

Memory Management

Complete virtual memory map with 4-level page tables

Note

- Negative addresses such as "-23 TB" are absolute addresses in bytes, counted down from the top of the 64-bit address space. It's easier to understand the layout when seen both in absolute addresses and in distance-from-top notation.

For example $0xffff900000000000 = -23$ TB, it's 23 TB lower than the top of the 64-bit address space (ffffffff).

Note that as we get closer to the top of the address space, the notation changes from TB to GB and then MB/KB.

- "16M TB" might look weird at first sight, but it's an easier way to visualize size notation than "16 EB", which few will recognize at first sight as 16 exabytes. It also shows it nicely how incredibly large 64-bit address space is.

Start addr	Offset	End addr	Size	VM area description
0000000000000000	0	00007fffffffffff	128 TB	user-space virtual memory, different per mm
0000800000000000	+128 TB	ffff7fffffffffff	~16M TB	... huge, almost 64 bits wide hole of non-canonical virtual memory addresses up to the -128 TB starting offset of kernel mappings.
				Kernel-space virtual memory, shared between all processes
ffff800000000000	-128 TB	ffff87ffffffffff	8 TB	... guard hole, also reserved for hypervisor
ffff880000000000	-120 TB	ffff887ffffffffff	0.5 TB	LDT remap for PTI
ffff888000000000	-119.5 TB	ffffc87ffffffffff	64 TB	direct mapping of all physical memory (page_offset)
ffffc88000000000	-55.5 TB	ffffc8fffffffffff	0.5 TB	... unused hole
ffffc90000000000	-55 TB	ffffe8fffffffffff	32 TB	vmalloc/ioremap space (vmalloc_base)
ffffe90000000000	-23 TB	ffffe9fffffffffff	1 TB	... unused hole
ffffea0000000000	-22 TB	ffffeaffffffffffff	1 TB	virtual memory map (vmemmap_base)
ffffeb0000000000	-21 TB	ffffebfffffffffff	1 TB	... unused hole
ffffec0000000000	-20 TB	fffffbfffffffffff	16 TB	KASAN shadow memory
				Identical layout to the 56-bit one from here on:
fffffc0000000000	-4 TB	fffffdfffffffffff	2 TB	... unused hole
fffffe0000000000	-2 TB	fffffe7ffffffffff	0.5 TB	vaddr_end for KASLR
fffffe8000000000	-1.5 TB	fffffefffffffffff	0.5 TB	cpu_entry_area mapping
ffffff0000000000	-1 TB	ffffff7ffffffffff	0.5 TB	... unused hole
ffffff8000000000	-512 GB	ffffffeeffffffffff	444 GB	%esp fixup stacks
ffffffe800000000	-68 GB	ffffffeeffffffffff	64 GB	... unused hole
ffffffe900000000	-4 GB	ffffffe97fffffffff	2 GB	EFI region mapping space
ffffffe980000000	-2 GB	ffffffe99fffffffff	512 MB	... unused hole
ffffffe980000000	-2048 MB			kernel text mapping, mapped to physical address 0
ffffffe9fa000000	-1536 MB	ffffffe9fefffffffff	1520 MB	module mapping space
ffffffe9ff000000	-16 MB			
FIXADDR_START	~-11 MB	ffffffe9ff5fffff	~0.5 MB	kernel-internal fixmap range, variable size and of
ffffffe9fff60000	-10 MB	ffffffe9fff600fff	4 kB	legacy vsyscall ABI
ffffffe9fffe0000	-2 MB	ffffffe9fffefffffff	2 MB	... unused hole

Complete virtual memory map with 5-level page tables

Note

- With 56-bit addresses, user-space memory gets expanded by a factor of 512x, from 0.125 PB to 64 PB. All kernel mappings shift down to the -64 PB starting offset and many of the regions expand to support the much larger physical memory supported.

Start addr	Offset	End addr	Size	VM area description
0000000000000000	0	00ffffffffffffff	64 PB	user-space virtual memory, different per mm
0100000000000000	+64 PB	feffffffffffffff	~16K PB	... huge, still almost 64 bits wide hole of non-canonical virtual memory addresses up to the -64 PB starting offset of kernel mappings.

					Kernel-space virtual memory, shared between all pr
ff00000000000000	-64 PB	ff0fffffffffffffff	4 PB	...	guard hole, also reserved for hypervisor
ff10000000000000	-60 PB	ff10fffffffffffffff	0.25 PB	LDT remap for PTI	
ff11000000000000	-59.75 PB	ff90fffffffffffffff	32 PB	direct mapping of all physical memory (page_offset	
ff91000000000000	-27.75 PB	ff9fffffffffffffff	3.75 PB	...	unused hole
ffa0000000000000	-24 PB	ffd1fffffffffffffff	12.5 PB	vmalloc/ioremap space (vmalloc_base)	
ffd2000000000000	-11.5 PB	ffd3fffffffffffffff	0.5 PB	...	unused hole
ffd4000000000000	-11 PB	ffd5fffffffffffffff	0.5 PB	virtual memory map (vmemmap_base)	
ffd6000000000000	-10.5 PB	ffdefffffffffffffff	2.25 PB	...	unused hole
ffdf000000000000	-8.25 PB	fffffbffffffffffff	~8 PB	KASAN shadow memory	
					Identical layout to the 47-bit one from here on:
fffffc0000000000	-4 TB	fffffdffffffffffff	2 TB	...	unused hole
fffffe0000000000	-2 TB	fffffe7ffffffffffff	0.5 TB	vaddr_end for KASLR	
fffffe8000000000	-1.5 TB	fffffeffffffffffff	0.5 TB	cpu_entry_area mapping	
fffffe9000000000	-1 TB	fffffe7ffffffffffff	0.5 TB	...	unused hole
fffffe8000000000	-512 GB	fffffeefffffffffffff	444 GB	%esp fixup stacks	
fffffe9000000000	-68 GB	fffffeffffffffffff	64 GB	...	unused hole
fffffe9000000000	-4 GB	fffffe7ffffffffffff	2 GB	EFI region mapping space	
fffffe9000000000	-2 GB	fffffe79ffffffffffff	512 MB	...	unused hole
fffffe9000000000	-2048 MB	fffffe79ffffffffffff	512 MB	kernel text mapping, mapped to physical address 0	
fffffe9000000000	-1536 MB	fffffe79ffffffffffff	1520 MB	module mapping space	
fffffe9000000000	-16 MB	fffffe79ffffffffffff	1520 MB		
FIXADDR_START	~-11 MB	fffffe79ffffffffffff	~0.5 MB	kernel-internal fixmap range, variable size and of	
fffffe9000000000	-10 MB	fffffe79ffffffffffff	4 kB	legacy vsyscall ABI	
fffffe9000000000	-2 MB	fffffe79ffffffffffff	2 MB	...	unused hole

Architecture defines a 64-bit virtual address. Implementations can support less. Currently supported are 48- and 57-bit virtual addresses. Bits 63 through to the most-significant implemented bit are sign extended. This causes hole between user space and kernel addresses if you interpret them as unsigned.

The direct mapping covers all memory in the system up to the highest memory address (this means in some cases it can also include PCI memory holes).

We map EFI runtime services in the 'efi_pgd' PGD in a 64Gb large virtual memory window (this size is arbitrary, it can be raised later if needed). The mappings are not part of any other kernel PGD and are only available during EFI runtime calls.

Note that if CONFIG_RANDOMIZE_MEMORY is enabled, the direct mapping of all physical memory, vmalloc/ioremap space and virtual memory map are randomized. Their order is preserved but their base will be offset early at boot time.

Be very careful vs. KASLR when changing anything here. The KASLR address range must not overlap with anything except the KASAN shadow area, which is correct as KASAN disables KASLR.

For both 4- and 5-level layouts, the STACKLEAK_POISON value in the last 2MB hole: ffffffff4111