counter 0
      Bit 0 = Unused

1110 - Latch timer status
      Latch combination of counter mode into corresponding ports
      Bit 3 = Counter 2
      Bit 2 = Counter 1
      Bit 1 = Counter 0

      The output of ports 0x40-0x42 following this command will be:

      Bit 7 = Output pin
      Bit 6 = Count loaded (0 if timer has expired)
      Bit 5-4 = Read / Write mode
                01 = MSB only
                10 = LSB only
                11 = LSB / MSB (16-bit)
      Bit 3-1 = Mode
      Bit 0 = Binary (0) / BCD mode (1)
```

## 2.2. RTC

The second device which was available in the original PC was the MC146818 real time clock. The original device is now obsolete, and usually emulated by the system chipset, sometimes by an HPET and some frankenstein IRQ routing.

The RTC is accessed through CMOS variables, which uses an index register to control which bytes are read. Since there is only one index register, read of the CMOS and read of the RTC require lock protection (in addition, it is dangerous to allow userspace utilities such as `hwclock` to have direct RTC access, as they could corrupt kernel reads and writes of CMOS memory).

The RTC generates an interrupt which is usually routed to IRQ 8. The interrupt can function as a periodic timer, an additional once a day alarm, and can issue interrupts after an update of the CMOS registers by the MC146818 is complete. The type of interrupt is signalled in the RTC status registers.

The RTC will update the current time fields by battery power even while the system is off. The current time fields should not be read while an update is in progress, as indicated in the status register.

The clock uses a 32.768kHz crystal, so bits 6-4 of register A should be programmed



to a 32kHz divider if the RTC is to count seconds.

This is the RAM map originally used for the RTC/CMOS:

Location	Size	Description
-----		
00h	byte	Current second (BCD)
01h	byte	Seconds alarm (BCD)
02h	byte	Current minute (BCD)
03h	byte	Minutes alarm (BCD)
04h	byte	Current hour (BCD)
05h	byte	Hours alarm (BCD)
06h	byte	Current day of week (BCD)
07h	byte	Current day of month (BCD)
08h	byte	Current month (BCD)
09h	byte	Current year (BCD)
0Ah	byte	Register A bit 7 = Update in progress bit 6-4 = Divider for clock 000 = 4.194 MHz 001 = 1.049 MHz 010 = 32 kHz 10X = test modes 110 = reset / disable 111 = reset / disable bit 3-0 = Rate selection for periodic interrupt 000 = periodic timer disabled 001 = 3.90625 uS 010 = 7.8125 uS 011 = .122070 mS 100 = .244141 mS ... 1101 = 125 mS 1110 = 250 mS 1111 = 500 mS
0Bh	byte	Register B bit 7 = Run (0) / Halt (1) bit 6 = Periodic interrupt enable bit 5 = Alarm interrupt enable bit 4 = Update-ended interrupt enable bit 3 = Square wave interrupt enable bit 2 = BCD calendar (0) / Binary (1) bit 1 = 12-hour mode (0) / 24-hour mode (1) bit 0 = 0 (DST off) / 1 (DST enabled)
0Ch	byte	Register C (read only) bit 7 = interrupt request flag (IRQF) bit 6 = periodic interrupt flag (PF) bit 5 = alarm interrupt flag (AF) bit 4 = update interrupt flag (UF) bit 3-0 = reserved
0Dh	byte	Register D (read only)

(continues on next page)

(continued from previous page)

```

                                bit 7   = RTC has power
                                bit 6-0 = reserved
32h          byte   Current century BCD (*)
(*) location vendor specific and now determined from ACPI global_
→ tables

```

## 2.3. APIC

On Pentium and later processors, an on-board timer is available to each CPU as part of the Advanced Programmable Interrupt Controller. The APIC is accessed through memory-mapped registers and provides interrupt service to each CPU, used for IPIs and local timer interrupts.

Although in theory the APIC is a safe and stable source for local interrupts, in practice, many bugs and glitches have occurred due to the special nature of the APIC CPU-local memory-mapped hardware. Beware that CPU errata may affect the use of the APIC and that workarounds may be required. In addition, some of these workarounds pose unique constraints for virtualization - requiring either extra overhead incurred from extra reads of memory-mapped I/O or additional functionality that may be more computationally expensive to implement.

Since the APIC is documented quite well in the Intel and AMD manuals, we will avoid repetition of the detail here. It should be pointed out that the APIC timer is programmed through the LVT (local vector timer) register, is capable of one-shot or periodic operation, and is based on the bus clock divided down by the programmable divider register.

## 2.4. HPET

HPET is quite complex, and was originally intended to replace the PIT / RTC support of the X86 PC. It remains to be seen whether that will be the case, as the de facto standard of PC hardware is to emulate these older devices. Some systems designated as legacy free may support only the HPET as a hardware timer device.

The HPET spec is rather loose and vague, requiring at least 3 hardware timers, but allowing implementation freedom to support many more. It also imposes no fixed rate on the timer frequency, but does impose some extremal values on frequency, error and slew.

In general, the HPET is recommended as a high precision (compared to PIT /RTC) time source which is independent of local variation (as there is only one HPET in any given system). The HPET is also memory-mapped, and its presence is indicated through ACPI tables by the BIOS.

Detailed specification of the HPET is beyond the current scope of this document, as it is also very well documented elsewhere.

## 2.5. Offboard Timers

Several cards, both proprietary (watchdog boards) and commonplace (e1000) have timing chips built into the cards which may have registers which are accessible to kernel or user drivers. To the author's knowledge, using these to generate a clocksource for a Linux or other kernel has not yet been attempted and is in general frowned upon as not playing by the agreed rules of the game. Such a timer device would require additional support to be virtualized properly and is not considered important at this time as no known operating system does this.

### 1.14.3 3. TSC Hardware

The TSC or time stamp counter is relatively simple in theory; it counts instruction cycles issued by the processor, which can be used as a measure of time. In practice, due to a number of problems, it is the most complicated timekeeping device to use.

The TSC is represented internally as a 64-bit MSR which can be read with the RDMSR, RDTSC, or RDTSCP (when available) instructions. In the past, hardware limitations made it possible to write the TSC, but generally on old hardware it was only possible to write the low 32-bits of the 64-bit counter, and the upper 32-bits of the counter were cleared. Now, however, on Intel processors family 0Fh, for models 3, 4 and 6, and family 06h, models e and f, this restriction has been lifted and all 64-bits are writable. On AMD systems, the ability to write the TSC MSR is not an architectural guarantee.

The TSC is accessible from CPL-0 and conditionally, for CPL > 0 software by means of the CR4.TSD bit, which when enabled, disables CPL > 0 TSC access.

Some vendors have implemented an additional instruction, RDTSCP, which returns atomically not just the TSC, but an indicator which corresponds to the processor number. This can be used to index into an array of TSC variables to determine offset information in SMP systems where TSCs are not synchronized. The presence of this instruction must be determined by consulting CPUID feature bits.

Both VMX and SVM provide extension fields in the virtualization hardware which allows the guest visible TSC to be offset by a constant. Newer implementations promise to allow the TSC to additionally be scaled, but this hardware is not yet widely available.

### 3.1. TSC synchronization

The TSC is a CPU-local clock in most implementations. This means, on SMP platforms, the TSCs of different CPUs may start at different times depending on when the CPUs are powered on. Generally, CPUs on the same die will share the same clock, however, this is not always the case.

The BIOS may attempt to resynchronize the TSCs during the poweron process and the operating system or other system software may attempt to do this as well. Several hardware limitations make the problem worse - if it is not possible to write the full 64-bits of the TSC, it may be impossible to match the TSC in newly arriving CPUs to that of the rest of the system, resulting in unsynchronized TSCs. This

may be done by BIOS or system software, but in practice, getting a perfectly synchronized TSC will not be possible unless all values are read from the same clock, which generally only is possible on single socket systems or those with special hardware support.

### 3.2. TSC and CPU hotplug

As touched on already, CPUs which arrive later than the boot time of the system may not have a TSC value that is synchronized with the rest of the system. Either system software, BIOS, or SMM code may actually try to establish the TSC to a value matching the rest of the system, but a perfect match is usually not a guarantee. This can have the effect of bringing a system from a state where TSC is synchronized back to a state where TSC synchronization flaws, however small, may be exposed to the OS and any virtualization environment.

### 3.3. TSC and multi-socket / NUMA

Multi-socket systems, especially large multi-socket systems are likely to have individual clocksources rather than a single, universally distributed clock. Since these clocks are driven by different crystals, they will not have perfectly matched frequency, and temperature and electrical variations will cause the CPU clocks, and thus the TSCs to drift over time. Depending on the exact clock and bus design, the drift may or may not be fixed in absolute error, and may accumulate over time.

In addition, very large systems may deliberately slew the clocks of individual cores. This technique, known as spread-spectrum clocking, reduces EMI at the clock frequency and harmonics of it, which may be required to pass FCC standards for telecommunications and computer equipment.

It is recommended not to trust the TSCs to remain synchronized on NUMA or multiple socket systems for these reasons.

### 3.4. TSC and C-states

C-states, or idling states of the processor, especially C1E and deeper sleep states may be problematic for TSC as well. The TSC may stop advancing in such a state, resulting in a TSC which is behind that of other CPUs when execution is resumed. Such CPUs must be detected and flagged by the operating system based on CPU and chipset identifications.

The TSC in such a case may be corrected by catching it up to a known external clocksource.

### 3.5. TSC frequency change / P-states

To make things slightly more interesting, some CPUs may change frequency. They may or may not run the TSC at the same rate, and because the frequency change may be staggered or slewed, at some points in time, the TSC rate may not be known other than falling within a range of values. In this case, the TSC will not be a stable time source, and must be calibrated against a known, stable, external clock to be a usable source of time.

Whether the TSC runs at a constant rate or scales with the P-state is model dependent and must be determined by inspecting CPUID, chipset or vendor specific MSR fields.

In addition, some vendors have known bugs where the P-state is actually compensated for properly during normal operation, but when the processor is inactive, the P-state may be raised temporarily to service cache misses from other processors. In such cases, the TSC on halted CPUs could advance faster than that of non-halted processors. AMD Turion processors are known to have this problem.

### 3.6. TSC and STPCLK / T-states

External signals given to the processor may also have the effect of stopping the TSC. This is typically done for thermal emergency power control to prevent an overheating condition, and typically, there is no way to detect that this condition has happened.

### 3.7. TSC virtualization - VMX

VMX provides conditional trapping of RDTSC, RDMSR, WRMSR and RDTSCP instructions, which is enough for full virtualization of TSC in any manner. In addition, VMX allows passing through the host TSC plus an additional TSC\_OFFSET field specified in the VMCS. Special instructions must be used to read and write the VMCS field.

### 3.8. TSC virtualization - SVM

SVM provides conditional trapping of RDTSC, RDMSR, WRMSR and RDTSCP instructions, which is enough for full virtualization of TSC in any manner. In addition, SVM allows passing through the host TSC plus an additional offset field specified in the SVM control block.

### 3.9. TSC feature bits in Linux

In summary, there is no way to guarantee the TSC remains in perfect synchronization unless it is explicitly guaranteed by the architecture. Even if so, the TSCs in multi-sockets or NUMA systems may still run independently despite being locally consistent.

The following feature bits are used by Linux to signal various TSC attributes, but they can only be taken to be meaningful for UP or single node systems.

X86_FEATURE_TSC	The TSC is available in hardware
X86_FEATURE_RDTSCP	The RDTSCP instruction is available
X86_FEATURE_CONSTANT_TSC	The TSC rate is unchanged with P-states
X86_FEATURE_NONSTOP_TSC	The TSC does not stop in C-states
X86_FEATURE_TSC_RELIABLE	TSC sync checks are skipped (VMware)

## 1.14.4 4. Virtualization Problems

Timekeeping is especially problematic for virtualization because a number of challenges arise. The most obvious problem is that time is now shared between the host and, potentially, a number of virtual machines. Thus the virtual operating system does not run with 100% usage of the CPU, despite the fact that it may very well make that assumption. It may expect it to remain true to very exacting bounds when interrupt sources are disabled, but in reality only its virtual interrupt sources are disabled, and the machine may still be preempted at any time. This causes problems as the passage of real time, the injection of machine interrupts and the associated clock sources are no longer completely synchronized with real time.

This same problem can occur on native hardware to a degree, as SMM mode may steal cycles from the naturally on X86 systems when SMM mode is used by the BIOS, but not in such an extreme fashion. However, the fact that SMM mode may cause similar problems to virtualization makes it a good justification for solving many of these problems on bare metal.

### 4.1. Interrupt clocking

One of the most immediate problems that occurs with legacy operating systems is that the system timekeeping routines are often designed to keep track of time by counting periodic interrupts. These interrupts may come from the PIT or the RTC, but the problem is the same: the host virtualization engine may not be able to deliver the proper number of interrupts per second, and so guest time may fall behind. This is especially problematic if a high interrupt rate is selected, such as 1000 HZ, which is unfortunately the default for many Linux guests.

There are three approaches to solving this problem; first, it may be possible to simply ignore it. Guests which have a separate time source for tracking ‘wall clock’ or ‘real time’ may not need any adjustment of their interrupts to maintain proper time. If this is not sufficient, it may be necessary to inject additional interrupts into the guest in order to increase the effective interrupt rate. This approach

leads to complications in extreme conditions, where host load or guest lag is too much to compensate for, and thus another solution to the problem has risen: the guest may need to become aware of lost ticks and compensate for them internally. Although promising in theory, the implementation of this policy in Linux has been extremely error prone, and a number of buggy variants of lost tick compensation are distributed across commonly used Linux systems.

Windows uses periodic RTC clocking as a means of keeping time internally, and thus requires interrupt slewing to keep proper time. It does use a low enough rate (ed: is it 18.2 Hz?) however that it has not yet been a problem in practice.

## 4.2. TSC sampling and serialization

As the highest precision time source available, the cycle counter of the CPU has aroused much interest from developers. As explained above, this timer has many problems unique to its nature as a local, potentially unstable and potentially unsynchronized source. One issue which is not unique to the TSC, but is highlighted because of its very precise nature is sampling delay. By definition, the counter, once read is already old. However, it is also possible for the counter to be read ahead of the actual use of the result. This is a consequence of the superscalar execution of the instruction stream, which may execute instructions out of order. Such execution is called non-serialized. Forcing serialized execution is necessary for precise measurement with the TSC, and requires a serializing instruction, such as CPUID or an MSR read.

Since CPUID may actually be virtualized by a trap and emulate mechanism, this serialization can pose a performance issue for hardware virtualization. An accurate time stamp counter reading may therefore not always be available, and it may be necessary for an implementation to guard against “backwards” reads of the TSC as seen from other CPUs, even in an otherwise perfectly synchronized system.

## 4.3. Timespec aliasing

Additionally, this lack of serialization from the TSC poses another challenge when using results of the TSC when measured against another time source. As the TSC is much higher precision, many possible values of the TSC may be read while another clock is still expressing the same value.

That is, you may read  $(T, T+10)$  while external clock  $C$  maintains the same value. Due to non-serialized reads, you may actually end up with a range which fluctuates - from  $(T-1.. T+10)$ . Thus, any time calculated from a TSC, but calibrated against an external value may have a range of valid values. Re-calibrating this computation may actually cause time, as computed after the calibration, to go backwards, compared with time computed before the calibration.

This problem is particularly pronounced with an internal time source in Linux, the kernel time, which is expressed in the theoretically high resolution timespec - but which advances in much larger granularity intervals, sometimes at the rate of jiffies, and possibly in catchup modes, at a much larger step.

This aliasing requires care in the computation and recalibration of `kvmclock` and any other values derived from TSC computation (such as TSC virtualization itself).



### 4.4. Migration

Migration of a virtual machine raises problems for timekeeping in two ways. First, the migration itself may take time, during which interrupts cannot be delivered, and after which, the guest time may need to be caught up. NTP may be able to help to some degree here, as the clock correction required is typically small enough to fall in the NTP-correctable window.

An additional concern is that timers based off the TSC (or HPET, if the raw bus clock is exposed) may now be running at different rates, requiring compensation in some way in the hypervisor by virtualizing these timers. In addition, migrating to a faster machine may preclude the use of a passthrough TSC, as a faster clock cannot be made visible to a guest without the potential of time advancing faster than usual. A slower clock is less of a problem, as it can always be caught up to the original rate. KVM clock avoids these problems by simply storing multipliers and offsets against the TSC for the guest to convert back into nanosecond resolution values.

### 4.5. Scheduling

Since scheduling may be based on precise timing and firing of interrupts, the scheduling algorithms of an operating system may be adversely affected by virtualization. In theory, the effect is random and should be universally distributed, but in contrived as well as real scenarios (guest device access, causes of virtualization exits, possible context switch), this may not always be the case. The effect of this has not been well studied.

In an attempt to work around this, several implementations have provided a paravirtualized scheduler clock, which reveals the true amount of CPU time for which a virtual machine has been running.

### 4.6. Watchdogs

Watchdog timers, such as the lock detector in Linux may fire accidentally when running under hardware virtualization due to timer interrupts being delayed or misinterpretation of the passage of real time. Usually, these warnings are spurious and can be ignored, but in some circumstances it may be necessary to disable such detection.

### 4.7. Delays and precision timing

Precise timing and delays may not be possible in a virtualized system. This can happen if the system is controlling physical hardware, or issues delays to compensate for slower I/O to and from devices. The first issue is not solvable in general for a virtualized system; hardware control software can't be adequately virtualized without a full real-time operating system, which would require an RT aware virtualization platform.

The second issue may cause performance problems, but this is unlikely to be a significant issue. In many cases these delays may be eliminated through configuration or paravirtualization.



## 4.8. Covert channels and leaks

In addition to the above problems, time information will inevitably leak to the guest about the host in anything but a perfect implementation of virtualized time. This may allow the guest to infer the presence of a hypervisor (as in a red-pill type detection), and it may allow information to leak between guests by using CPU utilization itself as a signalling channel. Preventing such problems would require completely isolated virtual time which may not track real time any longer. This may be useful in certain security or QA contexts, but in general isn't recommended for real-world deployment scenarios.

## 1.15 KVM VCPU Requests

### 1.15.1 Overview

KVM supports an internal API enabling threads to request a VCPU thread to perform some activity. For example, a thread may request a VCPU to flush its TLB with a VCPU request. The API consists of the following functions:

```
/* Check if any requests are pending for VCPU @vcpu. */
bool kvm_request_pending(struct kvm_vcpu *vcpu);

/* Check if VCPU @vcpu has request @req pending. */
bool kvm_test_request(int req, struct kvm_vcpu *vcpu);

/* Clear request @req for VCPU @vcpu. */
void kvm_clear_request(int req, struct kvm_vcpu *vcpu);

/*
 * Check if VCPU @vcpu has request @req pending. When the request is
 * pending it will be cleared and a memory barrier, which pairs with
 * another in kvm_make_request(), will be issued.
 */
bool kvm_check_request(int req, struct kvm_vcpu *vcpu);

/*
 * Make request @req of VCPU @vcpu. Issues a memory barrier, which
 * pairs
 * with another in kvm_check_request(), prior to setting the
 * request.
 */
void kvm_make_request(int req, struct kvm_vcpu *vcpu);

/* Make request @req of all VCPUs of the VM with struct kvm @kvm. */
bool kvm_make_all_cpus_request(struct kvm *kvm, unsigned int req);
```

Typically a requester wants the VCPU to perform the activity as soon as possible after making the request. This means most requests (`kvm_make_request()` calls) are followed by a call to `kvm_vcpu_kick()`, and `kvm_make_all_cpus_request()` has the kicking of all VCPUs built into it.

## VCPU Kicks

The goal of a VCPU kick is to bring a VCPU thread out of guest mode in order to perform some KVM maintenance. To do so, an IPI is sent, forcing a guest mode exit. However, a VCPU thread may not be in guest mode at the time of the kick. Therefore, depending on the mode and state of the VCPU thread, there are two other actions a kick may take. All three actions are listed below:

- 1) Send an IPI. This forces a guest mode exit.
- 2) Waking a sleeping VCPU. Sleeping VCPUs are VCPU threads outside guest mode that wait on waitqueues. Waking them removes the threads from the waitqueues, allowing the threads to run again. This behavior may be suppressed, see `KVM_REQUEST_NO_WAKEUP` below.
- 3) Nothing. When the VCPU is not in guest mode and the VCPU thread is not sleeping, then there is nothing to do.

## VCPU Mode

VCPUs have a mode state, `vcpu->mode`, that is used to track whether the guest is running in guest mode or not, as well as some specific outside guest mode states. The architecture may use `vcpu->mode` to ensure VCPU requests are seen by VCPUs (see “Ensuring Requests Are Seen”), as well as to avoid sending unnecessary IPIs (see “IPI Reduction”), and even to ensure IPI acknowledgements are waited upon (see “Waiting for Acknowledgements”). The following modes are defined:

### OUTSIDE\_GUEST\_MODE

The VCPU thread is outside guest mode.

### IN\_GUEST\_MODE

The VCPU thread is in guest mode.

### EXITING\_GUEST\_MODE

The VCPU thread is transitioning from `IN_GUEST_MODE` to `OUTSIDE_GUEST_MODE`.

### READING\_SHADOW\_PAGE\_TABLES

The VCPU thread is outside guest mode, but it wants the sender of certain VCPU requests, namely `KVM_REQ_TLB_FLUSH`, to wait until the VCPU thread is done reading the page tables.

## 1.15.2 VCPU Request Internals

VCPU requests are simply bit indices of the `vcpu->requests` bitmap. This means general bitops, like those documented in [\[atomic-ops\]](#) could also be used, e.g.

```
clear_bit(KVM_REQ_UNHALT & KVM_REQUEST_MASK, &vcpu->requests);
```

However, VCPU request users should refrain from doing so, as it would break the abstraction. The first 8 bits are reserved for architecture independent requests, all additional bits are available for architecture dependent requests.

## Architecture Independent Requests

### KVM\_REQ\_TLB\_FLUSH

KVM's common MMU notifier may need to flush all of a guest's TLB entries, calling `kvm_flush_remote_tlbs()` to do so. Architectures that choose to use the common `kvm_flush_remote_tlbs()` implementation will need to handle this VCPU request.

### KVM\_REQ\_MMU\_RELOAD

When shadow page tables are used and memory slots are removed it's necessary to inform each VCPU to completely refresh the tables. This request is used for that.

### KVM\_REQ\_PENDING\_TIMER

This request may be made from a timer handler run on the host on behalf of a VCPU. It informs the VCPU thread to inject a timer interrupt.

### KVM\_REQ\_UNHALT

This request may be made from the KVM common function `kvm_vcpu_block()`, which is used to emulate an instruction that causes a CPU to halt until one of an architectural specific set of events and/or interrupts is received (determined by checking `kvm_arch_vcpu_runnable()`). When that event or interrupt arrives `kvm_vcpu_block()` makes the request. This is in contrast to when `kvm_vcpu_block()` returns due to any other reason, such as a pending signal, which does not indicate the VCPU's halt emulation should stop, and therefore does not make the request.

## KVM\_REQUEST\_MASK

VCPU requests should be masked by `KVM_REQUEST_MASK` before using them with bitops. This is because only the lower 8 bits are used to represent the request's number. The upper bits are used as flags. Currently only two flags are defined.

## VCPU Request Flags

### KVM\_REQUEST\_NO\_WAKEUP

This flag is applied to requests that only need immediate attention from VCPUs running in guest mode. That is, sleeping VCPUs do not need to be awoken for these requests. Sleeping VCPUs will handle the requests when they are awoken later for some other reason.

### KVM\_REQUEST\_WAIT

When requests with this flag are made with `kvm_make_all_cpus_request()`, then the caller will wait for each VCPU to acknowledge its IPI before proceeding. This flag only applies to VCPUs that would receive IPIs. If, for example, the VCPU is sleeping, so no IPI is necessary, then the requesting thread does not wait. This means that this flag may be safely combined with `KVM_REQUEST_NO_WAKEUP`.

See “Waiting for Acknowledgements” for more information about requests with `KVM_REQUEST_WAIT`.

### 1.15.3 VCPU Requests with Associated State

Requesters that want the receiving VCPU to handle new state need to ensure the newly written state is observable to the receiving VCPU thread’s CPU by the time it observes the request. This means a write memory barrier must be inserted after writing the new state and before setting the VCPU request bit. Additionally, on the receiving VCPU thread’s side, a corresponding read barrier must be inserted after reading the request bit and before proceeding to read the new state associated with it. See scenario 3, Message and Flag, of [lwn-mb] and the kernel documentation [memory-barriers].

The pair of functions, `kvm_check_request()` and `kvm_make_request()`, provide the memory barriers, allowing this requirement to be handled internally by the API.

### 1.15.4 Ensuring Requests Are Seen

When making requests to VCPUs, we want to avoid the receiving VCPU executing in guest mode for an arbitrary long time without handling the request. We can be sure this won’t happen as long as we ensure the VCPU thread checks `kvm_request_pending()` before entering guest mode and that a kick will send an IPI to force an exit from guest mode when necessary. Extra care must be taken to cover the period after the VCPU thread’s last `kvm_request_pending()` check and before it has entered guest mode, as kick IPIs will only trigger guest mode exits for VCPU threads that are in guest mode or at least have already disabled interrupts in order to prepare to enter guest mode. This means that an optimized implementation (see “IPI Reduction”) must be certain when it’s safe to not send the IPI. One solution, which all architectures except s390 apply, is to:

- set `vcpu->mode` to `IN_GUEST_MODE` between disabling the interrupts and the last `kvm_request_pending()` check;
- enable interrupts atomically when entering the guest.

This solution also requires memory barriers to be placed carefully in both the requesting thread and the receiving VCPU. With the memory barriers we can exclude the possibility of a VCPU thread observing `!kvm_request_pending()` on its last check and then not receiving an IPI for the next request made of it, even if the request is made immediately after the check. This is done by way of the Dekker memory barrier pattern (scenario 10 of [lwn-mb]). As the Dekker pattern requires two variables, this solution pairs `vcpu->mode` with `vcpu->requests`. Substituting them into the pattern gives:

CPU1	CPU2
=====	=====
<code>local_irq_disable();</code>	
<code>WRITE_ONCE(vcpu-&gt;mode, IN_GUEST_MODE);</code>	<code>kvm_make_request(REQ, vcpu);</code>
<code>smp_mb();</code>	<code>smp_mb();</code>
<code>if (kvm_request_pending(vcpu)) {</code>	<code>if (READ_ONCE(vcpu-&gt;mode) ==</code>

(continues on next page)

(continued from previous page)

<pre> ...abort guest entry... } </pre>	<pre> IN_GUEST_MODE) { ...send IPI... } </pre>
--	--

As stated above, the IPI is only useful for VCPU threads in guest mode or that have already disabled interrupts. This is why this specific case of the Dekker pattern has been extended to disable interrupts before setting `vcpu->mode` to `IN_GUEST_MODE`. `WRITE_ONCE()` and `READ_ONCE()` are used to pedantically implement the memory barrier pattern, guaranteeing the compiler doesn't interfere with `vcpu->mode`'s carefully planned accesses.

## IPI Reduction

As only one IPI is needed to get a VCPU to check for any/all requests, then they may be coalesced. This is easily done by having the first IPI sending kick also change the VCPU mode to something `!IN_GUEST_MODE`. The transitional state, `EXITING_GUEST_MODE`, is used for this purpose.

## Waiting for Acknowledgements

Some requests, those with the `KVM_REQUEST_WAIT` flag set, require IPIs to be sent, and the acknowledgements to be waited upon, even when the target VCPU threads are in modes other than `IN_GUEST_MODE`. For example, one case is when a target VCPU thread is in `READING_SHADOW_PAGE_TABLES` mode, which is set after disabling interrupts. To support these cases, the `KVM_REQUEST_WAIT` flag changes the condition for sending an IPI from checking that the VCPU is `IN_GUEST_MODE` to checking that it is not `OUTSIDE_GUEST_MODE`.

## Request-less VCPU Kicks

As the determination of whether or not to send an IPI depends on the two-variable Dekker memory barrier pattern, then it's clear that request-less VCPU kicks are almost never correct. Without the assurance that a non-IPI generating kick will still result in an action by the receiving VCPU, as the final `kvm_request_pending()` check does for request-accompanying kicks, then the kick may not do anything useful at all. If, for instance, a request-less kick was made to a VCPU that was just about to set its mode to `IN_GUEST_MODE`, meaning no IPI is sent, then the VCPU thread may continue its entry without actually having done whatever it was the kick was meant to initiate.

One exception is x86's posted interrupt mechanism. In this case, however, even the request-less VCPU kick is coupled with the same `local_irq_disable()` + `smp_mb()` pattern described above; the ON bit (Outstanding Notification) in the posted interrupt descriptor takes the role of `vcpu->requests`. When sending a posted interrupt, `PIR.ON` is set before reading `vcpu->mode`; dually, in the VCPU thread, `vmx_sync_pir_to_irr()` reads `PIR` after setting `vcpu->mode` to `IN_GUEST_MODE`.

### 1.15.5 Additional Considerations

#### Sleeping VCPUs

VCPU threads may need to consider requests before and/or after calling functions that may put them to sleep, e.g. `kvm_vcpu_block()`. Whether they do or not, and, if they do, which requests need consideration, is architecture dependent. `kvm_vcpu_block()` calls `kvm_arch_vcpu_runnable()` to check if it should awaken. One reason to do so is to provide architectures a function where requests may be checked if necessary.

#### Clearing Requests

Generally it only makes sense for the receiving VCPU thread to clear a request. However, in some circumstances, such as when the requesting thread and the receiving VCPU thread are executed serially, such as when they are the same thread, or when they are using some form of concurrency control to temporarily execute synchronously, then it's possible to know that the request may be cleared immediately, rather than waiting for the receiving VCPU thread to handle the request in VCPU RUN. The only current examples of this are `kvm_vcpu_block()` calls made by VCPUs to block themselves. A possible side-effect of that call is to make the `KVM_REQ_UNHALT` request, which may then be cleared immediately when the VCPU returns from the call.

### 1.15.6 References

## 1.16 Review checklist for kvm patches

1. The patch must follow `Documentation/process/coding-style.rst` and `Documentation/process/submitting-patches.rst`.
2. Patches should be against `kvm.git` master branch.
3. If the patch introduces or modifies a new userspace API: - the API must be documented in *The Definitive KVM (Kernel-based Virtual Machine) API Documentation* - the API must be discoverable using `KVM_CHECK_EXTENSION`
4. New state must include support for save/restore.
5. New features must default to off (userspace should explicitly request them). Performance improvements can and should default to on.
6. New cpu features should be exposed via `KVM_GET_SUPPORTED_CPUID2`
7. Emulator changes should be accompanied by unit tests for `qemu-kvm.git` `kvm/test` directory.
8. Changes should be vendor neutral when possible. Changes to common code are better than duplicating changes to vendor code.
9. Similarly, prefer changes to arch independent code than to arch dependent code.

10. User/kernel interfaces and guest/host interfaces must be 64-bit clean (all variables and sizes naturally aligned on 64-bit; use specific types only - u64 rather than ulong).
11. New guest visible features must either be documented in a hardware manual or be accompanied by documentation.
12. Features must be robust against reset and kexec - for example, shared host/guest memory must be unshared to prevent the host from writing to guest memory that the guest has not reserved for this purpose.

## 1.17 ARM

### 1.17.1 Internal ABI between the kernel and HYP

This file documents the interaction between the Linux kernel and the hypervisor layer when running Linux as a hypervisor (for example KVM). It doesn't cover the interaction of the kernel with the hypervisor when running as a guest (under Xen, KVM or any other hypervisor), or any hypervisor-specific interaction when the kernel is used as a host.

Note: KVM/arm has been removed from the kernel. The API described here is still valid though, as it allows the kernel to kexec when booted at HYP. It can also be used by a hypervisor other than KVM if necessary.

On arm and arm64 (without VHE), the kernel doesn't run in hypervisor mode, but still needs to interact with it, allowing a built-in hypervisor to be either installed or torn down.

In order to achieve this, the kernel must be booted at HYP (arm) or EL2 (arm64), allowing it to install a set of stubs before dropping to SVC/EL1. These stubs are accessible by using a 'hvc #0' instruction, and only act on individual CPUs.

Unless specified otherwise, any built-in hypervisor must implement these functions (see arch/arm{,64}/include/asm/virt.h):

- `r0/x0 = HVC_SET_VECTORS`  
`r1/x1 = vectors`

Set HVBAR/VBAR\_EL2 to 'vectors' to enable a hypervisor. 'vectors' must be a physical address, and respect the alignment requirements of the architecture. Only implemented by the initial stubs, not by Linux hypervisors.

- `r0/x0 = HVC_RESET_VECTORS`

Turn HYP/EL2 MMU off, and reset HVBAR/VBAR\_EL2 to the initial stubs' exception vector value. This effectively disables an existing hypervisor.

- `r0/x0 = HVC_SOFT_RESTART`  
`r1/x1 = restart address`  
`x2 = x0's value when entering the next payload (arm64)`  
`x3 = x1's value when entering the next payload (arm64)`  
`x4 = x2's value when entering the next payload (arm64)`



Mask all exceptions, disable the MMU, clear I+D bits, move the arguments into place (arm64 only), and jump to the restart address while at HYP/EL2. This hypercall is not expected to return to its caller.

Any other value of r0/x0 triggers a hypervisor-specific handling, which is not documented here.

The return value of a stub hypercall is held by r0/x0, and is 0 on success, and HVC\_STUB\_ERR on error. A stub hypercall is allowed to clobber any of the caller-saved registers (x0-x18 on arm64, r0-r3 and ip on arm). It is thus recommended to use a function call to perform the hypercall.

### 1.17.2 Power State Coordination Interface (PSCI)

KVM implements the PSCI (Power State Coordination Interface) specification in order to provide services such as CPU on/off, reset and power-off to the guest.

The PSCI specification is regularly updated to provide new features, and KVM implements these updates if they make sense from a virtualization point of view.

This means that a guest booted on two different versions of KVM can observe two different “firmware” revisions. This could cause issues if a given guest is tied to a particular PSCI revision (unlikely), or if a migration causes a different PSCI version to be exposed out of the blue to an unsuspecting guest.

In order to remedy this situation, KVM exposes a set of “firmware pseudo-registers” that can be manipulated using the GET/SET\_ONE\_REG interface. These registers can be saved/restored by userspace, and set to a convenient value if required.

The following register is defined:

- **KVM\_REG\_ARM\_PSCI\_VERSION:**
  - Only valid if the vcpu has the KVM\_ARM\_VCPU\_PSCI\_0\_2 feature set (and thus has already been initialized)
  - Returns the current PSCI version on GET\_ONE\_REG (defaulting to the highest PSCI version implemented by KVM and compatible with v0.2)
  - Allows any PSCI version implemented by KVM and compatible with v0.2 to be set with SET\_ONE\_REG
  - Affects the whole VM (even if the register view is per-vcpu)
- **KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_1:**

Holds the state of the firmware support to mitigate CVE-2017-5715, as offered by KVM to the guest via a HVC call. The workaround is described under SMCCC\_ARCH\_WORKAROUND\_1 in [1].

Accepted values are:

**KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_1\_NOT\_AVAIL:**  
KVM does not offer firmware support for the workaround. The mitigation status for the guest is unknown.

**KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_1\_AVAIL:**  
The workaround HVC call is available to the guest and required for the mitigation.



**KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_1\_NOT\_REQUIRED:**

The workaround HVC call is available to the guest, but it is not needed on this VCPU.

- **KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_2:**

Holds the state of the firmware support to mitigate CVE-2018-3639, as offered by KVM to the guest via a HVC call. The workaround is described under SMCCC\_ARCH\_WORKAROUND\_2 in<sup>1</sup>.

Accepted values are:

**KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_2\_NOT\_AVAIL:**

A workaround is not available. KVM does not offer firmware support for the workaround.

**KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_2\_UNKNOWN:**

The workaround state is unknown. KVM does not offer firmware support for the workaround.

**KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_2\_AVAIL:**

The workaround is available, and can be disabled by a vCPU. If KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_2\_ENABLED is set, it is active for this vCPU.

**KVM\_REG\_ARM\_SMCCC\_ARCH\_WORKAROUND\_2\_NOT\_REQUIRED:**

The workaround is always active on this vCPU or it is not needed.

### 1.17.3 Paravirtualized time support for arm64

Arm specification DEN0057/A defines a standard for paravirtualised time support for AArch64 guests:

<https://developer.arm.com/docs/den0057/a>

KVM/arm64 implements the stolen time part of this specification by providing some hypervisor service calls to support a paravirtualized guest obtaining a view of the amount of time stolen from its execution.

Two new SMCCC compatible hypercalls are defined:

- PV\_TIME\_FEATURES: 0xC5000020
- PV\_TIME\_ST: 0xC5000021

These are only available in the SMC64/HVC64 calling convention as paravirtualized time is not available to 32 bit Arm guests. The existence of the PV\_FEATURES hypercall should be probed using the SMCCC 1.1 ARCH\_FEATURES mechanism before calling it.

#### PV\_TIME\_FEATURES

<sup>1</sup> [https://developer.arm.com/-/media/developer/pdf/ARM\\_DEN\\_0070A\\_Firmware\\_interfaces\\_for\\_mitigating\\_CVE-2017-5715.pdf](https://developer.arm.com/-/media/developer/pdf/ARM_DEN_0070A_Firmware_interfaces_for_mitigating_CVE-2017-5715.pdf)

Function ID:	(uint32) 0xC5000020
PV_call_i	(uint32) The function to query for support. Currently only PV_TIME_ST is supported.
Return value:	(int64) NOT_SUPPORTED (-1) or SUCCESS (0) if the relevant PV-time feature is supported by the hypervisor.

## PV\_TIME\_ST

Function ID:	(uint32) 0xC5000021
Return value:	(int64) IPA of the stolen time data structure for this VCPU. On failure: NOT_SUPPORTED (-1)

The IPA returned by PV\_TIME\_ST should be mapped by the guest as normal memory with inner and outer write back caching attributes, in the inner shareable domain. A total of 16 bytes from the IPA returned are guaranteed to be meaningfully filled by the hypervisor (see structure below).

PV\_TIME\_ST returns the structure for the calling VCPU.

## Stolen Time

The structure pointed to by the PV\_TIME\_ST hypercall is as follows:

Field	Byte Length	Byte Offset	Description
Revision	4	0	Must be 0 for version 1.0
Attributes	4	4	Must be 0
Stolen time	8	8	Stolen time in unsigned nanoseconds indicating how much time this VCPU thread was involuntarily not running on a physical CPU.

All values in the structure are stored little-endian.

The structure will be updated by the hypervisor prior to scheduling a VCPU. It will be present within a reserved region of the normal memory given to the guest. The guest should not attempt to write into this memory. There is a structure per VCPU of the guest.

It is advisable that one or more 64k pages are set aside for the purpose of these structures and not used for other purposes, this enables the guest to map the region using 64k pages and avoids conflicting attributes with other memory.

For the user space interface see [Generic vcpu interface](#) section “3. GROUP: KVM\_ARM\_VCPU\_PVTIME\_CTRL” .

## 1.18 Devices

### 1.18.1 ARM Virtual Interrupt Translation Service (ITS)

#### Device types supported:

KVM\_DEV\_TYPE\_ARM\_VGIC\_ITS ARM Interrupt Translation Service Controller

The ITS allows MSI(-X) interrupts to be injected into guests. This extension is optional. Creating a virtual ITS controller also requires a host GICv3 (see arm-vgic-v3.txt), but does not depend on having physical ITS controllers.

There can be multiple ITS controllers per guest, each of them has to have a separate, non-overlapping MMIO region.

#### Groups

#### KVM\_DEV\_ARM\_VGIC\_GRP\_ADDR

##### Attributes:

##### KVM\_VGIC\_ITS\_ADDR\_TYPE (rw, 64-bit)

Base address in the guest physical address space of the GICv3 ITS control register frame. This address needs to be 64K aligned and the region covers 128K.

Errors:

-E2BIG	Address outside of addressable IPA range
-EINVAL	Incorrectly aligned address
-EEXIST	Address already configured
-EFAULT	Invalid user pointer for attr->addr.
-ENODEV	Incorrect attribute or the ITS is not supported.

#### KVM\_DEV\_ARM\_VGIC\_GRP\_CTRL

##### Attributes:

##### KVM\_DEV\_ARM\_VGIC\_CTRL\_INIT

request the initialization of the ITS, no additional parameter in kvm\_device\_attr.addr.

##### KVM\_DEV\_ARM\_ITS\_CTRL\_RESET

reset the ITS, no additional parameter in kvm\_device\_attr.addr. See “ITS Reset State” section.

**KVM\_DEV\_ARM\_ITS\_SAVE\_TABLES**

save the ITS table data into guest RAM, at the location provided by the guest in corresponding registers/table entries.

The layout of the tables in guest memory defines an ABI. The entries are laid out in little endian format as described in the last paragraph.

**KVM\_DEV\_ARM\_ITS\_RESTORE\_TABLES**

restore the ITS tables from guest RAM to ITS internal structures.

The GICV3 must be restored before the ITS and all ITS registers but the GITS\_CTLR must be restored before restoring the ITS tables.

The GITS\_IIDR read-only register must also be restored before calling KVM\_DEV\_ARM\_ITS\_RESTORE\_TABLES as the IIDR revision field encodes the ABI revision.

The expected ordering when restoring the GICv3/ITS is described in section “ITS Restore Sequence” .

Errors:

- ENXIO	ITS not properly configured as required prior to setting this attribute
- ENOME	Memory shortage when allocating ITS internal data
- EINVAL	Inconsistent restored data
- EFAULT	Invalid guest ram access
- EBUSY	One or more VCPUS are running
- EACCES	The virtual ITS is backed by a physical GICv4 ITS, and the state is not available

**KVM\_DEV\_ARM\_VGIC\_GRP\_ITS\_REGS****Attributes:**

The attr field of kvm\_device\_attr encodes the offset of the ITS register, relative to the ITS control frame base address (ITS\_base).

kvm\_device\_attr.addr points to a \_\_u64 value whatever the width of the addressed register (32/64 bits). 64 bit registers can only be accessed with full length.

Writes to read-only registers are ignored by the kernel except for:

- GITS\_CREADR. It must be restored otherwise commands in the queue will be re-executed after restoring CWRITER. GITS\_CREADR must be restored before restoring the GITS\_CTLR which is likely to enable the ITS. Also it must

be restored after `GITS_CBASER` since a write to `GITS_CBASER` resets `GITS_CREADR`.

- `GITS_IIDR`. The Revision field encodes the table layout ABI revision. In the future we might implement direct injection of virtual LPIs. This will require an upgrade of the table layout and an evolution of the ABI. `GITS_IIDR` must be restored before calling `KVM_DEV_ARM_ITS_RESTORE_TABLES`.

For other registers, getting or setting a register has the same effect as reading/writing the register on real hardware.

Errors:

<code>-ENXIO</code>	Offset does not correspond to any supported register
<code>-EFAULT</code>	Invalid user pointer for <code>attr-&gt;addr</code>
<code>-EINVAL</code>	Offset is not 64-bit aligned
<code>-EBUSY</code>	one or more VCPUS are running

### ITS Restore Sequence:

The following ordering must be followed when restoring the GIC and the ITS:

- restore all guest memory and create vcpus
- restore all redistributors
- provide the ITS base address (`KVM_DEV_ARM_VGIC_GRP_ADDR`)
- restore the ITS in the following order:
  - Restore `GITS_CBASER`
  - Restore all other `GITS_` registers, except `GITS_CTLR`!
  - Load the ITS table data (`KVM_DEV_ARM_ITS_RESTORE_TABLES`)
  - Restore `GITS_CTLR`

Then vcpus can be started.

### ITS Table ABI REV0:

Revision 0 of the ABI only supports the features of a virtual GICv3, and does not support a virtual GICv4 with support for direct injection of virtual interrupts for nested hypervisors.

The device table and ITT are indexed by the DeviceID and EventID, respectively. The collection table is not indexed by CollectionID, and the entries in the collection are listed in no particular order. All entries are 8 bytes.

Device Table Entry (DTE):

bits:	63	62 ... 49	48 ... 5	4 ... 0
values:	V	next	ITT_addr	Size

where:

- V indicates whether the entry is valid. If not, other fields are not meaningful.
- next: equals to 0 if this entry is the last one; otherwise it corresponds to the DeviceID offset to the next DTE, capped by  $2^{14} - 1$ .
- ITT\_addr matches bits [51:8] of the ITT address (256 Byte aligned).
- Size specifies the supported number of bits for the EventID, minus one

Collection Table Entry (CTE):

bits:	63	62 .. 52	51 ... 16	15 ... 0
values:	V	RES0	RDBase	ICID

where:

- V indicates whether the entry is valid. If not, other fields are not meaningful.
- RES0: reserved field with Should-Be-Zero-or-Preserved behavior.
- RDBase is the PE number (GICR\_TYPER.Processor\_Number semantic),
- ICID is the collection ID

Interrupt Translation Entry (ITE):

bits:	63 ... 48	47 ... 16	15 ... 0
values:	next	pINTID	ICID

where:

- next: equals to 0 if this entry is the last one; otherwise it corresponds to the EventID offset to the next ITE capped by  $2^{16} - 1$ .
- pINTID is the physical LPI ID; if zero, it means the entry is not valid and other fields are not meaningful.
- ICID is the collection ID

**ITS Reset State:**

RESET returns the ITS to the same state that it was when first created and initialized. When the RESET command returns, the following things are guaranteed:

- The ITS is not enabled and quiescent `GITS_CTLR.Enabled = 0` .Quiescent=1
- There is no internally cached state
- No collection or device table are used `GITS_BASER<n>.Valid = 0`
- `GITS_CBASER = 0`, `GITS_CREADR = 0`, `GITS_CWRITER = 0`
- The ABI version is unchanged and remains the one set when the ITS device was first created.

**1.18.2 ARM Virtual Generic Interrupt Controller v2 (VGIC)**

Device types supported:

- `KVM_DEV_TYPE_ARM_VGIC_V2` ARM Generic Interrupt Controller v2.0

Only one VGIC instance may be instantiated through either this API or the legacy `KVM_CREATE_IRQCHIP` API. The created VGIC will act as the VM interrupt controller, requiring emulated user-space devices to inject interrupts to the VGIC instead of directly to CPUs.

GICv3 implementations with hardware compatibility support allow creating a guest GICv2 through this interface. For information on creating a guest GICv3 device and guest ITS devices, see `arm-vgic-v3.txt`. It is not possible to create both a GICv3 and GICv2 device on the same VM.

**Groups:****`KVM_DEV_ARM_VGIC_GRP_ADDR`**

Attributes:

**`KVM_VGIC_V2_ADDR_TYPE_DIST` (rw, 64-bit)**

Base address in the guest physical address space of the GIC distributor register mappings. Only valid for `KVM_DEV_TYPE_ARM_VGIC_V2`. This address needs to be 4K aligned and the region covers 4 KByte.

**`KVM_VGIC_V2_ADDR_TYPE_CPU` (rw, 64-bit)**

Base address in the guest physical address space of the GIC virtual cpu interface register mappings. Only valid for `KVM_DEV_TYPE_ARM_VGIC_V2`. This address needs to be 4K aligned and the region covers 4 KByte.

Errors:

-	Address outside of addressable IPA range
E2BIG	
-	Incorrectly aligned address
EINVA	
-	Address already configured
EEXIS	
-	The group or attribute is unknown/unsupported for this device or hardware support is missing.
ENXIC	
-	Invalid user pointer for attr->addr.
EFAUI	

## KVM\_DEV\_ARM\_VGIC\_GRP\_DIST\_REGS

Attributes:

The attr field of `kvm_device_attr` encodes two values:

bits:	63	....	40	39 .. 32	31	....	┐
→ 0							
values:		reserved		vcpu_index		offset	┐
→							

All distributor regs are (rw, 32-bit)

The offset is relative to the “Distributor base address” as defined in the GICv2 specs. Getting or setting such a register has the same effect as reading or writing the register on the actual hardware from the cpu whose index is specified with the `vcpu_index` field. Note that most distributor fields are not banked, but return the same value regardless of the `vcpu_index` used to access the register.

`GICD_IIDR.Revision` is updated when the KVM implementation of an emulated GICv2 is changed in a way directly observable by the guest or userspace. Userspace should read `GICD_IIDR` from KVM and write back the read value to confirm its expected behavior is aligned with the KVM implementation. Userspace should set `GICD_IIDR` before setting any other registers (both `KVM_DEV_ARM_VGIC_GRP_DIST_REGS` and `KVM_DEV_ARM_VGIC_GRP_CPU_REGS`) to ensure the expected behavior. Unless `GICD_IIDR` has been set from userspace, writes to the interrupt group registers (`GICD_IGROUPR`) are ignored.

Errors:

-ENXIO	Getting or setting this register is not yet supported
-EBUSY	One or more VCPUs are running
-EINVAL	Invalid <code>vcpu_index</code> supplied

## KVM\_DEV\_ARM\_VGIC\_GRP\_CPU\_REGS

Attributes:

The attr field of `kvm_device_attr` encodes two values:



bits:	63	....	40	39 .. 32	31	....	↵
↵ 0							
values:		reserved		vcpu_index		offset	↵
↵							

All CPU interface regs are (rw, 32-bit)

The offset specifies the offset from the “CPU interface base address” as defined in the GICv2 specs. Getting or setting such a register has the same effect as reading or writing the register on the actual hardware.

The Active Priorities Registers APR<sub>n</sub> are implementation defined, so we set a fixed format for our implementation that fits with the model of a “GICv2 implementation without the security extensions” which we present to the guest. This interface always exposes four register APR[0-3] describing the maximum possible 128 preemption levels. The semantics of the register indicate if any interrupts in a given preemption level are in the active state by setting the corresponding bit.

Thus, preemption level X has one or more active interrupts if and only if:

$$\text{APR}_n[X \bmod 32] == 0b1, \text{ where } n = X / 32$$

Bits for undefined preemption levels are RAZ/WI.

Note that this differs from a CPU’s view of the APRs on hardware in which a GIC without the security extensions expose group 0 and group 1 active priorities in separate register groups, whereas we show a combined view similar to GICv2’s GICH\_APR.

For historical reasons and to provide ABI compatibility with userspace we export the GICC\_PMR register in the format of the GICH\_VMCR.VMPriMask field in the lower 5 bits of a word, meaning that userspace must always use the lower 5 bits to communicate with the KVM device and must shift the value left by 3 places to obtain the actual priority mask level.

Errors:

-ENXIO	Getting or setting this register is not yet supported
-EBUSY	One or more VCPUs are running
-EINVAL	Invalid vcpu_index supplied

## KVM\_DEV\_ARM\_VGIC\_GRP\_NR\_IRQS

Attributes:

A value describing the number of interrupts (SGI, PPI and SPI) for this GIC instance, ranging from 64 to 1024, in increments of 32.

Errors:

- EINVAI	Value set is out of the expected range
- EBUSY	Value has already be set, or GIC has already been initialized with default values.

**KVM\_DEV\_ARM\_VGIC\_GRP\_CTRL**

Attributes:

**KVM\_DEV\_ARM\_VGIC\_CTRL\_INIT**

request the initialization of the VGIC or ITS, no additional parameter in `kvm_device_attr.addr`.

Errors:

-ENXIO	VGIC not properly configured as required prior to calling this attribute
- ENODEV	no online VCPU
- ENOMEM	memory shortage when allocating vgic internal data

### 1.18.3 ARM Virtual Generic Interrupt Controller v3 and later (VGICv3)

**Device types supported:**

- KVM\_DEV\_TYPE\_ARM\_VGIC\_V3 ARM Generic Interrupt Controller v3.0

Only one VGIC instance may be instantiated through this API. The created VGIC will act as the VM interrupt controller, requiring emulated user-space devices to inject interrupts to the VGIC instead of directly to CPUs. It is not possible to create both a GICv3 and GICv2 on the same VM.

Creating a guest GICv3 device requires a host GICv3 as well.

**Groups:****KVM\_DEV\_ARM\_VGIC\_GRP\_ADDR**

Attributes:

**KVM\_VGIC\_V3\_ADDR\_TYPE\_DIST (rw, 64-bit)**

Base address in the guest physical address space of the GICv3 distributor register mappings. Only valid for KVM\_DEV\_TYPE\_ARM\_VGIC\_V3. This address needs to be 64K aligned and the region covers 64 KByte.

**KVM\_VGIC\_V3\_ADDR\_TYPE\_REDIST (rw, 64-bit)**

Base address in the guest physical address space of the GICv3 redistributor register mappings. There are two 64K pages for each VCPU and all of the redistributor pages are contiguous.

Only valid for KVM\_DEV\_TYPE\_ARM\_VGIC\_V3. This address needs to be 64K aligned.

#### **KVM\_VGIC\_V3\_ADDR\_TYPE\_REDIST\_REGION (rw, 64-bit)**

The attribute data pointed to by `kvm_device_attr.addr` is a `__u64` value:

bits:	63	....	52		51	....	16		15	-
↪ 12	11	-	0							
values:		count			base				↪	
↪ flags		index								

- index encodes the unique redistributor region index
- flags: reserved for future use, currently 0
- base field encodes bits [51:16] of the guest physical base address of the first redistributor in the region.
- count encodes the number of redistributors in the region. Must be greater than 0.

There are two 64K pages for each redistributor in the region and redistributors are laid out contiguously within the region. Regions are filled with redistributors in the index order. The sum of all region count fields must be greater than or equal to the number of VCPUs. Redistributor regions must be registered in the incremental index order, starting from index 0.

The characteristics of a specific redistributor region can be read by presetting the index field in the attr data. Only valid for KVM\_DEV\_TYPE\_ARM\_VGIC\_V3.

It is invalid to mix calls with KVM\_VGIC\_V3\_ADDR\_TYPE\_REDIST and KVM\_VGIC\_V3\_ADDR\_TYPE\_REDIST\_REGION attributes.

Errors:

- E2BIG	Address outside of addressable IPA range
- EINVA	Incorrectly aligned address, bad redistributor region count/index, mixed redistributor region attribute usage
- EEXIST	Address already configured
- ENOENT	Attempt to read the characteristics of a non existing redistributor region
- ENXIO	The group or attribute is unknown/unsupported for this device or hardware support is missing.
- EFAULT	Invalid user pointer for attr->addr.

#### **KVM\_DEV\_ARM\_VGIC\_GRP\_DIST\_REGS, KVM\_DEV\_ARM\_VGIC\_GRP\_REDIST\_REGS**

Attributes:

The `attr` field of `kvm_device_attr` encodes two values:

bits:	63	....	32	31	....	0	
values:			mpidr			offset	

All distributor regs are (rw, 32-bit) and `kvm_device_attr.addr` points to a `__u32` value. 64-bit registers must be accessed by separately accessing the lower and higher word.

Writes to read-only registers are ignored by the kernel.

`KVM_DEV_ARM_VGIC_GRP_DIST_REGS` accesses the main distributor registers. `KVM_DEV_ARM_VGIC_GRP_REDIST_REGS` accesses the redistributor of the CPU specified by the `mpidr`.

The offset is relative to the “[Re]Distributor base address” as defined in the GICv3/4 specs. Getting or setting such a register has the same effect as reading or writing the register on real hardware, except for the following registers: `GICD_STATUSR`, `GICR_STATUSR`, `GICD_ISPENDR`, `GICR_ISPENDR0`, `GICD_ICPENDR`, and `GICR_ICPENDR0`. These registers behave differently when accessed via this interface compared to their architecturally defined behavior to allow software a full view of the VGIC’s internal state.

The `mpidr` field is used to specify which redistributor is accessed. The `mpidr` is ignored for the distributor.

The `mpidr` encoding is based on the affinity information in the architecture defined MPIDR, and the field is encoded as follows:

63	....	56	55	....	48	47	....	40	39	....	32	
		Aff3			Aff2			Aff1			Aff0	

Note that distributor fields are not banked, but return the same value regardless of the `mpidr` used to access the register.

`GICD_IIDR.Revision` is updated when the KVM implementation is changed in a way directly observable by the guest or userspace. Userspace should read `GICD_IIDR` from KVM and write back the read value to confirm its expected behavior is aligned with the KVM implementation. Userspace should set `GICD_IIDR` before setting any other registers to ensure the expected behavior.

The `GICD_STATUSR` and `GICR_STATUSR` registers are architecturally defined such that a write of a clear bit has no effect, whereas a write with a set bit clears that value. To allow userspace to freely set the values of these two registers, setting the attributes with the register offsets for these two registers simply sets the non-reserved bits to the value written.

Accesses (reads and writes) to the `GICD_ISPENDR` register region and `GICR_ISPENDR0` registers get/set the value of the latched pending state for the interrupts.

This is identical to the value returned by a guest read from `ISPENDR` for an edge triggered interrupt, but may differ for level triggered interrupts. For edge triggered interrupts, once an interrupt becomes pending (whether because of an edge detected on the input line or because of a guest write to `ISPENDR`) this state is “latched”, and only cleared when either the interrupt is activated or when the guest writes to `ICPENDR`. A level triggered interrupt may be pending either because the level input is held high by a device, or because of a guest write to the `ISPENDR` register. Only `ISPENDR` writes are latched; if the device lowers the line level then the interrupt is no longer pending unless the guest also wrote to `ISPENDR`, and conversely writes to `ICPENDR` or activations of the interrupt do not clear the pending status if the line level is still being held high. (These rules are documented in the GICv3 specification descriptions of the `ICPENDR` and `ISPENDR` registers.) For a level triggered interrupt the value accessed here is that of the latch which is set by `ISPENDR` and cleared by `ICPENDR` or interrupt activation, whereas the value returned by a guest read from `ISPENDR` is the logical OR of the latch value and the input line level.

Raw access to the latch state is provided to userspace so that it can save and restore the entire GIC internal state (which is defined by the combination of the current input line level and the latch state, and cannot be deduced from purely the line level and the value of the `ISPENDR` registers).

Accesses to `GICD_ICPENDR` register region and `GICR_ICPENDR0` registers have RAZ/WI semantics, meaning that reads always return 0 and writes are always ignored.

Errors:

-ENXIO	Getting or setting this register is not yet supported
-EBUSY	One or more VCPUs are running

## KVM\_DEV\_ARM\_VGIC\_GRP\_CPU\_SYSREGS

Attributes:

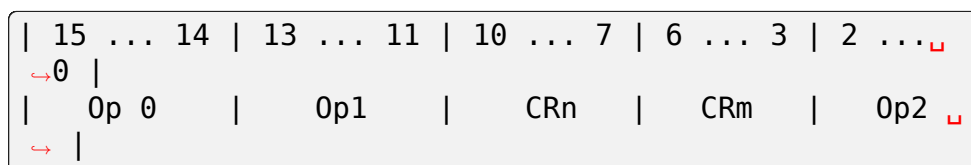
The `attr` field of `kvm_device_attr` encodes two values:

bits:	63	....	32		31	....	16		15	┐
→ ....	0									┐
values:		mpidr				RES				┐
→instr										

The `mpidr` field encodes the CPU ID based on the affinity information in the architecture defined MPIDR, and the field is encoded as follows:

	63	....	56		55	....	48		47	....	40		39	....	32	
		Aff3				Aff2				Aff1				Aff0		

The `instr` field encodes the system register to access based on the fields defined in the A64 instruction set encoding for system register access (RES means the bits are reserved for future use and should be zero):



All system regs accessed through this API are (rw, 64-bit) and `kvm_device_attr.addr` points to a `__u64` value.

`KVM_DEV_ARM_VGIC_GRP_CPU_SYSREGS` accesses the CPU interface registers for the CPU specified by the `mpidr` field.

CPU interface registers access is not implemented for AArch32 mode. Error `-ENXIO` is returned when accessed in AArch32 mode.

Errors:

-ENXIO	Getting or setting this register is not yet supported
-EBUSY	VCPU is running
-EINVAL	Invalid <code>mpidr</code> or register value supplied

## KVM\_DEV\_ARM\_VGIC\_GRP\_NR\_IRQS

Attributes:

A value describing the number of interrupts (SGI, PPI and SPI) for this GIC instance, ranging from 64 to 1024, in increments of 32.

`kvm_device_attr.addr` points to a `__u32` value.

Errors:

-EINVAL	Value set is out of the expected range
-EBUSY	Value has already be set.

## KVM\_DEV\_ARM\_VGIC\_GRP\_CTRL

Attributes:

### KVM\_DEV\_ARM\_VGIC\_CTRL\_INIT

request the initialization of the VGIC, no additional parameter in `kvm_device_attr.addr`.

### KVM\_DEV\_ARM\_VGIC\_SAVE\_PENDING\_TABLES

save all LPI pending bits into guest RAM pending tables.

The first kB of the pending table is not altered by this operation.

Errors:

-ENXIO	VGIC not properly configured as required prior to calling this attribute
-ENODEV	no online VCPU
-ENOMEM	memory shortage when allocating vgic internal data
-EFAULT	Invalid guest ram access
-EBUSY	One or more VCPUS are running

### KVM\_DEV\_ARM\_VGIC\_GRP\_LEVEL\_INFO

Attributes:

The attr field of `kvm_device_attr` encodes the following values:

bits:	63	....	32		31	....	10		↵
↵9	....	0							
values:			mpidr			info			↵
↵	vINTID								

The vINTID specifies which set of IRQs is reported on.

The info field specifies which information userspace wants to get or set using this interface. Currently we support the following info values:

#### VGIC\_LEVEL\_INFO\_LINE\_LEVEL:

Get/Set the input level of the IRQ line for a set of 32 contiguously numbered interrupts.

vINTID must be a multiple of 32.

`kvm_device_attr.addr` points to a `__u32` value which will contain a bitmap where a set bit means the interrupt level is asserted.

Bit[n] indicates the status for interrupt vINTID + n.

SGIs and any interrupt with a higher ID than the number of interrupts supported, will be RAZ/WI. LPIs are always edge-triggered and are therefore not supported by this interface.

PPIs are reported per VCPU as specified in the `mpidr` field, and SPIs are reported with the same value regardless of the `mpidr` specified.

The `mpidr` field encodes the CPU ID based on the affinity information in the architecture defined MPIDR, and the field is encoded as follows:

	63	....	56		55	....	48		47	....	40		39	....	32	
		Aff3				Aff2				Aff1				Aff0		

Errors:

- vINTID is not multiple of 32 or info field is not EINVAI VGIC_LEVEL_INFO_LINE_LEVEL
---

### 1.18.4 MPIC interrupt controller

Device types supported:

- KVM\_DEV\_TYPE\_FSL\_MPIC\_20 Freescale MPIC v2.0
- KVM\_DEV\_TYPE\_FSL\_MPIC\_42 Freescale MPIC v4.2

Only one MPIC instance, of any type, may be instantiated. The created MPIC will act as the system interrupt controller, connecting to each vcpu's interrupt inputs.

#### Groups:

##### **KVM\_DEV\_MPIC\_GRP\_MISC**

Attributes:

##### **KVM\_DEV\_MPIC\_BASE\_ADDR (rw, 64-bit)**

Base address of the 256 KiB MPIC register space. Must be naturally aligned. A value of zero disables the mapping. Reset value is zero.

##### **KVM\_DEV\_MPIC\_GRP\_REGISTER (rw, 32-bit)**

Access an MPIC register, as if the access were made from the guest. "attr" is the byte offset into the MPIC register space. Accesses must be 4-byte aligned.

MSIs may be signaled by using this attribute group to write to the relevant MSIIR.

##### **KVM\_DEV\_MPIC\_GRP\_IRQ\_ACTIVE (rw, 32-bit)**

IRQ input line for each standard openpic source. 0 is inactive and 1 is active, regardless of interrupt sense.

For edge-triggered interrupts: Writing 1 is considered an activating edge, and writing 0 is ignored. Reading returns 1 if a previously signaled edge has not been acknowledged, and 0 otherwise.

"attr" is the IRQ number. IRQ numbers for standard sources are the byte offset of the relevant IVPR from EIVPR0, divided by 32.

#### IRQ Routing:

The MPIC emulation supports IRQ routing. Only a single MPIC device can be instantiated. Once that device has been created, it's available as irqchip id 0.

This irqchip 0 has 256 interrupt pins, which expose the interrupts in the main array of interrupt sources (a.k.a. "SRC" interrupts).

The numbering is the same as the MPIC device tree binding - based on the register offset from the beginning of the sources array, without regard to any subdivisions in chip documentation such as "internal" or "external" interrupts.



Access to non-SRC interrupts is not implemented through IRQ routing mechanisms.

### 1.18.5 FLIC (floating interrupt controller)

FLIC handles floating (non per-cpu) interrupts, i.e. I/O, service and some machine check interruptions. All interrupts are stored in a per-vm list of pending interrupts. FLIC performs operations on this list.

Only one FLIC instance may be instantiated.

FLIC provides support to - add interrupts (`KVM_DEV_FLIC_ENQUEUE`) - inspect currently pending interrupts (`KVM_FLIC_GET_ALL_IRQS`) - purge all pending floating interrupts (`KVM_DEV_FLIC_CLEAR_IRQS`) - purge one pending floating I/O interrupt (`KVM_DEV_FLIC_CLEAR_IO_IRQ`) - enable/disable for the guest transparent async page faults - register and modify adapter interrupt sources (`KVM_DEV_FLIC_ADAPTER_*`) - modify AIS (adapter-interruption-suppression) mode state (`KVM_DEV_FLIC_AISM`) - inject adapter interrupts on a specified adapter (`KVM_DEV_FLIC_AIRQ_INJECT`) - get/set all AIS mode states (`KVM_DEV_FLIC_AISM_ALL`)

#### Groups:

##### **KVM\_DEV\_FLIC\_ENQUEUE**

Passes a buffer and length into the kernel which are then injected into the list of pending interrupts. `attr->addr` contains the pointer to the buffer and `attr->attr` contains the length of the buffer. The format of the data structure `kvm_s390_irq` as it is copied from userspace is defined in `usr/include/linux/kvm.h`.

##### **KVM\_DEV\_FLIC\_GET\_ALL\_IRQS**

Copies all floating interrupts into a buffer provided by userspace. When the buffer is too small it returns `-ENOMEM`, which is the indication for userspace to try again with a bigger buffer.

`-ENOBUFS` is returned when the allocation of a kernelspace buffer has failed.

`-EFAULT` is returned when copying data to userspace failed. All interrupts remain pending, i.e. are not deleted from the list of currently pending interrupts. `attr->addr` contains the userspace address of the buffer into which all interrupt data will be copied. `attr->attr` contains the size of the buffer in bytes.

##### **KVM\_DEV\_FLIC\_CLEAR\_IRQS**

Simply deletes all elements from the list of currently pending floating interrupts. No interrupts are injected into the guest.

##### **KVM\_DEV\_FLIC\_CLEAR\_IO\_IRQ**

Deletes one (if any) I/O interrupt for a subchannel identified by the subsystem identification word passed via the buffer specified by `attr->addr` (address) and `attr->attr` (length).

##### **KVM\_DEV\_FLIC\_APF\_ENABLE**

Enables async page faults for the guest. So in case of a major page fault the host is allowed to handle this async and continues the guest.

**KVM\_DEV\_FLIC\_APF\_DISABLE\_WAIT**

Disables async page faults for the guest and waits until already pending async page faults are done. This is necessary to trigger a completion interrupt for every init interrupt before migrating the interrupt list.

**KVM\_DEV\_FLIC\_ADAPTER\_REGISTER**

Register an I/O adapter interrupt source. Takes a `kvm_s390_io_adapter` describing the adapter to register:

```
struct kvm_s390_io_adapter {
    __u32 id;
    __u8 isc;
    __u8 maskable;
    __u8 swap;
    __u8 flags;
};
```

`id` contains the unique id for the adapter, `isc` the I/O interruption subclass to use, `maskable` whether this adapter may be masked (interrupts turned off), `swap` whether the indicators need to be byte swapped, and `flags` contains further characteristics of the adapter.

Currently defined values for ‘flags’ are:

- **KVM\_S390\_ADAPTER\_SUPPRESSIBLE**: adapter is subject to AIS (adapter-interrupt-suppression) facility. This flag only has an effect if the AIS capability is enabled.

Unknown flag values are ignored.

**KVM\_DEV\_FLIC\_ADAPTER\_MODIFY**

Modifies attributes of an existing I/O adapter interrupt source. Takes a `kvm_s390_io_adapter_req` specifying the adapter and the operation:

```
struct kvm_s390_io_adapter_req {
    __u32 id;
    __u8 type;
    __u8 mask;
    __u16 pad0;
    __u64 addr;
};
```

`id` specifies the adapter and `type` the operation. The supported operations are:

**KVM\_S390\_IO\_ADAPTER\_MASK**

mask or unmask the adapter, as specified in `mask`

**KVM\_S390\_IO\_ADAPTER\_MAP**

This is now a no-op. The mapping is purely done by the irq route.

**KVM\_S390\_IO\_ADAPTER\_UNMAP**

This is now a no-op. The mapping is purely done by the irq route.

**KVM\_DEV\_FLIC\_AISM**

modify the adapter-interrupt-suppression mode for a given `isc` if the

AIS capability is enabled. Takes a `kvm_s390_ais_req` describing:

```
struct kvm_s390_ais_req {
    __u8 isc;
    __u16 mode;
};
```

`isc` contains the target I/O interruption subclass, `mode` the target adapter-interruption-suppression mode. The following modes are currently supported:

- `KVM_S390_AIS_MODE_ALL`: ALL-Interruptions Mode, i.e. `airq` injection is always allowed;
- `KVM_S390_AIS_MODE_SINGLE`: SINGLE-Interruption Mode, i.e. `airq` injection is only allowed once and the following adapter interrupts will be suppressed until the mode is set again to ALL-Interruptions or SINGLE-Interruption mode.

### KVM\_DEV\_FLIC\_AIRQ\_INJECT

Inject adapter interrupts on a specified adapter. `attr->attr` contains the unique id for the adapter, which allows for adapter-specific checks and actions. For adapters subject to AIS, handle the `airq` injection suppression for an `isc` according to the adapter-interruption-suppression mode on condition that the AIS capability is enabled.

### KVM\_DEV\_FLIC\_AISM\_ALL

Gets or sets the adapter-interruption-suppression mode for all ISCs. Takes a `kvm_s390_ais_all` describing:

```
struct kvm_s390_ais_all {
    __u8 simm; /* Single-Interruption-Mode mask */
    __u8 nimm; /* No-Interruption-Mode mask */
};
```

`simm` contains Single-Interruption-Mode mask for all ISCs, `nimm` contains No-Interruption-Mode mask for all ISCs. Each bit in `simm` and `nimm` corresponds to an ISC (MSB0 bit 0 to ISC 0 and so on). The combination of `simm` bit and `nimm` bit presents AIS mode for a ISC.

`KVM_DEV_FLIC_AISM_ALL` is indicated by `KVM_CAP_S390_AIS_MIGRATION`.

Note: The `KVM_SET_DEVICE_ATTR/KVM_GET_DEVICE_ATTR` device ioctls executed on FLIC with an unknown group or attribute gives the error code `EINVAL` (instead of `ENXIO`, as specified in the API documentation). It is not possible to conclude that a FLIC operation is unavailable based on the error code resulting from a usage attempt.

---

**Note:** The `KVM_DEV_FLIC_CLEAR_IO_IRQ` ioctl will return `EINVAL` in case a zero `schid` is specified.

---

### 1.18.6 Generic vcpu interface

The virtual cpu “device” also accepts the ioctls `KVM_SET_DEVICE_ATTR`, `KVM_GET_DEVICE_ATTR`, and `KVM_HAS_DEVICE_ATTR`. The interface uses the same struct `kvm_device_attr` as other devices, but targets VCPU-wide settings and controls.

The groups and attributes per virtual cpu, if any, are architecture specific.

## 1. GROUP: `KVM_ARM_VCPU_PMU_V3_CTRL`

### Architectures

ARM64

### 1.1. ATTRIBUTE: `KVM_ARM_VCPU_PMU_V3_IRQ`

#### Parameters

in `kvm_device_attr.addr` the address for PMU overflow interrupt is a pointer to an int

Returns:

-	The PMU overflow interrupt is already set
<code>EBUSY</code>	
-	Error reading interrupt number
<code>EFAUL</code>	
-	PMUv3 not supported or the overflow interrupt not set when attempting to get it
<code>ENXIC</code>	
-	<code>KVM_ARM_VCPU_PMU_V3</code> feature missing from VCPU
<code>ENOD</code>	
-	Invalid PMU overflow interrupt number supplied or trying to set the IRQ number without using an in-kernel irqchip.
<code>EINVA</code>	

A value describing the PMUv3 (Performance Monitor Unit v3) overflow interrupt number for this vcpu. This interrupt could be a PPI or SPI, but the interrupt type must be same for each vcpu. As a PPI, the interrupt number is the same for all vcups, while as an SPI it must be a separate number per vcpu.

### 1.2 ATTRIBUTE: `KVM_ARM_VCPU_PMU_V3_INIT`

#### Parameters

no additional parameter in `kvm_device_attr.addr`

Returns:

- EEXIST	Interrupt number already used
- ENODEV	PMUv3 not supported or GIC not initialized
- ENXIO	PMUv3 not supported, missing VCPU feature or interrupt number not set
- EBUSY	PMUv3 already initialized

Request the initialization of the PMUv3. If using the PMUv3 with an in-kernel virtual GIC implementation, this must be done after initializing the in-kernel irqchip.

### 1.3 ATTRIBUTE: KVM\_ARM\_VCPU\_PMU\_V3\_FILTER

#### Parameters

in `kvm_device_attr.addr` the address for a PMU event filter is a pointer to a struct `kvm_pmu_event_filter`

#### Returns

- ENOD	PMUv3 not supported or GIC not initialized
- ENXIC	PMUv3 not properly configured or in-kernel irqchip not configured as required prior to calling this attribute
- EBUSY	PMUv3 already initialized
- EINVA	Invalid filter range

Request the installation of a PMU event filter described as follows:

```
struct kvm_pmu_event_filter {
    __u16    base_event;
    __u16    nevents;

#define KVM_PMU_EVENT_ALLOW 0
#define KVM_PMU_EVENT_DENY 1

    __u8     action;
    __u8     pad[3];
};
```

A filter range is defined as the range `[@base_event, @base_event + @nevents)`, together with an `@action` (`KVM_PMU_EVENT_ALLOW` or `KVM_PMU_EVENT_DENY`). The first registered range defines the global policy (global `ALLOW` if the first `@action` is `DENY`, global `DENY` if the first `@action` is `ALLOW`). Multiple ranges can be programmed, and must fit within the event space defined by the PMU architecture (10 bits on ARMv8.0, 16 bits from ARMv8.1 onwards).

Note: “Cancelling” a filter by registering the opposite action for the same range doesn’ t change the default action. For example, installing an ALLOW filter for event range [0:10) as the first filter and then applying a DENY action for the same range will leave the whole range as disabled.

Restrictions: Event 0 (SW\_INCR) is never filtered, as it doesn’ t count a hardware event. Filtering event 0x1E (CHAIN) has no effect either, as it isn’ t strictly speaking an event. Filtering the cycle counter is possible using event 0x11 (CPU\_CYCLES).

## 2. GROUP: KVM\_ARM\_VCPU\_TIMER\_CTRL

### Architectures

ARM, ARM64

### 2.1. ATTRIBUTES: KVM\_ARM\_VCPU\_TIMER\_IRQ\_VTIMER, KVM\_ARM\_VCPU\_TIMER\_IRQ\_PTIMER

#### Parameters

in `kvm_device_attr.addr` the address for the timer interrupt is a pointer to an int

Returns:

-EINVAL	Invalid timer interrupt number
-EBUSY	One or more VCPUs has already run

A value describing the architected timer interrupt number when connected to an in-kernel virtual GIC. These must be a PPI ( $16 \leq \text{intid} < 32$ ). Setting the attribute overrides the default values (see below).

KVM_ARM_VCPU_TIMER_IRQ_VTIME]	The EL1 virtual timer intid (default: 27)
KVM_ARM_VCPU_TIMER_IRQ_PTIME]	The EL1 physical timer intid (default: 30)

Setting the same PPI for different timers will prevent the VCPUs from running. Setting the interrupt number on a VCPU configures all VCPUs created at that time to use the number provided for a given timer, overwriting any previously configured values on other VCPUs. Userspace should configure the interrupt numbers on at least one VCPU after creating all VCPUs and before running any VCPUs.

### 3. GROUP: KVM\_ARM\_VCPU\_PVTIME\_CTRL

#### Architectures

ARM64

#### 3.1 ATTRIBUTE: KVM\_ARM\_VCPU\_PVTIME\_IPA

##### Parameters

64-bit base address

Returns:

-ENXIO	Stolen time not implemented
-EEXIST	Base address already set for this VCPU
-EINVAL	Base address not 64 byte aligned

Specifies the base address of the stolen time structure for this VCPU. The base address must be 64 byte aligned and exist within a valid guest memory region. See *Paravirtualized time support for arm64* for more information including the layout of the stolen time structure.

#### 1.18.7 VFIO virtual device

Device types supported:

- KVM\_DEV\_TYPE\_VFIO

Only one VFIO instance may be created per VM. The created device tracks VFIO groups in use by the VM and features of those groups important to the correctness and acceleration of the VM. As groups are enabled and disabled for use by the VM, KVM should be updated about their presence. When registered with KVM, a reference to the VFIO-group is held by KVM.

##### Groups:

KVM\_DEV\_VFIO\_GROUP

##### KVM\_DEV\_VFIO\_GROUP attributes:

**KVM\_DEV\_VFIO\_GROUP\_ADD: Add a VFIO group to VFIO-KVM device tracking**

kvm\_device\_attr.addr points to an int32\_t file descriptor for the VFIO group.

**KVM\_DEV\_VFIO\_GROUP\_DEL: Remove a VFIO group from VFIO-KVM device tracking**

kvm\_device\_attr.addr points to an int32\_t file descriptor for the VFIO group.

**KVM\_DEV\_VFIO\_GROUP\_SET\_SPAPR\_TCE: attaches a guest visible TCE table**

allocated by sPAPR KVM. kvm\_device\_attr.addr points to a struct:

```
struct kvm_vfio_spapr_tce {
    __s32    groupfd;
    __s32    tablefd;
};
```

where:

- @groupfd is a file descriptor for a VFIO group;
- @tablefd is a file descriptor for a TCE table allocated via KVM\_CREATE\_SPAPR\_TCE.

### **1.18.8 Generic vm interface**

The virtual machine “device” also accepts the ioctls KVM\_SET\_DEVICE\_ATTR, KVM\_GET\_DEVICE\_ATTR, and KVM\_HAS\_DEVICE\_ATTR. The interface uses the same struct `kvm_device_attr` as other devices, but targets VM-wide settings and controls.

The groups and attributes per virtual machine, if any, are architecture specific.

## **1. GROUP: KVM\_S390\_VM\_MEM\_CTRL**

### **Architectures**

s390

### **1.1. ATTRIBUTE: KVM\_S390\_VM\_MEM\_ENABLE\_CMMA**

#### **Parameters**

none

#### **Returns**

-EBUSY if a vcpu is already defined, otherwise 0

Enables Collaborative Memory Management Assist (CMMA) for the virtual machine.

### **1.2. ATTRIBUTE: KVM\_S390\_VM\_MEM\_CLR\_CMMA**

#### **Parameters**

none

#### **Returns**

-EINVAL if CMMA was not enabled; 0 otherwise

Clear the CMMA status for all guest pages, so any pages the guest marked as unused are again used any may not be reclaimed by the host.



### 1.3. ATTRIBUTE KVM\_S390\_VM\_MEM\_LIMIT\_SIZE

#### Parameters

in attr->addr the address for the new limit of guest memory

#### Returns

-EFAULT if the given address is not accessible; -EINVAL if the virtual machine is of type UCONTROL; -E2BIG if the given guest memory is to big for that machine; -EBUSY if a vcpu is already defined; -ENOMEM if not enough memory is available for a new shadow guest mapping; 0 otherwise.

Allows userspace to query the actual limit and set a new limit for the maximum guest memory size. The limit will be rounded up to 2048 MB, 4096 GB, 8192 TB respectively, as this limit is governed by the number of page table levels. In the case that there is no limit we will set the limit to KVM\_S390\_NO\_MEM\_LIMIT (U64\_MAX).

## 2. GROUP: KVM\_S390\_VM\_CPU\_MODEL

#### Architectures

s390

### 2.1. ATTRIBUTE: KVM\_S390\_VM\_CPU\_MACHINE (r/o)

Allows user space to retrieve machine and kvm specific cpu related information:

```
struct kvm_s390_vm_cpu_machine {
    __u64 cpuid;           # CPUID of host
    __u32 ibc;             # IBC level range offered by host
    __u8  pad[4];
    __u64 fac_mask[256];   # set of cpu facilities enabled by KVM
    __u64 fac_list[256];   # set of cpu facilities offered by host
}
```

#### Parameters

address of buffer to store the machine related cpu data of type struct kvm\_s390\_vm\_cpu\_machine\*

#### Returns

-EFAULT if the given address is not accessible from kernel space; -ENOMEM if not enough memory is available to process the ioctl; 0 in case of success.

## 2.2. ATTRIBUTE: KVM\_S390\_VM\_CPU\_PROCESSOR (r/w)

Allows user space to retrieve or request to change cpu related information for a vcpu:

```
struct kvm_s390_vm_cpu_processor {
    __u64 cpuid;          # CPUID currently (to be) used by this
    ↪vcpu
    __u16 ibc;            # IBC level currently (to be) used by
    ↪this vcpu
    __u8  pad[6];
    __u64 fac_list[256];  # set of cpu facilities currently (to
    ↪be) used
                        # by this vcpu
}
```

KVM does not enforce or limit the cpu model data in any form. Take the information retrieved by means of KVM\_S390\_VM\_CPU\_MACHINE as hint for reasonable configuration setups. Instruction interceptions triggered by additionally set facility bits that are not handled by KVM need to be implemented in the VM driver code.

### Parameters

address of buffer to store/set the processor related cpu data of type struct kvm\_s390\_vm\_cpu\_processor\*.

### Returns

-EBUSY in case 1 or more vcpus are already activated (only in write case); -EFAULT if the given address is not accessible from kernel space; -ENOMEM if not enough memory is available to process the ioctl; 0 in case of success.

## 2.3. ATTRIBUTE: KVM\_S390\_VM\_CPU\_MACHINE\_FEAT (r/o)

Allows user space to retrieve available cpu features. A feature is available if provided by the hardware and supported by kvm. In theory, cpu features could even be completely emulated by kvm.

```
struct kvm_s390_vm_cpu_feat {
    __u64 feat[16]; # Bitmap (1 = feature available), MSB 0 bit
    ↪numbering
};
```

### Parameters

address of a buffer to load the feature list from.

### Returns

-EFAULT if the given address is not accessible from kernel space; 0 in case of success.

## 2.4. ATTRIBUTE: KVM\_S390\_VM\_CPU\_PROCESSOR\_FEAT (r/w)

Allows user space to retrieve or change enabled cpu features for all VCPUs of a VM. Features that are not available cannot be enabled.

See [2.3. ATTRIBUTE: KVM\\_S390\\_VM\\_CPU\\_MACHINE\\_FEAT \(r/o\)](#) for a description of the parameter struct.

### Parameters

address of a buffer to store/load the feature list from.

### Returns

-EFAULT if the given address is not accessible from kernel space;  
 -EINVAL if a cpu feature that is not available is to be enabled; -  
 EBUSY if at least one VCPU has already been defined; 0 in case of  
 success.

## 2.5. ATTRIBUTE: KVM\_S390\_VM\_CPU\_MACHINE\_SUBFUNC (r/o)

Allows user space to retrieve available cpu subfunctions without any filtering done by a set IBC. These subfunctions are indicated to the guest VCPU via query or “test bit” subfunctions and used e.g. by cpacf functions, plo and ptff.

A subfunction block is only valid if KVM\_S390\_VM\_CPU\_MACHINE contains the STFL(E) bit introducing the affected instruction. If the affected instruction indicates subfunctions via a “query subfunction”, the response block is contained in the returned struct. If the affected instruction indicates subfunctions via a “test bit” mechanism, the subfunction codes are contained in the returned struct in MSB 0 bit numbering.

```
struct kvm_s390_vm_cpu_subfunc {
    u8 plo[32];           # always valid (ESA/390 feature)
    u8 ptff[16];          # valid with TOD-clock steering
    u8 kmac[16];          # valid with Message-Security-Assist
    u8 kmc[16];           # valid with Message-Security-Assist
    u8 km[16];            # valid with Message-Security-Assist
    u8 kind[16];          # valid with Message-Security-Assist
    u8 klmd[16];          # valid with Message-Security-Assist
    u8 pckmo[16];         # valid with Message-Security-Assist-
    ↪Extension 3
    u8 kmctr[16];         # valid with Message-Security-Assist-
    ↪Extension 4
    u8 kmf[16];           # valid with Message-Security-Assist-
    ↪Extension 4
    u8 kmo[16];           # valid with Message-Security-Assist-
    ↪Extension 4
    u8 pcc[16];           # valid with Message-Security-Assist-
    ↪Extension 4
    u8 ppno[16];          # valid with Message-Security-Assist-
    ↪Extension 5
    u8 kma[16];           # valid with Message-Security-Assist-
```

(continues on next page)

(continued from previous page)

```
↪Extension 8
    u8 kdsa[16];          # valid with Message-Security-Assist-
↪Extension 9
    u8 reserved[1792];    # reserved for future instructions
};
```

**Parameters**

address of a buffer to load the subfunction blocks from.

**Returns**

-EFAULT if the given address is not accessible from kernel space; 0 in case of success.

## 2.6. ATTRIBUTE: KVM\_S390\_VM\_CPU\_PROCESSOR\_SUBFUNC (r/w)

Allows user space to retrieve or change cpu subfunctions to be indicated for all VCPUs of a VM. This attribute will only be available if kernel and hardware support are in place.

The kernel uses the configured subfunction blocks for indication to the guest. A subfunction block will only be used if the associated STFL(E) bit has not been disabled by user space (so the instruction to be queried is actually available for the guest).

As long as no data has been written, a read will fail. The IBC will be used to determine available subfunctions in this case, this will guarantee backward compatibility.

See [2.5. ATTRIBUTE: KVM\\_S390\\_VM\\_CPU\\_MACHINE\\_SUBFUNC \(r/o\)](#) for a description of the parameter struct.

**Parameters**

address of a buffer to store/load the subfunction blocks from.

**Returns**

-EFAULT if the given address is not accessible from kernel space;  
-EINVAL when reading, if there was no write yet; -EBUSY if at least one VCPU has already been defined; 0 in case of success.

## 3. GROUP: KVM\_S390\_VM\_TOD

**Architectures**

s390

### 3.1. ATTRIBUTE: KVM\_S390\_VM\_TOD\_HIGH

Allows user space to set/get the TOD clock extension (u8) (superseded by KVM\_S390\_VM\_TOD\_EXT).

**Parameters**

address of a buffer in user space to store the data (u8) to

**Returns**

- EFAULT if the given address is not accessible from kernel space;
- EINVAL if setting the TOD clock extension to != 0 is not supported
- EOPNOTSUPP for a PV guest (TOD managed by the ultravisor)

### 3.2. ATTRIBUTE: KVM\_S390\_VM\_TOD\_LOW

Allows user space to set/get bits 0-63 of the TOD clock register as defined in the POP (u64).

**Parameters**

address of a buffer in user space to store the data (u64) to

**Returns**

- EFAULT if the given address is not accessible from kernel space
- EOPNOTSUPP for a PV guest (TOD managed by the ultravisor)

### 3.3. ATTRIBUTE: KVM\_S390\_VM\_TOD\_EXT

Allows user space to set/get bits 0-63 of the TOD clock register as defined in the POP (u64). If the guest CPU model supports the TOD clock extension (u8), it also allows user space to get/set it. If the guest CPU model does not support it, it is stored as 0 and not allowed to be set to a value != 0.

**Parameters**

address of a buffer in user space to store the data (kvm\_s390\_vm\_tod\_clock) to

**Returns**

- EFAULT if the given address is not accessible from kernel space;
- EINVAL if setting the TOD clock extension to != 0 is not supported
- EOPNOTSUPP for a PV guest (TOD managed by the ultravisor)

## 4. GROUP: KVM\_S390\_VM\_CRYPTO

**Architectures**

s390

### 4.1. ATTRIBUTE: KVM\_S390\_VM\_CRYPT0\_ENABLE\_AES\_KW (w/o)

Allows user space to enable aes key wrapping, including generating a new wrapping key.

**Parameters**

none

**Returns**

0

### 4.2. ATTRIBUTE: KVM\_S390\_VM\_CRYPT0\_ENABLE\_DEA\_KW (w/o)

Allows user space to enable dea key wrapping, including generating a new wrapping key.

**Parameters**

none

**Returns**

0

### 4.3. ATTRIBUTE: KVM\_S390\_VM\_CRYPT0\_DISABLE\_AES\_KW (w/o)

Allows user space to disable aes key wrapping, clearing the wrapping key.

**Parameters**

none

**Returns**

0

### 4.4. ATTRIBUTE: KVM\_S390\_VM\_CRYPT0\_DISABLE\_DEA\_KW (w/o)

Allows user space to disable dea key wrapping, clearing the wrapping key.

**Parameters**

none

**Returns**

0

## 5. GROUP: KVM\_S390\_VM\_MIGRATION

**Architectures**

s390

### 5.1. ATTRIBUTE: KVM\_S390\_VM\_MIGRATION\_STOP (w/o)

Allows userspace to stop migration mode, needed for PGSTE migration. Setting this attribute when migration mode is not active will have no effects.

**Parameters**

none

**Returns**

0

### 5.2. ATTRIBUTE: KVM\_S390\_VM\_MIGRATION\_START (w/o)

Allows userspace to start migration mode, needed for PGSTE migration. Setting this attribute when migration mode is already active will have no effects.

Dirty tracking must be enabled on all memslots, else -EINVAL is returned. When dirty tracking is disabled on any memslot, migration mode is automatically stopped.

**Parameters**

none

**Returns**

-ENOMEM if there is not enough free memory to start migration mode; -EINVAL if the state of the VM is invalid (e.g. no memory defined); 0 in case of success.

### 5.3. ATTRIBUTE: KVM\_S390\_VM\_MIGRATION\_STATUS (r/o)

Allows userspace to query the status of migration mode.

**Parameters**

address of a buffer in user space to store the data (u64) to; the data itself is either 0 if migration mode is disabled or 1 if it is enabled

**Returns**

-EFAULT if the given address is not accessible from kernel space; 0 in case of success.

## 1.18.9 XICS interrupt controller

Device type supported: KVM\_DEV\_TYPE\_XICS

**Groups:**

1. **KVM\_DEV\_XICS\_GRP\_SOURCES**

Attributes:

One per interrupt source, indexed by the source number.

2. **KVM\_DEV\_XICS\_GRP\_CTRL**

Attributes:

## 2.1 KVM\_DEV\_XICS\_NR\_SERVERS (write only)

The `kvm_device_attr.addr` points to a `__u32` value which is the number of interrupt server numbers (ie, highest possible vcpu id plus one).

Errors:

-EINVAL	Value greater than KVM_MAX_VCPU_ID.
-EFAULT	Invalid user pointer for <code>attr-&gt;addr</code> .
-EBUSY	A vcpu is already connected to the device.

This device emulates the XICS (eXternal Interrupt Controller Specification) defined in PAPR. The XICS has a set of interrupt sources, each identified by a 20-bit source number, and a set of Interrupt Control Presentation (ICP) entities, also called “servers” , each associated with a virtual CPU.

The ICP entities are created by enabling the `KVM_CAP_IRQ_ARCH` capability for each vcpu, specifying `KVM_CAP_IRQ_XICS` in `args[0]` and the interrupt server number (i.e. the vcpu number from the XICS’ s point of view) in `args[1]` of the `kvm_enable_cap` struct. Each ICP has 64 bits of state which can be read and written using the `KVM_GET_ONE_REG` and `KVM_SET_ONE_REG` ioctls on the vcpu. The 64 bit state word has the following bitfields, starting at the least-significant end of the word:

- Unused, 16 bits
- Pending interrupt priority, 8 bits Zero is the highest priority, 255 means no interrupt is pending.
- Pending IPI (inter-processor interrupt) priority, 8 bits Zero is the highest priority, 255 means no IPI is pending.
- Pending interrupt source number, 24 bits Zero means no interrupt pending, 2 means an IPI is pending
- Current processor priority, 8 bits Zero is the highest priority, meaning no interrupts can be delivered, and 255 is the lowest priority.

Each source has 64 bits of state that can be read and written using the `KVM_GET_DEVICE_ATTR` and `KVM_SET_DEVICE_ATTR` ioctls, specifying the `KVM_DEV_XICS_GRP_SOURCES` attribute group, with the attribute number being the interrupt source number. The 64 bit state word has the following bitfields, starting from the least-significant end of the word:

- Destination (server number), 32 bits

This specifies where the interrupt should be sent, and is the interrupt server number specified for the destination vcpu.

- Priority, 8 bits

This is the priority specified for this interrupt source, where 0 is the highest priority and 255 is the lowest. An interrupt with a priority of 255 will never be delivered.

- Level sensitive flag, 1 bit



This bit is 1 for a level-sensitive interrupt source, or 0 for edge-sensitive (or MSI).

- Masked flag, 1 bit

This bit is set to 1 if the interrupt is masked (cannot be delivered regardless of its priority), for example by the `ibm,int-off` RTAS call, or 0 if it is not masked.

- Pending flag, 1 bit

This bit is 1 if the source has a pending interrupt, otherwise 0.

Only one XICS instance may be created per VM.

### 1.18.10 POWER9 eXternal Interrupt Virtualization Engine (XIVE Gen1)

#### Device types supported:

- `KVM_DEV_TYPE_XIVE` POWER9 XIVE Interrupt Controller generation 1

This device acts as a VM interrupt controller. It provides the KVM interface to configure the interrupt sources of a VM in the underlying POWER9 XIVE interrupt controller.

Only one XIVE instance may be instantiated. A guest XIVE device requires a POWER9 host and the guest OS should have support for the XIVE native exploitation interrupt mode. If not, it should run using the legacy interrupt mode, referred as XICS (POWER7/8).

- Device Mappings

The KVM device exposes different MMIO ranges of the XIVE HW which are required for interrupt management. These are exposed to the guest in VMAs populated with a custom VM fault handler.

#### 1. Thread Interrupt Management Area (TIMA)

Each thread has an associated Thread Interrupt Management context composed of a set of registers. These registers let the thread handle priority management and interrupt acknowledgment. The most important are :

- Interrupt Pending Buffer (IPB)
- Current Processor Priority (CPPR)
- Notification Source Register (NSR)

They are exposed to software in four different pages each proposing a view with a different privilege. The first page is for the physical thread context and the second for the hypervisor. Only the third (operating system) and the fourth (user level) are exposed the guest.

#### 2. Event State Buffer (ESB)

Each source is associated with an Event State Buffer (ESB) with either a pair of even/odd pair of pages which provides commands to manage the source: to trigger, to EOI, to turn off the source for instance.

#### 3. Device pass-through

When a device is passed-through into the guest, the source interrupts are from a different HW controller (PHB4) and the ESB pages exposed to the guest should accommodate this change.

The `passthru_irq` helpers, `kvmppc_xive_set_mapped()` and `kvmppc_xive_clr_mapped()` are called when the device HW irqs are mapped into or unmapped from the guest IRQ number space. The KVM device extends these helpers to clear the ESB pages of the guest IRQ number being mapped and then lets the VM fault handler repopulate. The handler will insert the ESB page corresponding to the HW interrupt of the device being passed-through or the initial IPI ESB page if the device has being removed.

The ESB remapping is fully transparent to the guest and the OS device driver. All handling is done within VFIO and the above helpers in KVM-PPC.

- Groups:

## 1. KVM\_DEV\_XIVE\_GRP\_CTRL

Provides global controls on the device

### Attributes:

1.1 KVM\_DEV\_XIVE\_RESET (write only) Resets the interrupt controller configuration for sources and event queues. To be used by `kexec` and `kdump`.

Errors: none

1.2 KVM\_DEV\_XIVE\_EQ\_SYNC (write only) Sync all the sources and queues and mark the EQ pages dirty. This to make sure that a consistent memory state is captured when migrating the VM.

Errors: none

1.3 KVM\_DEV\_XIVE\_NR\_SERVERS (write only) The `kvm_device_attr.addr` points to a `__u32` value which is the number of interrupt server numbers (ie, highest possible vcpu id plus one).

Errors:

-EINVAL	Value greater than KVM_MAX_VCPU_ID.
-EFAULT	Invalid user pointer for <code>attr-&gt;addr</code> .
-EBUSY	A vCPU is already connected to the device.

## 2. KVM\_DEV\_XIVE\_GRP\_SOURCE (write only)

Initializes a new source in the XIVE device and mask it.

### Attributes:

Interrupt source number (64-bit)

The `kvm_device_attr.addr` points to a `__u64` value:

bits:	63	....	2	1	0
values:		unused		level	type

- type: 0:MSI 1:LSI
- level: assertion level in case of an LSI.

Errors:

-E2BIG	Interrupt source number is out of range
-ENOMEM	Could not create a new source block
-EFAULT	Invalid user pointer for attr->addr.
-ENXIO	Could not allocate underlying HW interrupt

### 3. KVM\_DEV\_XIVE\_GRP\_SOURCE\_CONFIG (write only)

Configures source targeting

#### Attributes:

Interrupt source number (64-bit)

The kvm\_device\_attr.addr points to a \_\_u64 value:

bits:	63	....	33	32	31 .. 3	2 .. 0
values:		eisn		mask	server	priority

- priority: 0-7 interrupt priority level
- server: CPU number chosen to handle the interrupt
- mask: mask flag (unused)
- eisn: Effective Interrupt Source Number

Errors:

-	Unknown source number
ENOENT	
-	Not initialized source number
EINVAL	
-	Invalid priority
EINVAL	
-	Invalid CPU number.
EINVAL	
-	Invalid user pointer for attr->addr.
EFAULT	
-	CPU event queues not configured or configuration of
ENXIO	the underlying HW interrupt failed
-	No CPU available to serve interrupt
EBUSY	

### 4. KVM\_DEV\_XIVE\_GRP\_EQ\_CONFIG (read-write)

Configures an event queue of a CPU

#### Attributes:

EQ descriptor identifier (64-bit)

The EQ descriptor identifier is a tuple (server, priority):

bits:	63	....	32	31 .. 3	2 .. 0
values:		unused		server	priority

The `kvm_device_attr.addr` points to:

```
struct kvm_ppc_xive_eq {
    __u32 flags;
    __u32 qshift;
    __u64 qaddr;
    __u32 qtoggle;
    __u32 qindex;
    __u8  pad[40];
};
```

- **flags: queue flags**

**KVM\_XIVE\_EQ\_ALWAYS\_NOTIFY (required)**

forces notification without using the coalescing mechanism provided by the XIVE END ESBs.

- `qshift`: queue size (power of 2)
- `qaddr`: real address of queue
- `qtoggle`: current queue toggle bit
- `qindex`: current queue index
- `pad`: reserved for future use

Errors:

-ENOENT	Invalid CPU number
-EINVAL	Invalid priority
-EINVAL	Invalid flags
-EINVAL	Invalid queue size
-EINVAL	Invalid queue address
-EFAULT	Invalid user pointer for <code>attr-&gt;addr</code> .
-EIO	Configuration of the underlying HW failed

5. **KVM\_DEV\_XIVE\_GRP\_SOURCE\_SYNC (write only)**

Synchronize the source to flush event notifications

**Attributes:**

Interrupt source number (64-bit)

Errors:

-ENOENT	Unknown source number
-EINVAL	Not initialized source number

- VCPU state

The XIVE IC maintains VP interrupt state in an internal structure called the NVT. When a VP is not dispatched on a HW processor thread, this structure can be updated by HW if the VP is the target of an event notification.

It is important for migration to capture the cached IPB from the NVT as it synthesizes the priorities of the pending interrupts. We capture a bit more to

report debug information.

KVM\_REG\_PPC\_VP\_STATE (2 \* 64bits):

bits:	63	....	32		31	....	0	
values:	TIMA word0				TIMA word1			
bits:	127	.....						64
values:	unused							

- Migration:

Saving the state of a VM using the XIVE native exploitation mode should follow a specific sequence. When the VM is stopped :

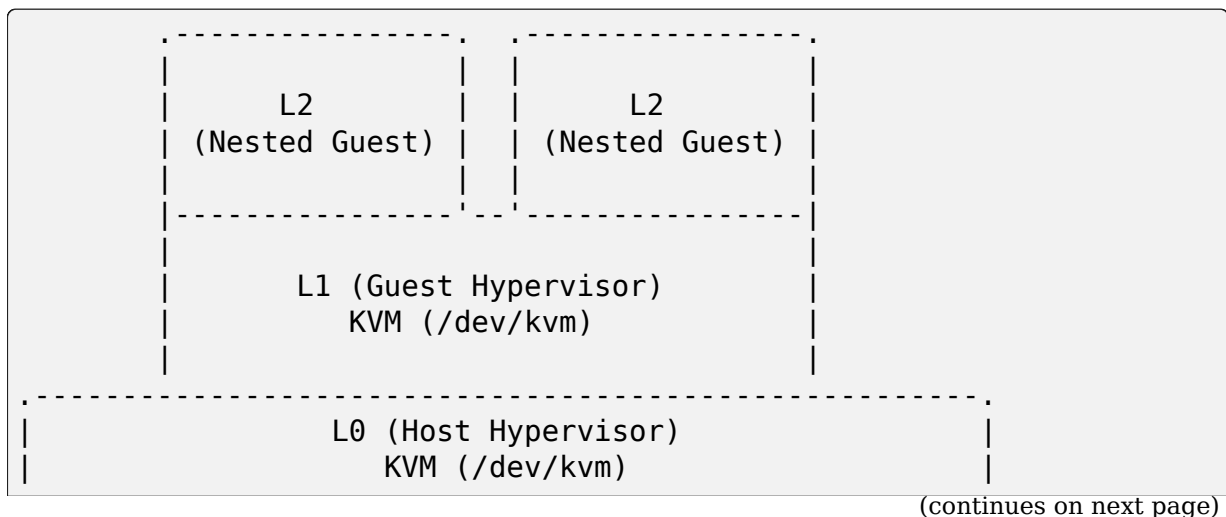
1. Mask all sources (PQ=01) to stop the flow of events.
2. Sync the XIVE device with the KVM control KVM\_DEV\_XIVE\_EQ\_SYNC to flush any in-flight event notification and to stabilize the EQs. At this stage, the EQ pages are marked dirty to make sure they are transferred in the migration sequence.
3. Capture the state of the source targeting, the EQs configuration and the state of thread interrupt context registers.

Restore is similar:

1. Restore the EQ configuration. As targeting depends on it.
2. Restore targeting
3. Restore the thread interrupt contexts
4. Restore the source states
5. Let the vCPU run

## 1.19 Running nested guests with KVM

A nested guest is the ability to run a guest inside another guest (it can be KVM-based or a different hypervisor). The straightforward example is a KVM guest that in turn runs on a KVM guest (the rest of this document is built on this example):



(continued from previous page)

-----  
Hardware (with virtualization extensions)  
-----

Terminology:

- L0 – level-0; the bare metal host, running KVM
- L1 – level-1 guest; a VM running on L0; also called the “guest hypervisor” , as it itself is capable of running KVM.
- L2 – level-2 guest; a VM running on L1, this is the “nested guest”

---

**Note:** The above diagram is modelled after the x86 architecture; s390x, ppc64 and other architectures are likely to have a different design for nesting.

For example, s390x always has an LPAR (LogicalPARTition) hypervisor running on bare metal, adding another layer and resulting in at least four levels in a nested setup –L0 (bare metal, running the LPAR hypervisor), L1 (host hypervisor), L2 (guest hypervisor), L3 (nested guest).

This document will stick with the three-level terminology (L0, L1, and L2) for all architectures; and will largely focus on x86.

---

### 1.19.1 Use Cases

There are several scenarios where nested KVM can be useful, to name a few:

- As a developer, you want to test your software on different operating systems (OSes). Instead of renting multiple VMs from a Cloud Provider, using nested KVM lets you rent a large enough “guest hypervisor” (level-1 guest). This in turn allows you to create multiple nested guests (level-2 guests), running different OSes, on which you can develop and test your software.
- Live migration of “guest hypervisors” and their nested guests, for load balancing, disaster recovery, etc.
- VM image creation tools (e.g. `virt-install`, etc) often run their own VM, and users expect these to work inside a VM.
- Some OSes use virtualization internally for security (e.g. to let applications run safely in isolation).

### 1.19.2 Enabling “nested” (x86)

From Linux kernel v4.19 onwards, the nested KVM parameter is enabled by default for Intel and AMD. (Though your Linux distribution might override this default.)

In case you are running a Linux kernel older than v4.19, to enable nesting, set the nested KVM module parameter to Y or 1. To persist this setting across reboots, you can add it in a config file, as shown below:

1. On the bare metal host (L0), list the kernel modules and ensure that the KVM modules:

```
$ lsmod | grep -i kvm
kvm_intel          133627  0
kvm                435079  1 kvm_intel
```

2. Show information for kvm\_intel module:

```
$ modinfo kvm_intel | grep -i nested
parm:          nested:bool
```

3. For the nested KVM configuration to persist across reboots, place the below in /etc/modprobe.d/kvm\_intel.conf (create the file if it doesn't exist):

```
$ cat /etc/modprobe.d/kvm_intel.conf
options kvm-intel nested=y
```

4. Unload and re-load the KVM Intel module:

```
$ sudo rmmod kvm-intel
$ sudo modprobe kvm-intel
```

5. Verify if the nested parameter for KVM is enabled:

```
$ cat /sys/module/kvm_intel/parameters/nested
Y
```

For AMD hosts, the process is the same as above, except that the module name is kvm-amd.

### 1.19.3 Additional nested-related kernel parameters (x86)

If your hardware is sufficiently advanced (Intel Haswell processor or higher, which has newer hardware virt extensions), the following additional features will also be enabled by default: “Shadow VMCS (Virtual Machine Control Structure)”, APIC Virtualization on your bare metal host (L0). Parameters for Intel hosts:

```
$ cat /sys/module/kvm_intel/parameters/enable_shadow_vmcs
Y

$ cat /sys/module/kvm_intel/parameters/enable_apicv
```

(continues on next page)

(continued from previous page)

```
Y
$ cat /sys/module/kvm_intel/parameters/ept
Y
```

---

**Note:** If you suspect your L2 (i.e. nested guest) is running slower, ensure the above are enabled (particularly `enable_shadow_vmcs` and `ept`).

---

### 1.19.4 Starting a nested guest (x86)

Once your bare metal host (L0) is configured for nesting, you should be able to start an L1 guest with:

```
$ qemu-kvm -cpu host [...]
```

The above will pass through the host CPU’ s capabilities as-is to the gues); or for better live migration compatibility, use a named CPU model supported by QEMU. e.g.:

```
$ qemu-kvm -cpu Haswell-noTSX-IBRS,vmx=on
```

then the guest hypervisor will subsequently be capable of running a nested guest with accelerated KVM.

### 1.19.5 Enabling “nested” (s390x)

1. On the host hypervisor (L0), enable the nested parameter on s390x:

```
$ rmmod kvm
$ modprobe kvm nested=1
```

---

**Note:** On s390x, the kernel parameter `hpage` is mutually exclusive with the nested paramter —i.e. to be able to enable nested, the `hpage` parameter *must* be disabled.

---

2. The guest hypervisor (L1) must be provided with the `sie` CPU feature —with QEMU, this can be done by using “host passthrough” (via the command-line `-cpu host`).
3. Now the KVM module can be loaded in the L1 (guest hypervisor):

```
$ modprobe kvm
```



### 1.19.6 Live migration with nested KVM

Migrating an L1 guest, with a *live* nested guest in it, to another bare metal host, works as of Linux kernel 5.3 and QEMU 4.2.0 for Intel x86 systems, and even on older versions for s390x.

On AMD systems, once an L1 guest has started an L2 guest, the L1 guest should no longer be migrated or saved (refer to QEMU documentation on “*savevm*” / “*loadvm*”) until the L2 guest shuts down. Attempting to migrate or save-and-load an L1 guest while an L2 guest is running will result in undefined behavior. You might see a kernel `BUG!` entry in `dmesg`, a kernel ‘oops’, or an outright kernel panic. Such a migrated or loaded L1 guest can no longer be considered stable or secure, and must be restarted. Migrating an L1 guest merely configured to support nesting, while not actually running L2 guests, is expected to function normally even on AMD systems but may fail once guests are started.

Migrating an L2 guest is always expected to succeed, so all the following scenarios should work even on AMD systems:

- Migrating a nested guest (L2) to another L1 guest on the *same* bare metal host.
- Migrating a nested guest (L2) to another L1 guest on a *different* bare metal host.
- Migrating a nested guest (L2) to a bare metal host.

### 1.19.7 Reporting bugs from nested setups

Debugging “nested” problems can involve sifting through log files across L0, L1 and L2; this can result in tedious back-n-forth between the bug reporter and the bug fixer.

- Mention that you are in a “nested” setup. If you are running any kind of “nesting” at all, say so. Unfortunately, this needs to be called out because when reporting bugs, people tend to forget to even *mention* that they’re using nested virtualization.
- Ensure you are actually running KVM on KVM. Sometimes people do not have KVM enabled for their guest hypervisor (L1), which results in them running with pure emulation or what QEMU calls it as “TCG”, but they think they’re running nested KVM. Thus confusing “nested Virt” (which could also mean, QEMU on KVM) with “nested KVM” (KVM on KVM).

#### Information to collect (generic)

The following is not an exhaustive list, but a very good starting point:

- Kernel, libvirt, and QEMU version from L0
- Kernel, libvirt and QEMU version from L1
- QEMU command-line of L1 – when using libvirt, you’ll find it here: `/var/log/libvirt/qemu/instance.log`

- QEMU command-line of L2 – as above, when using libvirt, get the complete libvirt-generated QEMU command-line
- `cat /sys/cpuinfo` from L0
- `cat /sys/cpuinfo` from L1
- `lscpu` from L0
- `lscpu` from L1
- Full `dmesg` output from L0
- Full `dmesg` output from L1

### **x86-specific info to collect**

Both the below commands, `x86info` and `dmidecode`, should be available on most Linux distributions with the same name:

- Output of: `x86info -a` from L0
- Output of: `x86info -a` from L1
- Output of: `dmidecode` from L0
- Output of: `dmidecode` from L1

### **s390x-specific info to collect**

Along with the earlier mentioned generic details, the below is also recommended:

- `/proc/sysinfo` from L1; this will also include the info from L0

## UML HOWTO

- *Introduction*
  - *How is UML Different from a VM using Virtualization package X?*
  - *Why Would I Want User Mode Linux?*
  - *Why not to run UML*
- *Building a UML instance*
  - *Creating an image*
  - *Edit key system files*
- *Setting Up UML Networking*
  - *Network configuration privileges*
  - *Configuring vector transports*
    - \* *Common options*
    - \* *Shared Options*
    - \* *tap transport*
    - \* *hybrid transport*
    - \* *raw socket transport*
    - \* *GRE socket transport*
    - \* *l2tpv3 socket transport*
    - \* *BESS socket transport*
  - *Configuring Legacy transports*
- *Running UML*
  - *Arguments*
    - \* *Mandatory Arguments:*
    - \* *Important Optional Arguments*
  - *Starting UML*
  - *Logging in*

- *The UML Management Console*
  - \* *version*
  - \* *help*
  - \* *halt and reboot*
  - \* *config*
  - \* *remove*
  - \* *sysrq*
  - \* *cad*
  - \* *stop*
  - \* *go*
  - \* *proc*
  - \* *stack*
- *Advanced UML Topics*
  - *Sharing Filesystems between Virtual Machines*
    - \* *Using layered block devices*
    - \* *Disk Usage*
    - \* *COW validity.*
    - \* *Cows can moo - uml\_moo : Merging a COW file with its backing file*
  - *Host file access*
    - \* *Using hostfs*
    - \* *hostfs as the root filesystem*
    - \* *Hostfs Caveats*
  - *Tuning UML*
- *Contributing to UML and Developing with UML*
  - *Tracing UML*
  - *Kernel debugging*
  - *Developing Device Drivers*
    - \* *Security Considerations*

## 2.1 Introduction

Welcome to User Mode Linux

User Mode Linux is the first Open Source virtualization platform (first release date 1991) and second virtualization platform for an x86 PC.

### 2.1.1 How is UML Different from a VM using Virtualization package X?

We have come to assume that virtualization also means some level of hardware emulation. In fact, it does not. As long as a virtualization package provides the OS with devices which the OS can recognize and has a driver for, the devices do not need to emulate real hardware. Most OSes today have built-in support for a number of “fake” devices used only under virtualization. User Mode Linux takes this concept to the ultimate extreme - there is not a single real device in sight. It is 100% artificial or if we use the correct term 100% paravirtual. All UML devices are abstract concepts which map onto something provided by the host - files, sockets, pipes, etc.

The other major difference between UML and various virtualization packages is that there is a distinct difference between the way the UML kernel and the UML programs operate. The UML kernel is just a process running on Linux - same as any other program. It can be run by an unprivileged user and it does not require anything in terms of special CPU features. The UML userspace, however, is a bit different. The Linux kernel on the host machine assists UML in intercepting everything the program running on a UML instance is trying to do and making the UML kernel handle all of its requests. This is different from other virtualization packages which do not make any difference between the guest kernel and guest programs. This difference results in a number of advantages and disadvantages of UML over let’ s say QEMU which we will cover later in this document.

### 2.1.2 Why Would I Want User Mode Linux?

- If User Mode Linux kernel crashes, your host kernel is still fine. It is not accelerated in any way (vhost, kvm, etc) and it is not trying to access any devices directly. It is, in fact, a process like any other.
- You can run a usermode kernel as a non-root user (you may need to arrange appropriate permissions for some devices).
- You can run a very small VM with a minimal footprint for a specific task (for example 32M or less).
- You can get extremely high performance for anything which is a “kernel specific task” such as forwarding, firewalling, etc while still being isolated from the host kernel.
- You can play with kernel concepts without breaking things.
- You are not bound by “emulating” hardware, so you can try weird and wonderful concepts which are very difficult to support when emulating real hardware

such as time travel and making your system clock dependent on what UML does (very useful for things like tests).

- It's fun.

### 2.1.3 Why not to run UML

- The syscall interception technique used by UML makes it inherently slower for any userspace applications. While it can do kernel tasks on par with most other virtualization packages, its userspace is **slow**. The root cause is that UML has a very high cost of creating new processes and threads (something most Unix/Linux applications take for granted).
- UML is strictly uniprocessor at present. If you want to run an application which needs many CPUs to function, it is clearly the wrong choice.

## 2.2 Building a UML instance

There is no UML installer in any distribution. While you can use off the shelf install media to install into a blank VM using a virtualization package, there is no UML equivalent. You have to use appropriate tools on your host to build a viable filesystem image.

This is extremely easy on Debian - you can do it using debootstrap. It is also easy on OpenWRT - the build process can build UML images. All other distros - YMMV.

### 2.2.1 Creating an image

Create a sparse raw disk image:

```
# dd if=/dev/zero of=disk_image_name bs=1 count=1 seek=16G
```

This will create a 16G disk image. The OS will initially allocate only one block and will allocate more as they are written by UML. As of kernel version 4.19 UML fully supports TRIM (as usually used by flash drives). Using TRIM inside the UML image by specifying discard as a mount option or by running `tune2fs -o discard /dev/ubdXX` will request UML to return any unused blocks to the OS.

Create a filesystem on the disk image and mount it:

```
# mkfs.ext4 ./disk_image_name && mount ./disk_image_name /mnt
```

This example uses ext4, any other filesystem such as ext3, btrfs, xfs, jfs, etc will work too.

Create a minimal OS installation on the mounted filesystem:

```
# debootstrap buster /mnt http://deb.debian.org/debian
```

debootstrap does not set up the root password, fstab, hostname or anything related to networking. It is up to the user to do that.

Set the root password -t he easiest way to do that is to chroot into the mounted image:

```
# chroot /mnt
# passwd
# exit
```

### 2.2.2 Edit key system files

UML block devices are called ubds. The fstab created by debootstrap will be empty and it needs an entry for the root file system:

```
/dev/ubd0    ext4    discard,errors=remount-ro    0    1
```

The image hostname will be set to the same as the host on which you are creating it image. It is a good idea to change that to avoid “Oh, bummer, I rebooted the wrong machine” .

UML supports two classes of network devices - the older `uml_net` ones which are scheduled for obsolescence. These are called `ethX`. It also supports the newer vector IO devices which are significantly faster and have support for some standard virtual network encapsulations like Ethernet over GRE and Ethernet over L2TPv3. These are called `vec0`.

Depending on which one is in use, `/etc/network/interfaces` will need entries like:

```
# legacy UML network devices
auto eth0
iface eth0 inet dhcp

# vector UML network devices
auto vec0
iface eth0 inet dhcp
```

We now have a UML image which is nearly ready to run, all we need is a UML kernel and modules for it.

Most distributions have a UML package. Even if you intend to use your own kernel, testing the image with a stock one is always a good start. These packages come with a set of modules which should be copied to the target filesystem. The location is distribution dependent. For Debian these reside under `/usr/lib/uml/modules`. Copy recursively the content of this directory to the mounted UML filesystem:

```
# cp -rax /usr/lib/uml/modules /mnt/lib/modules
```

If you have compiled your own kernel, you need to use the usual “install modules to a location” procedure by running:

```
# make install MODULES_DIR=/mnt/lib/modules
```

At this point the image is ready to be brought up.

## 2.3 Setting Up UML Networking

UML networking is designed to emulate an Ethernet connection. This connection may be either a point-to-point (similar to a connection between machines using a back-to-back cable) or a connection to a switch. UML supports a wide variety of means to build these connections to all of: local machine, remote machine(s), local and remote UML and other VM instances.

Transport	Type	Capabilities	Throughput
tap	vector	checksum, tso	> 8Gbit
hybrid	vector	checksum, tso, multipacket rx	> 6Gbit
raw	vector	checksum, tso, multipacket rx, tx”	> 6Gbit
EoGRE	vector	multipacket rx, tx	> 3Gbit
Eol2tpv3	vector	multipacket rx, tx	> 3Gbit
bess	vector	multipacket rx, tx	> 3Gbit
fd	vector	dependent on fd type	varies
tuntap	legacy	none	~ 500Mbit
daemon	legacy	none	~ 450Mbit
socket	legacy	none	~ 450Mbit
pcap	legacy	rx only	~ 450Mbit
ethertap	legacy	obsolete	~ 500Mbit
vde	legacy	obsolete	~ 500Mbit

- All transports which have tso and checksum offloads can deliver speeds approaching 10G on TCP streams.
- All transports which have multi-packet rx and/or tx can deliver pps rates of up to 1Mps or more.
- All legacy transports are generally limited to ~600-700MBit and 0.05Mps
- GRE and L2TPv3 allow connections to all of: local machine, remote machines, remote network devices and remote UML instances.
- Socket allows connections only between UML instances.
- Daemon and bess require running a local switch. This switch may be connected to the host as well.

### 2.3.1 Network configuration privileges

The majority of the supported networking modes need root privileges. For example, in the legacy tuntap networking mode, users were required to be part of the group associated with the tunnel device.

For newer network drivers like the vector transports, root privilege is required to fire an ioctl to setup the tun interface and/or use raw sockets where needed.

This can be achieved by granting the user a particular capability instead of running UML as root. In case of vector transport, a user can add the capability CAP\_NET\_ADMIN or CAP\_NET\_RAW, to the uml binary. Thenceforth, UML can be run with normal user privileges, along with full networking.



For example:

```
# sudo setcap cap_net_raw,cap_net_admin+ep linux
```

### 2.3.2 Configuring vector transports

All vector transports support a similar syntax:

If X is the interface number as in vec0, vec1, vec2, etc, the general syntax for options is:

```
vecX:transport="Transport Name",option=value,option=value,...,  
→option=value
```

#### Common options

These options are common for all transports:

- `depth=int` - sets the queue depth for vector IO. This is the amount of packets UML will attempt to read or write in a single system call. The default number is 64 and is generally sufficient for most applications that need throughput in the 2-4 Gbit range. Higher speeds may require larger values.
- `mac=XX:XX:XX:XX:XX` - sets the interface MAC address value.
- `gro=[0,1]` - sets GRO on or off. Enables receive/transmit offloads. The effect of this option depends on the host side support in the transport which is being configured. In most cases it will enable TCP segmentation and RX/TX check-summing offloads. The setting must be identical on the host side and the UML side. The UML kernel will produce warnings if it is not. For example, GRO is enabled by default on local machine interfaces (e.g. veth pairs, bridge, etc), so it should be enabled in UML in the corresponding UML transports (raw, tap, hybrid) in order for networking to operate correctly.
- `mtu=int` - sets the interface MTU
- `headroom=int` - adjusts the default headroom (32 bytes) reserved if a packet will need to be re-encapsulated into for instance VXLAN.
- `vec=0` - disable multipacket io and fall back to packet at a time mode

#### Shared Options

- `ifname=str` Transports which bind to a local network interface have a shared option - the name of the interface to bind to.
- `src, dst, src_port, dst_port` - all transports which use sockets which have the notion of source and destination and/or source port and destination port use these to specify them.
- `v6=[0,1]` to specify if a v6 connection is desired for all transports which operate over IP. Additionally, for transports that have some differences in the way they operate over v4 and v6 (for example EoL2TPv3), sets the correct mode of

operation. In the absense of this option, the socket type is determined based on what do the src and dst arguments resolve/parse to.

### tap transport

Example:

```
vecX:transport=tap,ifname=tap0,depth=128,gro=1
```

This will connect vec0 to tap0 on the host. Tap0 must already exist (for example created using `tunctl`) and UP.

tap0 can be configured as a point-to-point interface and given an ip address so that UML can talk to the host. Alternatively, it is possible to connect UML to a tap interface which is connected to a bridge.

While tap relies on the vector infrastructure, it is not a true vector transport at this point, because Linux does not support multi-packet IO on tap file descriptors for normal userspace apps like UML. This is a privilege which is offered only to something which can hook up to it at kernel level via specialized interfaces like vhost-net. A vhost-net like helper for UML is planned at some point in the future.

Privileges required: tap transport requires either:

- tap interface to exist and be created persistent and owned by the UML user using `tunctl`. Example `tunctl -u uml-user -t tap0`
- binary to have `CAP_NET_ADMIN` privilege

### hybrid transport

Example:

```
vecX:transport=hybrid,ifname=tap0,depth=128,gro=1
```

This is an experimental/demo transport which couples tap for transmit and a raw socket for receive. The raw socket allows multi-packet receive resulting in significantly higher packet rates than normal tap

Privileges required: hybrid requires `CAP_NET_RAW` capability by the UML user as well as the requirements for the tap transport.

### raw socket transport

Example:

```
vecX:transport=raw,ifname=p-veth0,depth=128,gro=1
```

This transport uses vector IO on raw sockets. While you can bind to any interface including a physical one, the most common use it to bind to the “peer” side of a veth pair with the other side configured on the host.

Example host configuration for Debian:

**/etc/network/interfaces:**

```

auto veth0
iface veth0 inet static
    address 192.168.4.1
    netmask 255.255.255.252
    broadcast 192.168.4.3
    pre-up ip link add veth0 type veth peer name p-veth0 && \
        ifconfig p-veth0 up

```

UML can now bind to p-veth0 like this:

```
vec0:transport=raw,ifname=p-veth0,depth=128,gro=1
```

If the UML guest is configured with 192.168.4.2 and netmask 255.255.255.0 it can talk to the host on 192.168.4.1

The raw transport also provides some support for offloading some of the filtering to the host. The two options to control it are:

- `bpffile=`str filename of raw bpf code to be loaded as a socket filter
- `bpfflash=`int 0/1 allow loading of bpf from inside User Mode Linux. This option allows the use of the `ethtool load firmware` command to load bpf code.

In either case the bpf code is loaded into the host kernel. While this is presently limited to legacy bpf syntax (not ebpf), it is still a security risk. It is not recommended to allow this unless the User Mode Linux instance is considered trusted.

Privileges required: raw socket transport requires `CAP_NET_RAW` capability.

**GRE socket transport**

Example:

```
vecX:transport=gre,src=$src_host,dst=$dst_host
```

This will configure an Ethernet over GRE (aka GRETAP or GREIRB) tunnel which will connect the UML instance to a GRE endpoint at host `dst_host`. GRE supports the following additional options:

- `rx_key=`int - GRE 32 bit integer key for rx packets, if set, `txkey` must be set too
- `tx_key=`int - GRE 32 bit integer key for tx packets, if set `rx_key` must be set too
- `sequence=[0,1]` - enable GRE sequence
- `pin_sequence=[0,1]` - pretend that the sequence is always reset on each packet (needed to interoperate with some really broken implementations)
- `v6=[0,1]` - force IPv4 or IPv6 sockets respectively
- GRE checksum is not presently supported

GRE has a number of caveats:

- You can use only one GRE connection per ip address. There is no way to multiplex connections as each GRE tunnel is terminated directly on the UML instance.
- The key is not really a security feature. While it was intended as such it's "security" is laughable. It is, however, a useful feature to ensure that the tunnel is not misconfigured.

An example configuration for a Linux host with a local address of 192.168.128.1 to connect to a UML instance at 192.168.129.1

### **/etc/network/interfaces:**

```
auto gt0
iface gt0 inet static
    address 10.0.0.1
    netmask 255.255.255.0
    broadcast 10.0.0.255
    mtu 1500
    pre-up ip link add gt0 type gretap local 192.168.128.1 \
        remote 192.168.129.1 || true
    down ip link del gt0 || true
```

Additionally, GRE has been tested versus a variety of network equipment.

Privileges required: GRE requires CAP\_NET\_RAW

### **l2tpv3 socket transport**

Warning. L2TPv3 has a "bug". It is the "bug" known as "has more options than GNU ls". While it has some advantages, there are usually easier (and less verbose) ways to connect a UML instance to something. For example, most devices which support L2TPv3 also support GRE.

Example:

```
vec0:transport=l2tpv3,udp=1,src=$src_host,dst=$dst_host,srcport=
↪$src_port,dstport=$dst_port,depth=128,rx_session=0xffffffff,tx_
↪session=0xffff
```

This will configure an Ethernet over L2TPv3 fixed tunnel which will connect the UML instance to a L2TPv3 endpoint at host \$dst\_host using the L2TPv3 UDP flavour and UDP destination port \$dst\_port.

L2TPv3 always requires the following additional options:

- rx\_session=int - l2tpv3 32 bit integer session for rx packets
- tx\_session=int - l2tpv3 32 bit integer session for tx packets

As the tunnel is fixed these are not negotiated and they are preconfigured on both ends.

Additionally, L2TPv3 supports the following optional parameters

- rx\_cookie=int - l2tpv3 32 bit integer cookie for rx packets - same functionality as GRE key, more to prevent misconfiguration than provide actual security

- tx\_cookie=int - l2tpv3 32 bit integer cookie for tx packets
- cookie64=[0,1] - use 64 bit cookies instead of 32 bit.
- counter=[0,1] - enable l2tpv3 counter
- pin\_counter=[0,1] - pretend that the counter is always reset on each packet (needed to interoperate with some really broken implementations)
- v6=[0,1] - force v6 sockets
- udp=[0,1] - use raw sockets (0) or UDP (1) version of the protocol

L2TPv3 has a number of caveats:

- you can use only one connection per ip address in raw mode. There is no way to multiplex connections as each L2TPv3 tunnel is terminated directly on the UML instance. UDP mode can use different ports for this purpose.

Here is an example of how to configure a linux host to connect to UML via L2TPv3:

**/etc/network/interfaces:**

```
auto l2tp1
iface l2tp1 inet static
    address 192.168.126.1
    netmask 255.255.255.0
    broadcast 192.168.126.255
    mtu 1500
    pre-up ip l2tp add tunnel remote 127.0.0.1 \
        local 127.0.0.1 encap udp tunnel_id 2 \
        peer_tunnel_id 2 udp_sport 1706 udp_dport 1707 && \
        ip l2tp add session name l2tp1 tunnel_id 2 \
        session_id 0xffffffff peer_session_id 0xffffffff
    down ip l2tp del session tunnel_id 2 session_id 0xffffffff && \
        ip l2tp del tunnel tunnel_id 2
```

Privileges required: L2TPv3 requires CAP\_NET\_RAW for raw IP mode and no special privileges for the UDP mode.

## BESS socket transport

BESS is a high performance modular network switch.

<https://github.com/NetSys/bess>

It has support for a simple sequential packet socket mode which in the more recent versions is using vector IO for high performance.

Example:

```
vecX:transport=bess,src=$unix_src,dst=$unix_dst
```

This will configure a BESS transport using the unix\_src Unix domain socket address as source and unix\_dst socket address as destination.

For BESS configuration and how to allocate a BESS Unix domain socket port please see the BESS documentation.

<https://github.com/NetSys/bess/wiki/Built-In-Modules-and-Ports>

BESS transport does not require any special privileges.

### 2.3.3 Configuring Legacy transports

Legacy transports are now considered obsolete. Please use the vector versions.

## 2.4 Running UML

This section assumes that either the user-mode-linux package from the distribution or a custom built kernel has been installed on the host.

These add an executable called `linux` to the system. This is the UML kernel. It can be run just like any other executable. It will take most normal linux kernel arguments as command line arguments. Additionally, it will need some UML specific arguments in order to do something useful.

### 2.4.1 Arguments

#### Mandatory Arguments:

- `mem=int[K,M,G]` - amount of memory. By default bytes. It will also accept K, M or G qualifiers.
- `ubdX[s,d,c,t]` = virtual disk specification. This is not really mandatory, but it is likely to be needed in nearly all cases so we can specify a root file system. The simplest possible image specification is the name of the image file for the filesystem (created using one of the methods described in [Creating an image](#))
  - UBD devices support copy on write (COW). The changes are kept in a separate file which can be discarded allowing a rollback to the original pristine image. If COW is desired, the UBD image is specified as: `cow_file, master_image`. Example: `ubd0=Filesystem.cow,Filesystem.img`
  - UBD devices can be set to use synchronous IO. Any writes are immediately flushed to disk. This is done by adding `s` after the `ubdX` specification
  - UBD performs some euristics on devices specified as a single filename to make sure that a COW file has not been specified as the image. To turn them off, use the `d` flag after `ubdX`
  - UBD supports TRIM - asking the Host OS to reclaim any unused blocks in the image. To turn it off, specify the `t` flag after `ubdX`
- `root=` root device - most likely `/dev/ubd0` (this is a Linux filesystem image)

## Important Optional Arguments

If UML is run as “linux” with no extra arguments, it will try to start an xterm for every console configured inside the image (up to 6 in most linux distributions). Each console is started inside an xterm. This makes it nice and easy to use UML on a host with a GUI. It is, however, the wrong approach if UML is to be used as a testing harness or run in a text-only environment.

In order to change this behaviour we need to specify an alternative console and wire it to one of the supported “line” channels. For this we need to map a console to use something different from the default xterm.

Example which will divert console number 1 to stdin/stdout:

```
con1=fd:0,fd:1
```

UML supports a wide variety of serial line channels which are specified using the following syntax

```
conX=channel_type:options[,channel_type:options]
```

If the channel specification contains two parts separated by comma, the first one is input, the second one output.

- The null channel - Discard all input or output. Example `con=null` will set all consoles to null by default.
- The fd channel - use file descriptor numbers for input/out. Example: `con1=fd:0,fd:1`.
- The port channel - listen on tcp port number. Example: `con1=port:4321`
- The pty and pts channels - use system pty/pts.
- The tty channel - bind to an existing system tty. Example: `con1=/dev/tty8` will make UML use the host 8th console (usually unused).
- The xterm channel - this is the default - bring up an xterm on this channel and direct IO to it. Note, that in order for xterm to work, the host must have the UML distribution package installed. This usually contains the port-helper and other utilities needed for UML to communicate with the xterm. Alternatively, these need to be compiled and installed from source. All options applicable to consoles also apply to UML serial lines which are presented as `ttyS` inside UML.

### 2.4.2 Starting UML

We can now run UML.

```
# linux mem=2048M umid=TEST \
  ubd0=Filesystem.img \
  vec0:transport=tap,ifname=tap0,depth=128,gro=1 \
  root=/dev/ubda con=null con0=null,fd:2 con1=fd:0,fd:1
```

This will run an instance with 2048M RAM, try to use the image file called `Filesystem.img` as root. It will connect to the host using `tap0`. All consoles except `con1` will be disabled and console 1 will use standard input/output making it appear in the same terminal it was started.

### 2.4.3 Logging in

If you have not set up a password when generating the image, you will have to shut down the UML instance, mount the image, `chroot` into it and set it - as described in the [Generating an Image](#) section. If the password is already set, you can just log in.

### 2.4.4 The UML Management Console

In addition to managing the image from “the inside” using normal `sysadmin` tools, it is possible to perform a number of low level operations using the UML management console. The UML management console is a low-level interface to the kernel on a running UML instance, somewhat like the i386 `SysRq` interface. Since there is a full-blown operating system under UML, there is much greater flexibility possible than with the `SysRq` mechanism.

There are a number of things you can do with the `mconsole` interface:

- get the kernel version
- add and remove devices
- halt or reboot the machine
- Send `SysRq` commands
- Pause and resume the UML
- Inspect processes running inside UML
- Inspect UML internal `/proc` state

You need the `mconsole` client (`uml_mconsole`) which is a part of the UML tools package available in most Linux distributions.

You also need `CONFIG_MCONSOLE` (under ‘General Setup’) enabled in the UML kernel. When you boot UML, you’ ll see a line like:

```
mconsole initialized on /home/jdike/.uml/umlNJ32yL/mconsole
```

If you specify a unique machine id on the UML command line, i.e. `umid=debian`, you’ ll see this:

```
mconsole initialized on /home/jdike/.uml/debian/mconsole
```

That file is the socket that `uml_mconsole` will use to communicate with UML. Run it with either the `umid` or the full path as its argument:

```
# uml_mconsole debian
```

or



```
# uml_mconsole /home/jdike/.uml/debian/mconsole
```

You'll get a prompt, at which you can run one of these commands:

- version
- help
- halt
- reboot
- config
- remove
- sysrq
- help
- cad
- stop
- go
- proc
- stack

### version

This command takes no arguments. It prints the UML version:

```
(mconsole) version
OK Linux OpenWrt 4.14.106 #0 Tue Mar 19 08:19:41 2019 x86_64
```

There are a couple actual uses for this. It's a simple no-op which can be used to check that a UML is running. It's also a way of sending a device interrupt to the UML. UML mconsole is treated internally as a UML device.

### help

This command takes no arguments. It prints a short help screen with the supported mconsole commands.

### halt and reboot

These commands take no arguments. They shut the machine down immediately, with no syncing of disks and no clean shutdown of userspace. So, they are pretty close to crashing the machine:

```
(mconsole) halt
OK
```

### config

“config” adds a new device to the virtual machine. This is supported by most UML device drivers. It takes one argument, which is the device to add, with the same syntax as the kernel command line:

```
(mconsole) config ubd3=/home/jdike/incoming/roots/root_fs_debian22
```

### remove

“remove” deletes a device from the system. Its argument is just the name of the device to be removed. The device must be idle in whatever sense the driver considers necessary. In the case of the ubd driver, the removed block device must not be mounted, swapped on, or otherwise open, and in the case of the network driver, the device must be down:

```
(mconsole) remove ubd3
```

### sysrq

This command takes one argument, which is a single letter. It calls the generic kernel’s SysRq driver, which does whatever is called for by that argument. See the SysRq documentation in Documentation/admin-guide/sysrq.rst in your favorite kernel tree to see what letters are valid and what they do.

### cad

This invokes the Ctl-Alt-Del action in the running image. What exactly this ends up doing is up to init, systemd, etc. Normally, it reboots the machine.

### stop

This puts the UML in a loop reading mconsole requests until a ‘go’ mconsole command is received. This is very useful as a debugging/snapshotting tool.

### go

This resumes a UML after being paused by a ‘stop’ command. Note that when the UML has resumed, TCP connections may have timed out and if the UML is paused for a long period of time, crond might go a little crazy, running all the jobs it didn’t do earlier.

## **proc**

This takes one argument - the name of a file in /proc which is printed to the mconsole standard output

## **stack**

This takes one argument - the pid number of a process. Its stack is printed to a standard output.

## **2.5 Advanced UML Topics**

### **2.5.1 Sharing Filesystems between Virtual Machines**

Don't attempt to share filesystems simply by booting two UMLs from the same file. That's the same thing as booting two physical machines from a shared disk. It will result in filesystem corruption.

#### **Using layered block devices**

The way to share a filesystem between two virtual machines is to use the copy-on-write (COW) layering capability of the ubd block driver. Any changed blocks are stored in the private COW file, while reads come from either device - the private one if the requested block is valid in it, the shared one if not. Using this scheme, the majority of data which is unchanged is shared between an arbitrary number of virtual machines, each of which has a much smaller file containing the changes that it has made. With a large number of UMLs booting from a large root filesystem, this leads to a huge disk space saving.

Sharing file system data will also help performance, since the host will be able to cache the shared data using a much smaller amount of memory, so UML disk requests will be served from the host's memory rather than its disks. There is a major caveat in doing this on multsocket NUMA machines. On such hardware, running many UML instances with a shared master image and COW changes may cause issues like NMIs from excess of inter-socket traffic.

If you are running UML on high end hardware like this, make sure to bind UML to a set of logical cpus residing on the same socket using the taskset command or have a look at the "tuning" section.

To add a copy-on-write layer to an existing block device file, simply add the name of the COW file to the appropriate ubd switch:

```
ubd0=root_fs_cow,root_fs_debian_22
```

where root\_fs\_cow is the private COW file and root\_fs\_debian\_22 is the existing shared filesystem. The COW file need not exist. If it doesn't, the driver will create and initialize it.

### Disk Usage

UML has TRIM support which will release any unused space in its disk image files to the underlying OS. It is important to use either `ls -ls` or `du` to verify the actual file size.

### COW validity.

Any changes to the master image will invalidate all COW files. If this happens, UML will *NOT* automatically delete any of the COW files and will refuse to boot. In this case the only solution is to either restore the old image (including its last modified timestamp) or remove all COW files which will result in their recreation. Any changes in the COW files will be lost.

### Cows can moo - `uml_moo` : Merging a COW file with its backing file

Depending on how you use UML and COW devices, it may be advisable to merge the changes in the COW file into the backing file every once in a while.

The utility that does this is `uml_moo`. Its usage is:

```
uml_moo COW_file new_backing_file
```

There's no need to specify the backing file since that information is already in the COW file header. If you're paranoid, boot the new merged file, and if you're happy with it, move it over the old backing file.

`uml_moo` creates a new backing file by default as a safety measure. It also has a destructive merge option which will merge the COW file directly into its current backing file. This is really only usable when the backing file only has one COW file associated with it. If there are multiple COWs associated with a backing file, a `-d` merge of one of them will invalidate all of the others. However, it is convenient if you're short of disk space, and it should also be noticeably faster than a non-destructive merge.

`uml_moo` is installed with the UML distribution packages and is available as a part of UML utilities.

### 2.5.2 Host file access

If you want to access files on the host machine from inside UML, you can treat it as a separate machine and either `nfs` mount directories from the host or copy files into the virtual machine with `scp`. However, since UML is running on the host, it can access those files just like any other process and make them available inside the virtual machine without the need to use the network. This is possible with the `hostfs` virtual filesystem. With it, you can mount a host directory into the UML filesystem and access the files contained in it just as you would on the host.

#### *SECURITY WARNING*

`Hostfs` without any parameters to the UML Image will allow the image to mount any part of the host filesystem and write to it. Always confine `hostfs` to a specific

“harmless” directory (for example /var/tmp) if running UML. This is especially important if UML is being run as root.

## Using hostfs

To begin with, make sure that hostfs is available inside the virtual machine with:

```
# cat /proc/filesystems
```

hostfs should be listed. If it's not, either rebuild the kernel with hostfs configured into it or make sure that hostfs is built as a module and available inside the virtual machine, and insmod it.

Now all you need to do is run mount:

```
# mount none /mnt/host -t hostfs
```

will mount the host's / on the virtual machine's /mnt/host. If you don't want to mount the host root directory, then you can specify a subdirectory to mount with the -o switch to mount:

```
# mount none /mnt/home -t hostfs -o /home
```

will mount the host's /home on the virtual machine's /mnt/home.

## hostfs as the root filesystem

It's possible to boot from a directory hierarchy on the host using hostfs rather than using the standard filesystem in a file. To start, you need that hierarchy. The easiest way is to loop mount an existing root\_fs file:

```
# mount root_fs uml_root_dir -o loop
```

You need to change the filesystem type of / in etc/fstab to be 'hostfs', so that line looks like this:

/dev/ubd/0	/	hostfs	defaults	1	1
------------	---	--------	----------	---	---

Then you need to chown to yourself all the files in that directory that are owned by root. This worked for me:

```
# find . -uid 0 -exec chown jdiike {} \;
```

Next, make sure that your UML kernel has hostfs compiled in, not as a module. Then run UML with the boot device pointing at that directory:

```
ubd0=/path/to/uml/root/directory
```

UML should then boot as it does normally.

### Hostfs Caveats

Hostfs does not support keeping track of host filesystem changes on the host (outside UML). As a result, if a file is changed without UML's knowledge, UML will not know about it and its own in-memory cache of the file may be corrupt. While it is possible to fix this, it is not something which is being worked on at present.

### 2.5.3 Tuning UML

UML at present is strictly uniprocessor. It will, however spin up a number of threads to handle various functions.

The UBD driver, SIGIO and the MMU emulation do that. If the system is idle, these threads will be migrated to other processors on a SMP host. This, unfortunately, will usually result in LOWER performance because of all of the cache/memory synchronization traffic between cores. As a result, UML will usually benefit from being pinned on a single CPU especially on a large system. This can result in performance differences of 5 times or higher on some benchmarks.

Similarly, on large multi-node NUMA systems UML will benefit if all of its memory is allocated from the same NUMA node it will run on. The OS will *NOT* do that by default. In order to do that, the sysadmin needs to create a suitable tmpfs ramdisk bound to a particular node and use that as the source for UML RAM allocation by specifying it in the TMP or TEMP environment variables. UML will look at the values of TMPDIR, TMP or TEMP for that. If that fails, it will look for shmfs mounted under /dev/shm. If everything else fails use /tmp/ regardless of the filesystem type used for it:

```
mount -t tmpfs -ompol=bind:X none /mnt/tmpfs-nodeX
TEMP=/mnt/tmpfs-nodeX taskset -cX linux options options options..
```

## 2.6 Contributing to UML and Developing with UML

UML is an excellent platform to develop new Linux kernel concepts - filesystems, devices, virtualization, etc. It provides unrivalled opportunities to create and test them without being constrained to emulating specific hardware.

Example - want to try how linux will work with 4096 “proper” network devices?

Not an issue with UML. At the same time, this is something which is difficult with other virtualization packages - they are constrained by the number of devices allowed on the hardware bus they are trying to emulate (for example 16 on a PCI bus in qemu).

If you have something to contribute such as a patch, a bugfix, a new feature, please send it to [linux-um@lists.infradead.org](mailto:linux-um@lists.infradead.org)

Please follow all standard Linux patch guidelines such as cc-ing relevant maintainers and run `./scripts/checkpatch.pl` on your patch. For more details see `Documentation/process/submitting-patches.rst`

Note - the list does not accept HTML or attachments, all emails must be formatted as plain text.

Developing always goes hand in hand with debugging. First of all, you can always run UML under gdb and there will be a whole section later on on how to do that. That, however, is not the only way to debug a linux kernel. Quite often adding tracing statements and/or using UML specific approaches such as ptracing the UML kernel process are significantly more informative.

### 2.6.1 Tracing UML

When running UML consists of a main kernel thread and a number of helper threads. The ones of interest for tracing are NOT the ones that are already ptraced by UML as a part of its MMU emulation.

These are usually the first three threads visible in a ps display. The one with the lowest PID number and using most CPU is usually the kernel thread. The other threads are the disk (ubd) device helper thread and the sigio helper thread. Running ptrace on this thread usually results in the following picture:

```
host$ strace -p 16566
--- SIGIO {si_signo=SIGIO, si_code=POLL_IN, si_band=65} ---
epoll_wait(4, [{EPOLLIN, {u32=3721159424, u64=3721159424}}], 64, 0)
↳= 1
epoll_wait(4, [], 64, 0) = 0
rt_sigreturn({mask=[PIPE]}) = 16967
ptrace(PTRACE_GETREGS, 16967, NULL, 0xd5f34f38) = 0
ptrace(PTRACE_GETREGSET, 16967, NT_X86_XSTATE, [{iov_
↳base=0xd5f35010, iov_len=832}]) = 0
ptrace(PTRACE_GETSIGINFO, 16967, NULL, {si_signo=SIGTRAP, si_
↳code=0x85, si_pid=16967, si_uid=0}) = 0
ptrace(PTRACE_SETREGS, 16967, NULL, 0xd5f34f38) = 0
ptrace(PTRACE_SETREGSET, 16967, NT_X86_XSTATE, [{iov_
↳base=0xd5f35010, iov_len=2696}]) = 0
ptrace(PTRACE_SYSEMU, 16967, NULL, 0) = 0
--- SIGCHLD {si_signo=SIGCHLD, si_code=CLD_TRAPPED, si_pid=16967,
↳si_uid=0, si_status=SIGTRAP, si_utime=65, si_stime=89} ---
wait4(16967, [{WIFSTOPPED(s) && WSTOPSIG(s) == SIGTRAP | 0x80}],
↳WSTOPPED|__WALL, NULL) = 16967
ptrace(PTRACE_GETREGS, 16967, NULL, 0xd5f34f38) = 0
ptrace(PTRACE_GETREGSET, 16967, NT_X86_XSTATE, [{iov_
↳base=0xd5f35010, iov_len=832}]) = 0
ptrace(PTRACE_GETSIGINFO, 16967, NULL, {si_signo=SIGTRAP, si_
↳code=0x85, si_pid=16967, si_uid=0}) = 0
timer_settime(0, 0, {it_interval={tv_sec=0, tv_nsec=0}, it_value=
↳{tv_sec=0, tv_nsec=2830912}}, NULL) = 0
getpid() = 16566
clock_nanosleep(CLOCK_MONOTONIC, 0, {tv_sec=1, tv_nsec=0}, NULL) = ?
↳ ERESTART_RESTARTBLOCK (Interrupted by signal)
--- SIGALRM {si_signo=SIGALRM, si_code=SI_TIMER, si_timerid=0, si_
↳overrun=0, si_value={int=1631716592, ptr=0x614204f0}} ---
```

(continues on next page)

(continued from previous page)

```
rt_sigreturn({mask=[PIPE]})           = -1 EINTR (Interrupted,
↪system call)
```

This is a typical picture from a mostly idle UML instance

- UML interrupt controller uses epoll - this is UML waiting for IO interrupts:  
`epoll_wait(4, [{EPOLLIN, {u32=3721159424, u64=3721159424}}], 64, 0) = 1`
- The sequence of ptrace calls is part of MMU emulation and running in the UML userspace
- `timer_settime` is part of the UML high res timer subsystem mapping timer requests from inside UML onto the host high resolution timers.
- `clock_nanosleep` is UML going into idle (similar to the way a PC will execute an ACPI idle).

As you can see UML will generate quite a bit of output even in idle. The output can be very informative when observing IO. It shows the actual IO calls, their arguments and return values.

## 2.6.2 Kernel debugging

You can run UML under gdb now, though it will not necessarily agree to be started under it. If you are trying to track a runtime bug, it is much better to attach gdb to a running UML instance and let UML run.

Assuming the same PID number as in the previous example, this would be:

```
# gdb -p 16566
```

This will STOP the UML instance, so you must enter *cont* at the GDB command line to request it to continue. It may be a good idea to make this into a gdb script and pass it to gdb as an argument.

## 2.6.3 Developing Device Drivers

Nearly all UML drivers are monolithic. While it is possible to build a UML driver as a kernel module, that limits the possible functionality to in-kernel only and non-UML specific. The reason for this is that in order to really leverage UML, one needs to write a piece of userspace code which maps driver concepts onto actual userspace host calls.

This forms the so called “user” portion of the driver. While it can reuse a lot of kernel concepts, it is generally just another piece of userspace code. This portion needs some matching “kernel” code which resides inside the UML image and which implements the Linux kernel part.

*Note: There are very few limitations in the way “kernel” and “user” interact.*

UML does not have a strictly defined kernel to host API. It does not try to emulate a specific architecture or bus. UML’s “kernel” and “user” can share memory, code



and interact as needed to implement whatever design the software developer has in mind. The only limitations are purely technical. Due to a lot of functions and variables having the same names, the developer should be careful which includes and libraries they are trying to refer to.

As a result a lot of userspace code consists of simple wrappers. F.e. `os_close_file()` is just a wrapper around `close()` which ensures that the userspace function `close` does not clash with similarly named function(s) in the kernel part.

## **Security Considerations**

Drivers or any new functionality should default to not accepting arbitrary filename, bpf code or other parameters which can affect the host from inside the UML instance. For example, specifying the socket used for IPC communication between a driver and the host at the UML command line is OK security-wise. Allowing it as a loadable module parameter isn't.

If such functionality is desirable for a particular application (e.g. loading BPF “firmware” for raw socket network transports), it should be off by default and should be explicitly turned on as a command line parameter at startup.

Even with this in mind, the level of isolation between UML and the host is relatively weak. If the UML userspace is allowed to load arbitrary kernel drivers, an attacker can use this to break out of UML. Thus, if UML is used in a production application, it is recommended that all modules are loaded at boot and kernel module loading is disabled afterwards.



## PARAVIRT\_OPS

Linux provides support for different hypervisor virtualization technologies. Historically different binary kernels would be required in order to support different hypervisors, this restriction was removed with `pv_ops`. Linux `pv_ops` is a virtualization API which enables support for different hypervisors. It allows each hypervisor to override critical operations and allows a single kernel binary to run on all supported execution environments including native machine - without any hypervisors.

`pv_ops` provides a set of function pointers which represent operations corresponding to low level critical instructions and high level functionalities in various areas. `pv_ops` allows for optimizations at run time by enabling binary patching of the low-ops critical operations at boot time.

`pv_ops` operations are classified into three categories:

- **simple indirect call**

These operations correspond to high level functionality where it is known that the overhead of indirect call isn't very important.

- **indirect call which allows optimization with binary patch**

Usually these operations correspond to low level critical instructions. They are called frequently and are performance critical. The overhead is very important.

- **a set of macros for hand written assembly code**

Hand written assembly codes (.S files) also need paravirtualization because they include sensitive instructions or some of code paths in them are very performance critical.



## **GUEST HALT POLLING**

The `cpuidle_haltpoll` driver, with the `haltpoll` governor, allows the guest vcpus to poll for a specified amount of time before halting.

This provides the following benefits to host side polling:

- 1) The `POLL` flag is set while polling is performed, which allows a remote vCPU to avoid sending an IPI (and the associated cost of handling the IPI) when performing a wakeup.
- 2) The VM-exit cost can be avoided.

The downside of guest side polling is that polling is performed even with other runnable tasks in the host.

The basic logic as follows: A global value, `guest_halt_poll_ns`, is configured by the user, indicating the maximum amount of time polling is allowed. This value is fixed.

Each vcpu has an adjustable `guest_halt_poll_ns` ( “per-cpu `guest_halt_poll_ns`” ), which is adjusted by the algorithm in response to events (explained below).

### **4.1 Module Parameters**

The `haltpoll` governor has 5 tunable module parameters:

- 1) `guest_halt_poll_ns`:

Maximum amount of time, in nanoseconds, that polling is performed before halting.

Default: 200000

- 2) `guest_halt_poll_shrink`:

Division factor used to shrink per-cpu `guest_halt_poll_ns` when wakeup event occurs after the global `guest_halt_poll_ns`.

Default: 2

- 3) `guest_halt_poll_grow`:

Multiplication factor used to grow per-cpu `guest_halt_poll_ns` when event occurs after per-cpu `guest_halt_poll_ns` but before global `guest_halt_poll_ns`.

Default: 2

- 4) `guest_halt_poll_grow_start`:

The per-cpu `guest_halt_poll_ns` eventually reaches zero in case of an idle system. This value sets the initial per-cpu `guest_halt_poll_ns` when growing. This can be increased from 10000, to avoid misses during the initial growth stage:

10k, 20k, 40k, ... (example assumes `guest_halt_poll_grow=2`).

Default: 50000

5) `guest_halt_poll_allow_shrink`:

Bool parameter which allows shrinking. Set to N to avoid it (per-cpu `guest_halt_poll_ns` will remain high once achieves global `guest_halt_poll_ns` value).

Default: Y

The module parameters can be set from the debugfs files in:

```
/sys/module/haltpoll/parameters/
```

## 4.2 Further Notes

- Care should be taken when setting the `guest_halt_poll_ns` parameter as a large value has the potential to drive the cpu usage to 100% on a machine which would be almost entirely idle otherwise.

## **NITRO ENCLAVES**

### **5.1 Overview**

Nitro Enclaves (NE) is a new Amazon Elastic Compute Cloud (EC2) capability that allows customers to carve out isolated compute environments within EC2 instances [1].

For example, an application that processes sensitive data and runs in a VM, can be separated from other applications running in the same VM. This application then runs in a separate VM than the primary VM, namely an enclave.

An enclave runs alongside the VM that spawned it. This setup matches low latency applications needs. The resources that are allocated for the enclave, such as memory and CPUs, are carved out of the primary VM. Each enclave is mapped to a process running in the primary VM, that communicates with the NE driver via an `ioctl` interface.

In this sense, there are two components:

1. An enclave abstraction process - a user space process running in the primary VM guest that uses the provided `ioctl` interface of the NE driver to spawn an enclave VM (that's 2 below).

There is a NE emulated PCI device exposed to the primary VM. The driver for this new PCI device is included in the NE driver.

The `ioctl` logic is mapped to PCI device commands e.g. the `NE_START_ENCLAVE` `ioctl` maps to an enclave start PCI command. The PCI device commands are then translated into actions taken on the hypervisor side; that's the Nitro hypervisor running on the host where the primary VM is running. The Nitro hypervisor is based on core KVM technology.

2. The enclave itself - a VM running on the same host as the primary VM that spawned it. Memory and CPUs are carved out of the primary VM and are dedicated for the enclave VM. An enclave does not have persistent storage attached.

The memory regions carved out of the primary VM and given to an enclave need to be aligned 2 MiB / 1 GiB physically contiguous memory regions (or multiple of this size e.g. 8 MiB). The memory can be allocated e.g. by using `hugetlbfs` from user space [2][3]. The memory size for an enclave needs to be at least 64 MiB. The enclave memory and CPUs need to be from the same NUMA node.

An enclave runs on dedicated cores. CPU 0 and its CPU siblings need to remain available for the primary VM. A CPU pool has to be set for NE purposes by an user

with admin capability. See the `cpu list` section from the kernel documentation [4] for how a CPU pool format looks.

An enclave communicates with the primary VM via a local communication channel, using `virtio-vsock` [5]. The primary VM has `virtio-pci vsock` emulated device, while the enclave VM has a `virtio-mmio vsock` emulated device. The `vsock` device uses `eventfd` for signaling. The enclave VM sees the usual interfaces - local APIC and IOAPIC - to get interrupts from `virtio-vsock` device. The `virtio-mmio` device is placed in memory below the typical 4 GiB.

The application that runs in the enclave needs to be packaged in an enclave image together with the OS ( e.g. kernel, ramdisk, init ) that will run in the enclave VM. The enclave VM has its own kernel and follows the standard Linux boot protocol [6].

The kernel `bzImage`, the kernel command line, the `ramdisk(s)` are part of the Enclave Image Format (EIF); plus an EIF header including metadata such as magic number, eif version, image size and CRC.

Hash values are computed for the entire enclave image (EIF), the kernel and `ramdisk(s)`. That' s used, for example, to check that the enclave image that is loaded in the enclave VM is the one that was intended to be run.

These crypto measurements are included in a signed attestation document generated by the Nitro Hypervisor and further used to prove the identity of the enclave; KMS is an example of service that NE is integrated with and that checks the attestation doc.

The enclave image (EIF) is loaded in the enclave memory at offset 8 MiB. The `init` process in the enclave connects to the `vsock CID` of the primary VM and a predefined port - 9000 - to send a heartbeat value - 0xb7. This mechanism is used to check in the primary VM that the enclave has booted. The CID of the primary VM is 3.

If the enclave VM crashes or gracefully exits, an interrupt event is received by the NE driver. This event is sent further to the user space enclave process running in the primary VM via a poll notification mechanism. Then the user space enclave process can exit.

[1] <https://aws.amazon.com/ec2/nitro/nitro-enclaves/> [2] <https://www.kernel.org/doc/html/latest/admin-guide/mm/hugetlbpage.html> [3] <https://lwn.net/Articles/807108/> [4] <https://www.kernel.org/doc/html/latest/admin-guide/kernel-parameters.html> [5] <https://man7.org/linux/man-pages/man7/vsock.7.html> [6] <https://www.kernel.org/doc/html/latest/x86/boot.html>



## BIBLIOGRAPHY

- [white-paper] [http://amd-dev.wpengine.netdna-cdn.com/wordpress/media/2013/12/AMD\\_Memory\\_Encryption\\_Whitepaper\\_v7-Public.pdf](http://amd-dev.wpengine.netdna-cdn.com/wordpress/media/2013/12/AMD_Memory_Encryption_Whitepaper_v7-Public.pdf)
- [api-spec] [https://support.amd.com/TechDocs/55766\\_SEV-KM\\_API\\_Specification.pdf](https://support.amd.com/TechDocs/55766_SEV-KM_API_Specification.pdf)
- [amd-apm] <https://support.amd.com/TechDocs/24593.pdf> (section 15.34)
- [kvm-forum] [https://www.linux-kvm.org/images/7/74/02x08A-Thomas\\_Lendacky-AMDs\\_Virtualizatoin\\_Memory\\_Encryption\\_Technology.pdf](https://www.linux-kvm.org/images/7/74/02x08A-Thomas_Lendacky-AMDs_Virtualizatoin_Memory_Encryption_Technology.pdf)
- [atomic-ops] Documentation/core-api/atomic\_ops.rst
- [memory-barriers] Documentation/memory-barriers.txt
- [lwn-mb] <https://lwn.net/Articles/573436/>