

Large Virtual Address support (52-bit) in ARM64 kernel

SPEAKER

Bhupesh Sharma

Linaro

<bhupesh.sharma@linaro.org>

Twitter: @bhupesh_sharma

\$whoami

- Part of Linaro Landing team.
- Been hacking on boot loaders & kernel since past 15 years.
- Contribute to:
 - Linux,
 - EFI/u-boot bootloader, and
 - User-space utilities like:
 - kexec-tools, and
 - makedumpfile.
- Co-maintain crash-utility tool



Outline

- Large VA support for arm64
 - What?
 - How?
- 52-bit VA kernel support arm64
- Flipping the arm64 kernel memory layout
- Impact on user-space applications
- 52-bit userspace VA
- How to test 52-bit VA support
- Next Steps

Large VA support for arm64 – What?

- 64-bit hardware is can address very large address space.
 - Upto 16 EiBs (16 × 10246 = 264 = 18,446,744,073,709,551,616 bytes)
 - Approx 18.4 exabytes of memory. 🎉
- Servers available with \Rightarrow 64 TiB (& upwards) of memory.
 - Use-cases \$\Rightarrow\$ requiring \$\Rightarrow\$ addressing spaces > 2^48 bytes.
- Limitations
 - Not all instruction sets, and
 - Not all processors,



• support a full 64-bit virtual or physical address space.

Large VA support for arm64 – What?

• x86_64

- Supports 5-level page tables in both <u>Hardware</u> & <u>Software</u>.

- Allows addressing address space = 2^57 bytes.
- Bumps limits to
 - \rightarrow 128 PiB of virtual address space,
 - \rightarrow 4 PiB of physical address space.
- arm64
 - Introduces 2 new architecture extensions
 - 52-bit addressing extensions
 - ARMv8.2 LVA, and
 - ARMv8.2 LPA
 - Allows addressing
 - \rightarrow 4 PiB of virtual address space,
 - \rightarrow 4 PiB of physical address space.









- ARMv8.2 LVA
 - Supports larger VA space
 - Each translation table base register of up to 52 bits
 - when using the 64KB translation granule.
- ARMv8.2 LPA
 - Allows larger intermediate physical address (IPA), and
 - PA space of up to 52 bits when using the 64KB translation granule.
 - Allows a level 1 block size where the block covers a 4TB address range for the 64KB translation granule if the implementation support 52 bits of PA.
 NOTE: These features are supported in AArch64 state only.
- Cortex-A processors with ARMv8.2 extension support:
 - Cortex A55
 - Cortex A75
 - Cortex A76

Hardware

Support



- arm64 Linux uses:
 - 4KB page size
 - 39-bit (512GB) virtual addresses \Rightarrow 3 translation tables levels
 - 48-bit (256TB) virtual addresses \Rightarrow 4 translation tables levels.

Translation table lookup with 4KB pages:

Reference : <u>Documentation/</u> <u>arm64/memory.</u> <u>rst</u>

+	+	+	+	+	+	+	+	+
63	56 55	48 47	40 39	32 31	24 23	16 15	8 7	0
+	+	+	+	+	+		+	+
1		1	1	1				
1		1		1		v	1	
1		1		1		[11:0]	in-page	e offset
1		1	1	1	+	-> [20:12]	L3 ind	ex
1		1	1	+-		-> [29:21]	L2 ind	ex
1		1	+			-> [38:30]	L1 ind	ex
1		+				-> [47:39]	L0 ind	ex
+						-> [63] TT	BR0/1	

Software

Support

- arm64 Linux uses:
 - 64KB page size
 - - but the memory layout is the same.
 - 48-bit (256TB) virtual addresses \Rightarrow 3 translation tables levels
 - 52-bit (4PiB) virtual addresses 🖒 3 translation tables levels with ARMv8.2 extension
 - expands number of descriptors in the first level of translation.

Translation table lookup with 64KB pages:

Reference : <u>Documentation/a</u> <u>rm64/memory.rs</u> <u>t</u>

+	+		+				+	+
63	56 55	48 47	40 39	32 31	24 23	16 15	8 7	0
+	+		+			+	+	+
1		1	1			1		
1		1	1			v		
1		1	1	1		[15:0]	in-page	offset
1		1	1	+		-> [28:16] L3 inde	x
1		1	+			-> [41:29]] L2 inde	x
1		+				-> [47:42]] L1 inde	x (48-bit)
1						[51:42]] Ll inde	x (52-bit
+						-> [63] T	TBR0/1	

Software

Support

• A sample arm64 translation table walk



Reference : ARMv8-A Address Translation from ARM

52-bit VA kernel support - arm64

- Design problem from a software support p-o-v
 - Older arm64 CPUs which don't support ARMv8.2 extensions.
 - New / Upcoming arm64 CPUs which support ARMv8.2 extensions.
- Selected design approach
 - Have a single kernel binary
 - At *early boot time* check if the ARMv8.2 hardware feature is present or not.



52-bit VA kernel support - arm64

- Single kernel binary for both 48-bit and 52-bit VA spaces.
- VMEMMAP constraints
 - must be sized large enough for 52-bit VAs, and
 - must be sized large enough to accommodate a fixed PAGE_OFFSET.
- VA bits related variables used by kernel code:

VA_BITS	Compile time constant	Maximum size of VA space, used for things like static array and region sizes
VA_BITS_MIN	Compile time constant	Minimum size of VA space, used to ensure pointers are addressable
VA_BITS_ACTUAL	Variable	The actual size of the kernel VA space

* vabits_actual

Reference : Documentation/arm64/memory.rst

Flipping the kernel memory layout - arm64



Direct Linear Map is in Higher Half of the VA space.

Reference : <u>Documentation/arm64/memory.rst</u>

- kernel text addresses is kept constant, even for 48 to 52-bits migration.
- We need to flip the VA space.

Impact on user-space applications - arm64

- User-space applications impacted due to flipped kernel memory layout
 - used to debug running / live kernels,
 - to analyze vmcore dumps.
 - for example: <u>kexec-tools</u>, <u>makedumpfile</u> & <u>crash-utility</u>.
- Debugging applications need to perform a va_to_pa() conversion
 - Walk the translation table(s) for determining the physical address
- These applications are broken upstream currently.
- Have proposed fixes for affected applications
 - some have been accepted upstream,
 - others are still pending
 - kexec-tools <u>fix</u>

52-bit userspace VA - arm64



⇒ What happens to *other* existing applications?

- To maintain backward compatibility
 - kernel will, *by default*, return virtual addresses to userspace from a 48-bit range.
- Opt-in model for willing user-space application(s)
 - Hint parameter is passed to mmap() calls to receive addresses in 52-bit range.

maybe_high_address = mmap(~0UL, size, prot, flags,...);

- How to build kernel which returns addresses in 52-bit range?
 - Enable CONFIG options:
 - CONFIG_EXPERT && CONFIG_ARM64_FORCE_52BIT

NOTE: Only intended for debugging + should **not** be used in production.

How to test 52-bit VA support

- What if I have *no* arm64 *hardware*?
- Use <u>qemu</u>

\$ virt-builder --arch aarch64 --root-password password:fedora fedora-30

Reference: Fedora AArch64 WiKi

fedora WIKI

\$ qemu-system-aarch64 -no-user-config -nodefaults -display none -m 2048 -cpu cortex-a57 -machine virt -smp 4 bios /usr/share/edk2.git/aarch64/QEMU_EFI.fd -device virtio-scsi-device,id=scsi -drive file=fedora-30.img,format=raw,if=none,id=hd0 -device scsi-hd,drive=hd0 -netdev user,id=usernet -device virtio-netdevice,netdev=usernet -boot efi -serial mon:stdio

Use <u>ARMv8 fast model</u> simulator for some quick checks

FVP terminal_0	×	FVP terminal_1	×			
Fedora (5.11.0) 30 (Thirty) Fedora (5.0.9-301.fc30.aarch64) 30 (Thirty) Fedora (0-rescue-5a17fec93c844eda9d2cc151033be19b) 30 (Thirty) System setup Use the ^ and v keys to change the selection. Press 'e' to edit the selected item, or 'c' for a command prompt.	1	0x16F4). add-symbol-file /home/zsun/Besktop/arm-52bit-simulator/fw/Build/ArmVExpress Ahrch64/DEBUG_6CC5/ARCCH64/OwmFrkg/VirtioBlkDxe/VirtioBlk/DEBUG/VirtioBlkD OxFAEC04000 Loading driver at 0x000FAE03000 EntryPoint=0x000FAE06468 VirtioBlkDxe,efi add-symbol-file /home/zsun/Besktop/arm-52bit-simulator/fw/Build/ArmVExpress AArch64/DEBUG_6CC5/ARCCH64/HdeHoduleFkg/Universal/FvSimpleFileSystemIxe/Fv FileSystemDxe/DEBUG/VSimpleFileSystem_d10 vXFADFA000 Loading driver at 0x000FADFS000 EntryPoint=0x000FADFE418 FvSimpleFileSystem FatOpenDevice: read of part_lba failed No Media Installed Fat filesystem on FADAEF98 [Bds]Booting UEFI Misc Device 3 [Bds]Booting UEFI Misc Device 3 [Bds]Booting UEFI Misc Device 5 [Bds]Booting UEFI Misc Device 5 [Bds]	s-FVP- xe.dll s-FVP- Simple m.efi			
Fast Models - CLCD RevC 2xAEM	Fast Models - CLCD RevC 2xAEMv8A Base RevC FVP ×					
<pre>+ON USERSW 18 ******* S6LED07 Daughter ******** +ON BODTSW 18 ******** Cluster0 **** Total Instr: 1,089,194,689 Total Ti Cluster1 **** Total Instr: 0</pre>	me: 20s	Rate Limit ON Grab mouse: LeftCtrl+LeftAlt				

Next Steps

- Fix broken userspace applications WIP.
- Create awareness about the flipped kernel address map.
- JIT and other applications need to made aware about the mmap() hint parameter and how they can receive addresses in 52-bit range.
- Ard Biesheuvel has posted a patchset to extend the Linear range for 52-bit configurations:
 - https://lore.kernel.org/linux-arm-kernel/20201008153602.9467-3-ardb@kernel.org/
- Test upstream kernel on both old CPUs (48-bit VA) and new CPUs with 52-bit VA)
 - In absence of a real 52-bit HW, you can use <u>ARMv8 fast model</u> simulator for some quick checks.

Slides can be found on github

Email: <bhupesh.sharma@linaro.org> Twitter: @bhupesh_sharma

The talk is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License.

Questions?