**Chapter 7** 

# Pagetables and the VMSA

This chapter is based, in part, on: Chapter D5 of the Arm Architecture Reference Manual DDI 0487H.a; and, Relaxed
 virtual memory in Armv8-A [54] by Ben Simner, Alasdair Armstrong, Jean Pichon-Pharabod, Christopher Pulte,
 Richard Grisenthwaite, and Peter Sewell, published in the proceedings of the 31st European Symposium on Programming
 (ESOP, 2022).

## 1277 **7.1 Introduction**

Modern computers heavily rely on *virtual memory* to enforce security boundaries: hypervisors and operating systems manage mappings from virtual to physical addresses in order to restrict the access individual processes and guest operating systems have to the underlying physical memory, and to memory-mapped devices. With the endemic use of memory-unsafe languages, even for critical infrastructure, understanding and verifying the programs which manage virtual memory mappings is more vital than ever, driving current interests in hypervisors. The virtual machines those hypervisors enable are the key pieces of software which have become solely responsible for implementing such critical security properties.

The following chapters focus on these aspects of the architecture, on virtual memory and virtualisation and the software they enable, with the aim of giving a precise formal semantics for the purpose of verifying real systems software which use those features.

I first give a description of the sequential behaviour of Arm's virtual memory (this chapter); then describe the *relaxed* behaviours and any open questions about Arm's virtual memory (§8); give our precise axiomatic semantics that capture these behaviours (§9); give an overview of the tooling and validation of the given model(s) (§10); and, finally, a sketch of an equivalent operational semantics (§11).

**This chapter overview** The remainder of this chapter will give: a brief overview of Arm's virtual memory systems architecture (§7.2); a detailed description of the Arm translation table format (§7.3); an overview of the multiple stages of translation (§7.4), and the different translation regimes (§7.5); a detailed explanation of the official Arm translation table walk pseudocode (§7.6); and finally a discussion on the existence and purpose of translation lookaside buffers (§7.7). This chapter does not present any new contributions or novel research, instead, it is a brief but necessary overview of the required architectural features.

# 1298 7.2 Virtual Memory

Armv8-A8's *virtual memory system architecture* (or *VMSA*) defines the virtual memory and virtulisation features
 of the Arm architecture. Its structure is described, in detail, in Chapter D5 of the Arm Architecture Reference
 Manual [1].

<sup>1302</sup> Conventionally, we think of memory as being a flat array of bytes, indexed by *physical addresses*. For smaller
 <sup>1303</sup> trusted devices, such as microcontrollers, this may be the end of the story. However larger 'application' class
 <sup>1304</sup> processors rely heavily on virtual addressing: interposing one or more layers of indirection between the accesses
 <sup>1305</sup> using the *virtual* addresses of the program and the 'true' physical addresses of memory. This indirection allows
 <sup>1306</sup> systems running on those processors to:

1271

1272

- partition the physical resources between different programs, giving access to only those resources that each
   program needs, and protecting those resources from other programs that do not need to access them;
- 1309 2. indirect accesses through specific ranges of addresses with convenient numeric values; and
- update those indirections at runtime to add, remove, or otherwise modify, the mappings to physical memory, to support techniques such as copy-on-write and paging.

To manage all this, typical operating systems splits the programs into distinct *processes* and associates each process with its own virtual to physical mapping. These mappings take the form of partial functions from the process' own (virtual) addresses to the real hardware physical addresses along with some permissions:

 $\texttt{translate}: \texttt{VirtualAddress} \rightarrow \texttt{PhysicalAddress} \times 2^{\{\texttt{Read},\texttt{Write},\texttt{Execute}\}}$ 

<sup>1315</sup> Note that this is a simplification. A more accurate translate function is given later on TODO: ?REF?.

Typically an operating system would create one such mapping for every process, partitioning the physical memory 1316 into disjoint subsets of physical addresses (the range of the translate function), and would allocate some convenient 1317 numeric values to be the virtual addresses the process interacts with (the domain of the translate function). Having 1318 this separation allows the processes to be given convienently aligned contiguous chunks of virtual address space 1319 even if the underlying physical resources are highly fragmented, or, in the case of paging, potentially not present 1320 in memory at all. Additionally, operating systems can provide many processes with mappings to the same physical 1321 resource (such as memory-mapped devices) and control which processes have access to such devices at any point 1322 in time. 1323

These mappings give rise to separate *address spaces* for each process. The diagram in Figure 7.1 illustrates an example, with two processes named P0 and P1 each with their own virtual address space. The left-hand side shows a representation of the 'memory' as the processes see it, with the memory split into *pages* (fixed-size blocks of contiguous addresses). The right-hand side is the equivalent representation of the actual physical memory, with each *physical* page of the available RAM. Note that this diagram shows the virtual address space as being smaller than the physical one, but in general, they may be the same size, or the virtual address space may be even larger than the physical space.

 $_{1331}$  If we assume each page has size 0x1000 then page 1 contains addresses 0x1000 to 0x1FFF inclusive, and we can interpret the diagram like so:

1333 ▷ For P0:

1334

1335

1338

- virtual addresses in pages 1, and 3 are unmapped.
- virtual addresses in pages 0 and 2 map to physical addresses in physical page 1.
- virtual addresses in page 4 map to physical addresses in physical page 5.
- 1337 ▷ For P1:
  - virtual addresses in pages 0 and 4 are unmapped.
- virtual addresses in page 1 map to physical addresses in physical page 5.
- virtual addresses in page 2 map to physical addresses in physical page 7.
- virtual addresses in page 3 map to physical addresses in physical page 8.

For example, if process P0 reads the address 0x2305, it will actually read from the physical location 0x1305, since virtual page 2 was mapped to physical page 1 in P0's address space.

Each address space corresponds to a distinct translate function. Note that these mappings may be: non-injective
(contain *aliasing* of multiple virtual addresses to the same physical address); partial (where some virtual addresses
do not map to a physical address at all); or overlapping with other processes' address spaces, in either the domain
(for example, the physical page 5 is mapped in both P0 and P1), or range (for example, the virtual page 2 is mapped
in both P0 and P1 but to different physical pages), or both.

Large application-class processor architectures, such as Armv8-A, often provide hardware support in the form of the *memory management unit* (the MMU), which, once configured by software, will automatically perform the translation from virtual to physical addresses. Software is then required to manage a set of translation functions, and is responsible for ensuring the correct translation function is being used by the MMU whenever a context

<sup>1353</sup> switch occurs, and handle any *faults* that the MMU generates.



Figure 7.1: Example virtual and physical address spaces for two processes.

# 1354 **7.3 Arm Translation Tables**

<sup>1355</sup> Software configures the MMU through the creation and modification of sets of *translation tables* (also referred to <sup>1356</sup> as *page tables*) for each of the translation functions.

The translation tables form an in-memory tree data structure which encode the (partial) translate function. Software creates and maintains these trees, and tells the MMU which tree (and so which translation function) to use at runtime. The hardware then reads from this tree structure to perform the translation, or from one of the

various caching structures described in TODO: ?REF?, whenever the process reads from, or writes to, memory.

A pointer to the root of the tree is stored in the TTBR ("Translation table base register") register (or rather, one of the various base registers described in more detail in **TODO**: **?REF?**), and this determines which translation function is currently in use by that processor.

Each node in the tree is a page-aligned chunk of memory which is treated as an array of 64-bit entries. Each entry
 is responsible for mapping some fixed part of the domain of the translation function, with the root table mapping
 the entire address space.

The tree is separated into different *levels*. with a root table pointed to by the base register and each subsequent child tree increases in level going deeper into the tree. Typically the root is at level 0 with a maximum depth of 4 (down to level 3), but the various configurations are discussed in the next section.

#### 1370 7.3.1 Translation table format

Arm's virtual memory system architecture is highly configurable. Writing to the SCTLR ("System control register") 1371 and TCR ("Translation control register") system registers allow the software developer to choose a configuration 1372 from a whole host of various options. To give a flavour of this configurability I list some of the configuration bits, 1373 some of which will be discussed in more detail in the next chapter; these include: the size of virtual addresses; the 1374 number of levels in the tree; the starting level; the size of a single page (or in Arm terminology, the size of the 1375 translation granule); the number of ASIDs and VMIDs; alignment requirements; memory attributes for hardware 1376 walks; enabling hardware management of access flags and dirty bits; write-execute-never permissions; and so on. 1377 To simplify things, in this dissertation, we consider just one common configuration, the one currently used by the 1378 Linux kernel: a tree of translation tables with maximum depth 4, with 4KiB pages with 48-bit addresses, unless 1379 explicitly stated otherwise. 1380

Each node is a table of 512 64-bit entries, bound as one 4096-byte block of memory. Each table controls the mapping of a fixed range of the virtual address space. This range is split into 512 equal slices, with each entry responsible for its slice. Each of those entries can be one of:

1. An *invalid* entry, which indicates that this slice of the domain is unmapped; 1384

2. A table entry, pointing to a next-level table (a child tree) which recursively maps this slice of the domain; or 1385

3. A page (last-level) or block (non-last-level) entry which defines a single fixed-size mapping for this slice of 1386 the domain.

1387

1391

1402

1407

1408

**Invalid entries** An invalid entry is defined by bit[0] of the entry being 0. The top 63 bits are ignored by 1388 hardware, and software is free to use those bits to store any metadata it wishes. Invalid entries may exist at any 1389 level in the tree. 1390

63		10
	ignored	0

**Block or page entries** Block and page entries are similar to each other; both create a mapping for a contiguous 1392 slice of the domain mapped by the entry, encoded as an output address (OA) with some metadata (including access 1393 permissions, memory type, and some software-defined bits). 1394

The OA is aligned to the size of the slice of the domain being mapped. For page entries, the OA is aligned on a 1395 page boundary. A block entry's OA at level 2 would be 2MiB aligned, and a block entry's OA at level 1 would 1396 be GiB aligned. This corresponds to the hardware reserving bits[n:12] of the entry to be 0 depending on how 1397 deep the entry is: at level 1 n=30; at level 2 n=21; and at level 3 n=12. 1398

Block entries can exist at levels 1 and 2. Page entries can only exist at level 3. 1399

For block entries bit[1] is 0, for page entries bit[1] is 1. 1400

Metadata (access permissions, shareability, memory type) are encoded into the attrs bits. 1401

_	63	50	47	n	(n-1)	12	11	210
	attrs	0	0	output address	iĮ	gnored	attrs	P1

**Table entries** A table entry contains a page-aligned pointer to a child table, but can also contain similar 1403 metadata as the block or page entry, including access permissions (read/write/execute), which are combined with 1404 any permissions from the child table. 1405

Table entries are allowed only at levels 0-2. 1406

63	50	47	12	11	210
atti	:s (	00	table pointer	Rest	$)^1$ 11

#### The Arm translation table walk 7.3.2

When the processor executes an instruction which takes an address, such as the Arm LDR or STR instructions, those addresses are virtual (addresses used by instructions are always virtual addresses). The hardware converts 1410 each virtual address to a physical address, and the MMU performs this conversion. 1411

To do this, the MMU reads the TTBR to get the currently in-use tree of translation tables. Then the MMU itself 1412 reads memory and walks the tree (except when it can read from a previously cached translation, as described in 1413 the next chapter) effectively computing the partial translate function the tree encodes, producing the physical 1414 address and any permissions, or reporting a fault back to the processor if the virtual address was unmapped, or if 1415 the permissions forbid the requested operation. 1416

Walk overview The hardware walker first slices up the input virtual address into chunks: the most-significant 1417 bit is used to determine which base register to use (see §7.5); the next 15 bits are typically ignored by hardware; 1418 the rest of the address is split into 9-bit fields which we refer to as fields a-d, with the remaining bits as field e. 1419 Fields a-d are used for indexing into the tables; and field e is the offset in the page, which is added to the final 1420 output address. 1421

<sup>&</sup>lt;sup>1</sup>The Arm architecture requires these bits are 0 and are reserved for future use.

							VA								
	63	62		48 47	39	38	30	29		1 2	0	12	11		0
	Ē		ignored		a		b		с		d			e	
The	e walk	then p	roceeds, wi	th the M	IMU takir	ng the	e followii	ng st	eps:						
	1 Rea	d the T	TBR registe	er.											
	2 Per	form a	64-bit singl	le-copy a	atomic rea	ad of	Mem[TTB	R+8*	a] to r	ead t	he ent	try in	the L	evel 0 ta	able. Ca
	rest	ılt L0ei	ntry.												
	C	ı If LØ	entry[0] i	s 0 (that	is, it's an	inva	lid entry	) the	n repor	t a fa	ault ba	ick to	the p	rocesso	r.
	l	b Othe	rwise if L0	entry[1	] is 0 the	n rep	ort a fau	lt ba	ck to th	e pr	ocesso	or (toj	p-leve	el tables	cannot
		block	c mappings	).											
	3 Per	form a	64-bit sing	le-copy	atomic re	ad of	Mem[L0	entr	y.tabl	e_p	ointe	r+8*b	] to 1	read the	entry i
	Lev	el 1 tal	ole, which v	ve will c	all L1ent	ry.									
	6	ı If L1e	entry[0] i	s 0 then	report a f	ault l	back to th	ne pr	rocesso	•					
	l	b If L10	entry[1] i	s 0 (it's a	ı block en	try):									
	<i>i</i> If the access is not permitted (See §7.3.2 "Access permissions"), report a fault to the processor														
		ii (	Otherwise,	return tl	ne output	addr	ess (See	§7.3.:	2 "Com	put	ing th	e fin	al ou	tput ad	dress")
		1	to the proce	essor.									_		
	4 Per	form a	64-bit sing	le-copy	atomic re	ad of	Mem[L1	entr	y.tabl	e_p	ointe	r+8*c	2] to 1	read the	entry i
	Lev	el 1 tat	ole, which v	ve will c	all L2ent	ry.									
	(	a If L20	entry[0] i	s 0 then	report a f	ault l	back to the	ne pr	rocesso	•					
	<i>b</i> If L2entry[1] is 0 (it's a block entry):														
		i	If the acces	s is not p	permitted	, repo	ort a faul	t to t	the proc	esso	or.				
	- D	и ( с	Otherwise,	return t	he output	addr	ess back	to th	he proc	esso	r. .,			1.1	
	5 Per	torm a	64-bit sing	le-copy	atomic re	ad of	Mem[L20	entr	y.tabi	e_p	ointei	r+8*0	1] to 1	read the	entry 1
	Lev	el 3 tat	ble, which v	ve will c	all L3ent	ry.	1								
	(		entry[0]1		report a I	ault I		ne pr	cocesso	: /11 ·		ı		1	1
	l	2 Else 1	II L 3 entry[	<b>I I I S O</b> , <b>I</b>	report a fa	ault d	ack to th	e pro	ocessor	tnis	encoc	iing is	s resei	rved and	i is treat
		invai	10).	(:+,											
	(	: Lsen	ury∐i]1S.I Ifthologoog	i (it s a p	age entry	): ron	wt o four	t ta t	ho nro	0000					
		:: .	Othornuico	s is not p	be outrout	, repo		ι ιυ Ι + ο +1	he proc	esso	)1. 				
		11	onierwise,	return t	ne output	audr	ess Dack		he proc	:550	ι.				
														Leve	13



1451

**Computing the final output address** The output address (OA) of the final descriptor is the start of the range 1452 mapped by the entry. The low order bits are all 0 in the output address, and need to be added on to compute the 1453 final output address of the translation. 1454

To compute this final output address the MMU takes the OA from the entry, and the level in the tree the entry is 1455 at, and 'completes' the address by bitwise appending the remaining fields to create the complete 48-bit output 1456 address. Recall that the OA field of the block mappings gets wider the deeper in the tree you are, and so for a 1GiB 1457 entry the OA field is only 18 bits wide but for a 4KiB page entry its OA field is the full 36 bits. 1458

▷ For a 1GiB (level 1) block entry; PA = OA::c::d::e 1459

▷ For a 2MiB (level 2) block entry; PA = OA::d::e 1460

1461  $\triangleright$  For a 4KiB (level 3) page entry; PA = OA::e

<sup>1462</sup> Note that this process means that the least-significant 12 bits of the input VA are unchanged and remain the same <sup>1463</sup> in the final output PA, regardless of how the translation function is configured.

**Access permissions** Once the walk is complete, and the final output address calculated, the MMU checks to see whether the requested access is permitted. Each level of the table can contain some access permissions and those permissions get combined at the end to calculate the final permissions.

<sup>1467</sup> For data accesses (reading and writing), table entries have an APTable field (bits[62:61]), and block/page entries <sup>1468</sup> have AP[2:1]<sup>1</sup> field (bits[7:6]). These fields can be decoded using the following table:

1469	

Field	When set (1)	When unset (0)
AP[2]	Read-only	Read&Write
AP[1]	Allow at EL1&0	Allow at EL1 only
APTable[1]	Force read-only	No effect on permissions.
APTable[0]	Force forbid access at EL0	No effect on EL0 permissions.

For executable permissions, which permit or forbid instruction fetching from some region of memory, there are no dedicated encodings of the access permission bits. Instead, all mappings are executable by default, unless one of the following applies: the region is mapped writeable at EL0, as writable EL0 regions are never executable at EL1; a global WXN ("Write-execute-never") configuration bit is set, and the entry was writeable; or, when one of the various translation table entry XN ("Execute-never") bits are set. For simplicity, this chapter assumes that execute-never bits are always disabled; see the full description in the Arm ARM TODO: ?REF?for more information.

To combine access permissions from the whole walk, the MMU takes the bitwise union of each of the APTable
fields from each table entry, and then intersects the result with the final AP[2:1] field to produce a final set of
permissions. Figure 7.2 contains a decoding table for a given table and leaf access permissions, for testing whether
a requested access is permitted. If the requested access is not permitted, then the MMU generates a permission
fault, which is reported back to the processor.

Faults The MMU may emit one of several fault types during a translation table walk (these are referred to by
 Arm as the *MMU fault* types):

- 1484  $\triangleright$  Translation fault.
- 1485These are caused by the mapping being invalid, either because bit[0] was 0, or because the descriptor1486encoding was reserved-as-invalid. Translation faults also result from trying to translate an address that is1487outside the 48-bit input address range.
- 1488  $\triangleright$  Permission fault.

1489

1490

1492

1498

- For when the mapping was valid, but the access permissions do not permit the requested access (for example, trying to write to a read-only address).
- $^{1491}$   $\triangleright$  Access flag fault.
  - These are generated when hardware management of access flags is disabled and the access flag bit is set.
- <sup>1493</sup> ▷ TLB Conflict aborts (see **TODO: ?REF?**).
- <sup>1494</sup> ▷ Alignment fault.
- Generated when an operation expects an aligned memory address, but is given a misaligned one, and alignment checking is enabled in the SCTLR.
- <sup>1497</sup> ▷ Address size fault.
  - For when the OA, or TTBR, has a value that is out of the physical address range.
- <sup>1499</sup> > Synchronous external abort on a translation table walk.
- These are *external aborts* (that come from the system not from the MMU) that happen due to accesses that the MMU generated. For example, if the next-level table field pointed to an address for which there was no memory or device, the system-on-chip would return a fault to the processor.

These faults lead to processor exceptions. The fault type is stored in the ESR\_ELn ("exception syndrome register")
register's EC ("exception class") field, and any supplementary information is stored in its ISS ("instruction specific
syndrome") field (such as which level in the tree the fault came from, whether the originating instruction was a
read or a write, and ). Exception handling code can read the ESR register to determine the fault type and cause,

<sup>&</sup>lt;sup>1</sup>Block/page entries do not store the entire AP field but only AP[2:1]. AP[0] is not present in AArch64.



**Figure 7.2:** Merging Access Permissions (Stage 1, EL1&0). Entries in **red** highlight differences from the APTable=00.

and can read the FAR\_ELn ("fault address register") to determine the virtual address which triggered the fault, and
 handle the fault appropriately.

Memory Attributes The processor does not necessarily know what is located at any physical address. There may be some dynamic random-access memory (DRAM, what one would generally consider 'memory'), but there may also be other memory-mapped devices, or non-volatile memory, or other peripherals, or possibly nothing at all.

To help accommodate this, hardware allows software to mark regions of memory as one of either *device* memory, *normal cacheable* memory, or normal *non-cacheable* memory, using the translation tables.

The desired memory type is determined from the AttrIndx field (bits[4:2]) in block and page entries. Instead of being directly encoded into this field, Arm chose to have the actual attributes stored in a separate register: the MAIR ("Memory attribute indirection register") register. The MAIR stores an array of eight 8-bit fields each of which contains an encoding of a memory type. The AttrIndx field in the entry is an integer in the range 0–7, which is the index of the field in the MAIR register to use.

This indirection means that the final result of translation depends not only on the value of the final leaf entry in memory, but on the value of certain system registers, such as the MAIR, at that time of the translation table walk.

Below are the three most common encodings for a MAIR field, and the ones that will be useful later when discussing
 tests:

- <sup>1524</sup> ▷ 0b0000\_0000: device memory.
- <sup>1525</sup> ▷ 0b0100\_0100: normal non-cacheable memory.
- <sup>1526</sup> ▷ 0b1111\_1111: normal cacheable memory, inner&outer write-back non-transient, read&write-allocating.

Memory locations marked as device tell the hardware that reads or writes to those locations may have side-effects.
This means hardware treats those locations differently: there will be no speculative instruction fetches, reads,
or writes to those locations; writes to those locations will not *gather* into larger writes; reads and writes to
those locations will not re-order with respect to others; those locations generally will not get cached; and other
thread-local optimizations get disabled. Note that Arm define a wide range of device memory types, allowing

the systems programmer to selectively re-enable some of the previously described behaviours to enable better
 performance where they deem it safe to do so.

For normal memory the software can choose between *cacheable* or non-*cacheable* memory. Arm provide a range of different options for the cacheability:

<sup>1536</sup> ▷ non-cacheable

<sup>1537</sup> ▷ write-back cacheable

1538  $\triangleright$  write-through cacheable

As with other features, there is a wide scope for configuration: separately configuring inner (L1,L2) and outer (L3) caches, and adding cache allocation hints (allocating on reads, writes or both).

As we will see later (TODO: ?REF?), the ability to change cacheability, or even have multiple aliases with different cacheability attributes, give rise to interesting behaviours and security considerations.

# **7.4** Virtualisation and a second stage of translation

<sup>1544</sup> So far this chapter has focused on operating systems and processes. However, modern systems isolate not just <sup>1545</sup> processes within an operating system but entire operating systems from one another, within a hypervisor.

To do this, software uses the virtual memory abstraction again, adding an extra layer. This layer, like the previous one, is supported by hardware. Processes use virtual addresses which are converted to *intermediate physical* (also sometimes known as *guest physical*) addresses using the operating system's configured translation tables but then these intermediate physical addresses (IPAs) go through another round of translation to convert those IPAs into the final physical address.

Arm calls these *stages* of translation, and the MMU supports both stages and can perform the full translation from virtual to physical (via the intermediate physical) address.

<sup>1553</sup> This means software must manage two sets of translation tables: operating systems manage the *stage 1* tables to

convert VAs to IPAs; and hypervisors manage *stage 2* tables to convert those IPAs to PAs; this gives two separate
 translate functions, which the MMU composes together at runtime:

 $\label{eq:translate_stage1:VirtualAddress} \rightharpoonup \mbox{IPA} \times \mbox{Permissions} \times \mbox{MemoryType} \\ \mbox{translate_stage2:IPA} \rightarrow \mbox{PhysicalAddress} \times \mbox{Permissions} \times \mbox{MemoryType} \\ \label{eq:translate_stage1}$ 

Hypervisors (running at EL2) can configure the stage 2 translate function by creating translation tables with a
similar format as before and then storing a pointer to the root of this tree in the VTTBR ("Virtualization translation
table base register") register. The MMU will read the VTTBR whenever it needs to perform a second-stage translation
to convert an IPA to a PA, and will do the translation table walk over that tree in much the same way as described
earlier for (what we can now call) the first-stage translation.

This results in two address spaces, a virtual address space and an intermediate-physical address space. Figure 7.3 contains an example layout of these address spaces for a machine running three processes (P0,P1,P2) in two operating systems (0S0,0S1). As with the earlier diagram in Figure 7.1, each column is a (set of) address spaces, with transformations between them defined by their respective translation functions. On the left-hand side are the virtual address spaces of the various processes, whose virtual addresses are translated (using the translation tables pointed to by the TTBR register) into intermediate-physical addresses in the central address spaces (for the respective OS). Those IPAs are then translated (using the VTTBR) into the final physical address.

<sup>1568</sup> Concretely, if P1 reads from address 0x1001, it will be translated into the IPA 0x3001 in 0S0's address space, <sup>1569</sup> which then gets translated again, and the processor will actually read from RAM at location 0x6001.

**Differences in the translation table format from stage 1** Stage 2 translation tables are similar to their stage 1 counterparts, but there are some minor differences:

<sup>1572</sup> ▷ Stage 2 table entries do not have any additional attributes, and so do not have an APTable field.

▷ Stage 2 AP field (called S2AP) has a slightly different (and simpler) format, see Figure 7.4.

 $_{1574}$   $\triangleright$  Stage 2 block and page entries do not have a MemAttrIndx field but rather encode the memory type directly

into the MemAttr field bits[5:2] (see the full description in the Arm ARM [1, D5-4874] for all possible
 encodings):



**Figure 7.3:** Example virtual, intermediate physical, and physical address spaces for three processes running on two operating systems.

Field	When set (1)	When unset (0)
S2AP[1]	Writeable	not Writeable
S2AP[0]	Readable	not Readable

Figure 7.4: S2AP field encoding.

- <sup>1577</sup> Øb0000: Device memory.
- 0b0101: Normal non-cacheable.
- 0b1111: Normal write-back inner&outer cacheable.

These are interesting as they mean that the stage 1 and stage 2 attributes (permissions and memory types) must be *combined* in order to produce the final output. This combination is not just a case of letting stage 2 overrule the stage 1 settings but rather that both stages get a veto: if stage 1 sets the memory type to be device or non-cacheable then it overrules what stage 2 sets. Similarly, if stage 1 permissions forbid an access then the stage 2 permissions cannot overrule that.

Second-stage translations during a first-stage walk There is a complication with the story so far. The stage 1 tables are created by the operating system, which is using an intermediate physical address space, not a physical one. The writes the OS does to the tables will be translated, as they are normal data writes. But, the tables themselves contain references to other tables, and those entries will be intermediate physical addresses, and so, they must also be translated, including the value of the TTBR itself.

In our assumed configuration of 4KiB pages and 4 levels of translation, this leads to a maximum of 24 memory accesses to perform the translation: 4 reads of stage 1 translation tables, 16 reads of stage 2 translation tables during those stage 1 walks, and a final 4 reads of the stage 2 translation tables to translate the output IPA into the final PA.

# **7.5** Translation regimes

As mentioned earlier, there are multiple translation table base registers. Each of them defines a translation function, pointing to the root of the tree of translation tables which define it. These translation functions are then composed together into various translation *regimes*, each defining the set of translation functions (and therefore which translation table base registers) which will be used for translations done by the processor.

Arm define a set of these translation regimes. Figure 7.5 gives an overview of three of the most common regimes,
 which are:

- <sup>1601</sup> ▷ EL1&0 (two-stage)
  - For programs executing at EL0 or EL1 when virtualisation (at EL2) is enabled.
    - VAs with the high bit set are translated into IPAs using the EL1-configured register, TTBR1\_EL1.
  - VAs are typically split into 'high' and 'low' regions with different translations, primarily used for separate kernel and user address spaces.
    - VAs without the high bit set are translated into IPAs using the EL1-configured register, TTBR0\_EL1.
    - IPAs are translated to PAs using the EL2-configured VTTBR\_EL2 register.
- <sup>1608</sup> ▷ EL1&0 (single-stage)
  - For programs executing at EL0 or EL1 when virtualisation (at EL2) is disabled.
  - VAs with the high bit set are translated into PAs using the EL1-configured register, TTBR1\_EL1.
  - VAs without the high bit set are translated into PAs using the EL1-configured register, TTBR0\_EL1.
- 1612 ▷ EL2

1602

160

1605

1606

1607

1609

1610

1611

1613

1614

1615

- For programs executing at EL2.
- VAs without the high bit set are translated into PAs using the EL2-configured register, TTBR0\_EL2.
- VAs with the high bit set are always unmapped.

<sup>1616</sup> Which translation regime is being used is defined by various system registers and the current system state.

- <sup>1617</sup> > Translations at EL1 or EL0 use one of the EL1&0 regimes.
- 1618  $\triangleright$  Translations at EL2 use the EL2 regime.
- <sup>1619</sup> ▷ TCR\_EL2 (set at EL2) determines whether the EL1&0 is a single-stage or two-stage regime.
- ▶ TTBR0\_EL1, TTBR1\_EL1 determine the stage 1 of the EL1&0 regimes, and can be set at EL1 or higher.
- <sup>1621</sup> ▷ TTBR0\_EL2 determines the stage 1 of the EL2 regime, and can only be set at EL2 or higher.

<sup>1622</sup> > VTTBR\_EL2 determines the stage 2 of the EL1&0 regime, and can only be set at EL2 or higher.

Arm define a wide range of other regimes, see the Arm ARM TODO: ?REF?. For simplicity, we ignore secure modes, including all of EL3.



Figure 7.5: Translation regimes that apply to EL0, EL1, and EL2.

# 1625 **7.6 Arm pseudocode**

1626 It is now useful to examine the official Arm pseudocode, especially those parts that relate to memory events.

We will do this in three steps: first, by looking at the pseudocode that is executed for an Arm store instruction; following the memory accesses that it performs down to any translations it performs; finally looking at the Arm translation table walker in full. There is a lot of detail infused throughout the Arm psueocode, so in this section we shall focus on the most pertient parts, and give some idea of what detail is omitted.

### <sup>1631</sup> 7.6.1 The lifecycle of a store

Arm give precise executable semantics for every instruction in their domain-specific Architecture Specification
 Language (ASL). This ASL code defines the sequential intra-instruction behaviour of each instruction, including
 memory accesses, and any translation table walks they perform.

```
1
          bits(64) address:
2
          bits(datasize) data;
3
4
6
7
          if n == 31 then
               if memop != MemOp_PREFETCH then CheckSPAlignment();
8
9
               address = SP[];
10
          else
               address = X[n];
11
12
13
15
16
18
21
                       data = X[t];
22
                   Mem[address, datasize DIV 8, acctype] = data;
23
43
```

Figure 7.6: Arm "STR (immediate)" ASL code.

# TODO: the importance of the ASL, and of sequential v concurrent behaviour will already be explained, but recap here anyway?

Figure 7.6 shows the Arm ASL for the "STR (Immediate)" instruction: STR Xt, [Xn]. This instruction writes the value contained in register Xt into the memory location stored in register Xn. The figure has some uninteresting (for this thesis) parts greyed out: those parts that deal with optional extensions such as memory tagging; unknown register values; register writeback; and, the load and prefetch instructions which use the same ASL code.

The ASL code first reads the virtual address either from the stack pointer (line 9) or by reading register Xn (line 11). It then reads the data from the register Xt (line 21), which will be written to memory. Finally, it performs the store itself using the Mem[] function (line 22).

#### 7.6.2 Writes to memory

1644

<sup>1645</sup> The Mem[] function is responsible for checking alignment and performing each memory access the instruction <sup>1646</sup> does. The ASL for Mem[] can be found in Figure 7.7. It does some alignment checks, and then calls MemSingle[] once for each single copy atomic write the access
 performs.

For example, for a fully aligned store, it calls MemSingle[] just once (lines 37 or 57), and, for a misaligned store, it will call MemSingle[] once for each byte (line 51).

<sup>1651</sup> The MemSingle[] call then performs the translation, and (if successful), the actual write to memory. Its ASL can be

<sup>1652</sup> found in Figure 7.8, with parts for extensions and store pair greyed out. On line 12, it calls AArch64. TranslateAddress

to do the translation table walk. If the translation succeeds, then the code calls PhysMemWrite (on line 40), an

<sup>1654</sup> uninterpreted function with no behaviour in ASL, which represents the actual write to memory. After perhaps

<sup>1655</sup> handling any external aborts from the write, the function returns.

```
1 Mem[bits(64) address, integer size, AccType acctype, boolean ispair] = bits(size*8)
       value_in
2
       boolean iswrite = TRUE;
3
       constant halfsize = size DIV 2;
       bits(size*8) value = value_in;
4
5
       bits(halfsize*8) lowhalf, highhalf;
6
       boolean atomic;
7
       boolean aligned;
8
10
11
       if ispair then
            // check alignment on size of element accessed, not overall access size
12
13
            aligned = AArch64.CheckAlignment(address, halfsize, acctype, iswrite);
14
       else
            aligned = AArch64.CheckAlignment(address, size, acctype, iswrite);
15
       if ispair then
16
17
           atomic = CheckAllInAlignedQuantity(address, size, 16);
18
       elsif size != 16 || !(acctype IN {AccType_VEC, AccType_VECSTREAM}) then
20
               atomic = aligned;
21
30
       if !atomic && ispair && address == Align(address, halfsize) then
31
32
            single_is_aligned = TRUE;
33
            <highhalf, lowhalf> = value;
34
            AArch64.MemSingle[address, halfsize, acctype, single_is_aligned, ispair] =
       lowhalf;
35
           AArch64.MemSingle[address + halfsize, halfsize, acctype, single_is_aligned,
       ispair] = highhalf;
36
       elsif atomic && ispair then
37
           AArch64.MemSingle[address, size, acctype, aligned, ispair] = value;
38
       elsif !atomic then
39
           assert size > 1:
40
           AArch64.MemSingle[address, 1, acctype, aligned] = value<7:0>;
41
           // For subsequent bytes it is CONSTRAINED UNPREDICTABLE whether an unaligned
42
        Device memory
           // access will generate an Alignment Fault, as to get this far means the
43
       first byte did
44
           // not, so we must be changing to a new translation page.
45
            if !aligned then
                c = ConstrainUnpredictable(Unpredictable_DEVPAGE2);
46
47
                assert c IN {Constraint_FAULT, Constraint_NONE};
                if c == Constraint_NONE then aligned = TRUE;
48
49
50
           for i = 1 to size-1
51
                AArch64.MemSingle[address+i, 1, acctype, aligned] = value<8*i+7:8*i>;
       else
           AArch64.MemSingle[address, size, acctype, aligned, ispair] = value;
57
58
       return;
59
```

```
1
     AArch64.MemSingle[bits(64) address, integer size, AccType acctype, boolean aligned
       , boolean ispair] = bits(size*8) value
2
       assert size IN {1, 2, 4, 8, 16};
3
       constant halfsize = size DIV 2;
4
7
            assert address == Align(address, size);
8
9
       AddressDescriptor memaddrdesc;
10
       iswrite = TRUE;
11
12
       memaddrdesc = AArch64.TranslateAddress(address, acctype, iswrite, aligned, size)
13
       // Check for aborts or debug exceptions
       if IsFault(memaddrdesc) then
14
15
           AArch64.Abort(address, memaddrdesc.fault);
16
17
       // Effect on exclusives
       if memaddrdesc.memattrs.shareability != Shareability_NSH then
18
19
            ClearExclusiveByAddress(memaddrdesc.paddress, ProcessorID(), size);
20
21
       // Memory array access
22
       AccessDescriptor accdesc;
23
29
            accdesc = CreateAccessDescriptor(acctype);
30
31
36
37
       PhysMemRetStatus memstatus;
38
       (atomic, splitpair) = CheckSingleAccessAttributes(address, memaddrdesc.memattrs,
        size, acctype, iswrite, aligned, ispair);
39
40
           memstatus = PhysMemWrite(memaddrdesc, size, accdesc, value);
41
            if IsFault(memstatus) then
               HandleExternalWriteAbort(memstatus, memaddrdesc, size, accdesc);
42
43
61
       return;
62
```

#### 1656 7.6.3 Translation table walks

It is the AArch64.TranslateAddress function which begins the process that performs the actual translation table walk, converting the input virtual address to the physical one. The full ASL code is too much to contain in a single figure, and so it can be found in §7.8 at the end of this chapter. This section will reference the relevant lines from the translation table walk ASL.

Figure 7.9 is an example trace of the execution of the STR Xt, [Xn] instruction, as it would happen if we were to execute it from EL1 in the EL1&0 two-stage regime. Each node represents an event in the trace (a memory or register access), and the arrows between them represent control flow. TODO: Generate from an actual isla trace rather than by hand? at least to be proper... TODO: Give labels to each event?

As described before, the instruction starts by reading the Xt and Xn registers, before beginning the call to AArch64.TranslateAddress.

The events drawn inside the dotted box come from accesses during the call to the translation table walk functions. 1667 It first calls FullTranslate (in AArch64. TranslateAddress, page 61, line 2), which calls S1Translate (in 1668 AArch64.FullTranslate, page 62, line 12), which calls S1Walk (in AArch64.S1Translate, page 63, line 29) to 1669 do the actual first-stage translation table walk. It begins by reading the relevant TTBR register to get the root 1670 table address (in AArch64.S1Walk, page 66, line 9). This is stored in a walkstate struct, which the ASL code 1671 uses to keep track of the state that changes as the walk progresses, notably, the next-level table address and 1672 any accumulated permissions. It then begins the loop to do the walk, starting from the table address read from 1673 the TTBR. On each iteration of the loop, the intermediate-physical address of the entry to be read is computed 1674 (in AArch64.S1Walk, page 66, line 38), and passed through a second stage of translation (in AArch64.S1Walk, 1675 page 66, line 47). 1676

This second stage translation calls S2Walk, which behaves similarly to the S1Walk function, taking the following steps: it reads the VTTBR (in AArch64.S2Walk, page 70, line 11); computes the (now) physical address of the entry to read (in AArch64.S2Walk, page 70, line 41); and reads it (in AArch64.S2Walk, page 70, line 44), eventually calling PhysMemRead (in AArch64.FetchDescriptor, page 72, line 23), which appears as the first R S2 L0 node in Figure 7.9.

S2Walk continues to loop, each time updating the running walkstate with the next-level table address from
the decoded descriptor (in AArch64.S2Walk, page 70, line 53), until a leaf entry is found. It is either invalid (in
AArch64.S2Walk, page 71, line 65), or, a valid page or block entry (in AArch64.S2Walk, page 71, line 70). These
correspond to the next three R S2 Ln events in the figure.

Assuming the walk did not fail with a fault, the S2Translate function returns with the physical address of the stage 1 level 0 table. S1Walk can continue, performing a read of the physical memory in the table (in AArch64.S1Walk, page 66, line 52). From there, S1Walk continues in much the same way as the stage 2 walk did: computing the current table intermediate-physical address, translating it to get the physical address, performing the read of memory to get the descriptor, until a leaf entry is found.

<sup>1691</sup> This process generates all the events up to, and including, the final stage 1 entry read (the R S1 L3 event), <sup>1692</sup> returning the intermediate-physical address that S1Walk computed.

Finally, FullTranslate calls S2Translate one last time (in AArch64. FullTranslate, page 62, line 22) on the
 intermediate-physical address, generating the last Rreg(VTTBR) and R S2 Ln events, and producing the final PA
 of the translation.

This output PA is what is passed to the PhysMemWrite of the MemSingle[] call we saw earlier, generating the final W [pa]=data event in the trace.

# 1698 **7.7** Caching in TLBs

Hardware does not simply perform the (up to) 24 additional memory accesses for every instruction-fetch, read,
or write. This would have an unacceptable performance penalty. Instead, the results of previous translations
of the same address are cached, in specialised structures called *Translation Lookaside Buffers*, or simply TLBs.
These TLBs can store whole translation results, or the separate virtual and intermediate-physical mappings, or
individual translation table entries, or a mix of the above, which we will explore more in the next chapter.

Rreg(Xt)=data					
Rreg(Xn)=va			ch64.TranslateAdd	Iress	
Rreg(TTBR) →	Rreg(VTTBR)	→ Rreg(VTTBR)	► Rreg(VTTBR)	→ Rreg(VTTBR)	→ Rreg(VTTBR)
	R S2 L0	R S2 L0	R S2 L0	R S2 L0	R S2 L0
	R S2 L1	R S2 L1	R S2 L1	R S2 L1	R S2 L1
	R S2 L2	R S2 L2	R S2 L2	R S2 L2	R S2 L2
	R S2 L3	R S2 L3	R S2 L3	R S2 L3	R S2 L3
	R S1 L0 —	R S1 L1 —	R S1 L2	R S1 L3	W [pa]=data

Figure 7.9: Memory and register accesses during a 'STR Xt, [Xn]' instruction.

When the processor translates a virtual address, it first looks for it in the TLB. If there is no entry, then this is called a *TLB miss* and a translation table walk must be performed. The results of this walk are typically then cached in the TLB, so future translations of the same address can directly grab the physical address, memory attributes, and permissions, without needing to do another translation table walk. This process and the various microarchitectural structures are explored more in §8.3.1.

<sup>1709</sup> If there is an entry, this is referred to as a *TLB hit*. In this case, the result can be taken directly from the TLB.

Under normal circumstances, the TLB is invisible to userspace programs. However, systems code is expected to manage the TLBs explicitly, using a set of instructions which Arm provide specifically for this purpose: the family of TLBI TLB-maintenance instructions. When context switching, the systems software must manually manage the TLB, invalidating stale entries for old mappings out of the cache. The behaviours that arise from reading from potentially stale TLB entries are explored in detail in §8.5.

**Address space identifiers** TLB misses and TLB maintenance are both expensive operations, and so to reduce the burden, Arm provide a mechanism to permit multiple processes' address spaces to be loaded into the TLB at the same time, by allowing the software to mark each address space with a numeric label. Arm call these *address space identifiers* (or ASIDs).

Entries in the TLB are tagged with the current ASID, and so only that process will see entries in the TLB with that ASID.

The current ASID is encoded in the high order bits of the current TTBR. During a context switch, the system software needs only switch to the new translation tables for the new address space of the other process, without doing TLB maintenance, so long as it ensures the ASIDs are distinct.

There are only finitely many ASIDs available (typically it is an 8-bit field), and so eventually TLB maintenance is required to re-use a previously allocated ASID for a new address space. But this happens far less frequently than the context switches themselves. The provided TLB maintenance instructions can target specific ASIDs, avoiding the need to over-invalidate other cached address space translations, preventing a cascade of TLB misses in other processes, further improving the runtime performance for a small amount of additional effort on the software side.

VMIDs Address space identifiers are used only for stage 1 translations. Stage 2 has virtual machine identifiers
 (VMIDs).

As before, the current VMID is encoded in the VTTBR\_EL2 register, and the TLB entries are additionally tagged
 with the current VMID (as well as the ASID), and a translation will only use TLB entries that match the current
 ASID and VMID.

1734 **TLB maintenace instructions** Arm define a whole family of instructions under the TLBI mnemonic.

<sup>1735</sup> The format for a TLBI instruction is a product of fields:

```
TLBI <type><level><broadcast>{,<reg>}
1736
1737
       <type> =
1738
          ALL | VMALL | ASID | VA{A|L} | IPAS2
1739
1740
       <level> =
         E1 | E2
1741
       <broadcast> =
1742
1743
          \{IS\}
       <reg> =
1744
         X0 | X1 | ... | X30
1745
```

Again, see the full description in the Arm manual for a more complete description [1, D5-4915].

<sup>1747</sup> The most common, and the ones that will be discussed in the following chapters, are as follows:

TLBI VAE1, Xn: Invalidate this CPU's cached copies of entries used to translate the virtual address in register
 Xn, for the EL1&0 regime, for the current ASID and VMID.

- TLBI VALE1, Xn: Invalidate this CPU's cached copies of any last-level entries used to translate the virtual address in register Xn, for the EL1&0 regime, for the current ASID and VMID.
- TLBI VAAE1, Xn: Invalidate this CPU's cached copies of any last-level entries used to translate the virtual
   address in register Xn, for the EL1&0 regime, for the current VMID, for any ASID.
- TLBI VAE1IS, Xn: Invalidate all CPU's cached copies of entries used to translate the virtual address in register Xn, for the EL1&0 regime, for the current ASID and VMID.
   (...and equivalent TLBI VAE2, TLBI VALE2, TLBI VAE2IS instructions for virtual addresses in the EL2 regime)
- TLBI IPAS2E1, Xn: Invalidate this CPU's cached copies of entries used to translate the intermediate physical address in register Xn, for the EL1&0 regime, for the current VMID.
- TLBI IPAS2LE1, Xn: Invalidate this CPU's cached copies of any last-level entries used to translate the intermediate physical address in register Xn, for the EL1&0 regime, for the current VMID.
- TLBI IPAS2E1IS, Xn: Invalidate all CPU's cached copies of entries used to translate the intermediate
   physical address in register Xn, for the EL1&0 regime, for the current VMID.
- 1764 > TLBI VMALLE1: Invalidate this CPU's cached copies of entries for the EL1&0 regime, for the current VMID.
- 1765 D TLBI VMALLEIIS: Invalidate all CPU's cached copies of entries for the EL1&0 regime, for the current VMID.
- 1766 > TLBI ALLE1: Invalidate this CPU's cached copies of entries for the EL1&0 regime, for any ASID or VMID.
- TLBI ALLE1IS: Invalidate all CPU's cached copies of entries for the EL1&0 regime, for any ASID or VMID.
   (...and equivalent TLBI ALLE2, and TLBI ALLE2IS instructions for the EL2 regime)
- TLBI ASIDE1, Xn: Invalidate this CPU's cached copies of entries for the EL1&0 regime, for the ASID specified
   in register Xn.
- TLBI ASIDE1IS, Xn: Invalidate this CPU's cached copies of entries for the EL1&0 regime, for the ASID
   specified in register Xn.
- 1773 (Note that the EL2 regime does not have ASIDs)

# 1774 **7.8 Arm ASL Reference**

Here I include the actual Arm ASL for the various parts of the translation machinery. This listing contains a
 verbatim subset of the ASL pseudocode for the translation table walk.

The sources have per-function line numbers and are annotated to direct the reader to those parts highlighted in §7.6.3. Lines which handle out-of-scope features (access flags, dirty bits, shareability domains, debugging, realms, secure states, atomics) are greyed out. Key lines have <u>coloured annotations</u>.

The ASL code listed here (minus the annotations) is copyright 2022 Arm Limited, company 02557590 registered in England. The ASL code is publicly available on Arm's webpage [55], we reproduce here only those parts of the ASL being discussed here (the translation table walk), for the purpose of critism, review, and quotation [56, s. 30].

#### 7.8.1 AArch64.TranslateAddress

### 7.8.2 AArch64.FullTranslate

```
1 AddressDescriptor AArch64.FullTranslate(bits(64) va, AccType acctype, boolean
       iswrite, boolean aligned)
2
3
     fault = NoFault();
     fault.acctype = acctype;
4
     fault.write = iswrite;
5
6
7
8
     regime = TranslationRegime(PSTATE.EL, acctype);
9
10
                                                Do the first stage of translation
     AddressDescriptor ipa;
11
     (fault, ipa) = AArch64.S1Translate tault, regime, ss, va, acctype, aligned,
12
       iswrite, ispriv);
13
     if fault.statuscode != Fault_None then ← Check for stage 1 translation fault
14
       return CreateFaultyAddressDescriptor(va, fault);
15
16
17
     if regime == Regime_EL10 && EL2Enabled() then
18
19
       s1aarch64 = TRUE;
20
       s2fs1walk = FALSE;
                                                 Do the second stage of translation
21
       AddressDescriptor pa;
       (fault, pa) = AArch64.S2Translate fault, ipa, s1aarch64, ss, s2fs1walk, acctype,
22
        aligned, iswrite, ispriv);
23
       if fault.statuscode != Fault_None then ← ____ Check for stage 2 translation fault
24
         return CreateFaultyAddressDescriptor(va, fault);
25
26
       else
27
         return pa;
28
     else
29
       return ipa;
```

7.8.3 AArch64.S1Translate

```
1 (FaultRecord, AddressDescriptor) AArch64.S1Translate(FaultRecord fault_in, Regime
       regime, SecurityState ss, bits(64) va, AccType acctype, boolean aligned_in,
       boolean iswrite_in, boolean ispriv)
2
     FaultRecord fault = fault_in;
     boolean aligned = aligned_in;
3
     boolean iswrite = iswrite_in;
4
     // Prepare fault fields in case a fault is detected
5
6
     fault.secondstage = FALSE;
     fault.s2fs1walk = FALSE;
8
9
     if !AArch64.S1Enabled(regime) then
       return AArch64.S1DisabledOutput(fault, regime, ss, va, acctype, aligned);
10
11
12
     walkparams = AArch64.GetS1TTWParams(regime, va);
13
     if (AArch64.S1InvalidTxSZ(walkparams) ||
14
         18
19
       fault.statuscode = Fault_Translation;
20
       fault.level = 0;
21
       return (fault, AddressDescriptor UNKNOWN);
22
23
     AddressDescriptor descaddress;
24
     TTWState walkstate;
25
     bits(64) descriptor;
26
     bits(64) new_desc;
27
     bits(64) mem_desc;
                                                Do the translation table walk
28
       (fault, descaddress, walkstate, descriptor) = AArch64. S1Walk(fault, walkparams,
29
       va, regime, ss, acctype, iswrite, ispriv);
30
31
       if fault.statuscode != Fault_None then ← <u>Check for S1 translation fault</u>
32
         return (fault, AddressDescriptor UNKNOWN);
33
34
37
38
       elsif AArch64.S1HasPermissionsFault(regime, ss, walkstate, walkparams, ispriv,
49
       acctype, iswrite) then ← Check for permission fault
         fault.statuscode = Fault_Permission;
50
51
52
       new_desc = descriptor;
53
56
```

63

```
57
69
70
92
93
      // Output Address
94
       oa = StageOA(va, walkparams.tgx, walkstate); ← Compute IPA
95
       MemoryAttributes memattrs;
96
112
         memattrs = walkstate.memattrs;
113
114
```

#### 7.8.4 AArch64.S1Walk

```
1 (FaultRecord, AddressDescriptor, TTWState, bits(64)) AArch64.S1Walk(FaultRecord
       fault_in, S1TTWParams walkparams, bits(64) va, Regime regime, SecurityState ss,
       AccType acctype, boolean iswrite_in, boolean ispriv)
2
     FaultRecord fault = fault_in;
     boolean iswrite = iswrite_in;
3
4
8
9
     walkstate = AArch64.S1InitialTTWState(walkparams, va, regime, ss); ← <u>read TTBR</u>
10
11
     // Detect Address Size Fault by TTB
12
     if AArch64.OAOutOfRange(walkstate, walkparams.ps, walkparams.tgx, va) then
13
       fault.statuscode = Fault_AddressSize;
14
       fault.level
                       = 0:
15
       return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64) UNKNOWN);
16
17
     bits(64) descriptor;
18
19
     AddressDescriptor walkaddress;
20
     walkaddress.vaddress = va;
21
25
       walkaddress.memattrs = walkstate.memattrs;
26
27
34
35
     DescriptorType desctype;
                     __ For each level in {0,1,2,3}
36
     repeat -
37
       fault.level = walkstate.level;
       FullAddress descaddress = AArch64.TTEntryAddress(walkstate.level, walkparams.tgx
38
       , walkparams.txsz, va, walkstate.baseaddress);
30
                                                                Get IPA of entry to read
40
       walkaddress.paddress = descaddress;
41
42
       if regime == Regime_EL10 && EL2Enabled() then
43
          s1aarch64 = TRUE;
          s2fs1walk = TRUE;
                                                                Do S2 translation to get the
44
                                                                PA of the entry
45
          aligned
                   = TRUE;
46
          iswrite
                    = FALSE;
          (s2fault, s2walkaddress) = AArch64.S2Translate (fault, walkaddress, s1aarch64,
47
       ss, s2fs1walk, AccType_TTW, aligned, iswrite, ispriv);
48
49
          if s2fault.statuscode != Fault_None then ← Check for S2 fault
            return (s2fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64)
50
       UNKNOWN);
51
52
          (fault, descriptor) = FetchDescriptor(walkparams.ee, s2walkaddress, fault);
53
                                    Read memory to get descriptor
       else
54
55
          (fault, descriptor) = FetchDescriptor(walkparams.ee, walkaddress, fault);
56
       if fault.statuscode != Fault_None then - Check for external abort
57
58
          return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64) UNKNOWN);
```

```
59
       desctype = AArch64.DecodeDescriptorType(descriptor, walkparams.ds, walkparams.
       tgx, walkstate.level);
60
61
       case desctype of
         when DescriptorType_Table
62
           walkstate = AArch64.S1NextWalkStateTable(walkstate, regime, walkparams,
63
       descriptor);
                                                      Extract next level table address
64
65
            // Detect Address Size Fault by table descriptor
            if AArch64.0AOutOfRange(walkstate, walkparams.ps, walkparams.tgx, va) then
66
              fault.statuscode = Fault_AddressSize;
67
              return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64)
68
       UNKNOWN);
69
          when DescriptorType_Page, DescriptorType_Block
70
            walkstate = AArch64.S1NextWalkStateLast(walkstate, regime, ss, walkparams,
71
       descriptor);
                                                     \Extract page start address
72
73
         when DescriptorType_Invalid
            fault.statuscode = Fault_Translation;
74
            return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64) UNKNOWN
75
       );
                      Return fault if invalid
76
77
          otherwise
78
79
            Unreachable();
80
     until desctype IN {DescriptorType_Page, DescriptorType_Block};
81
82
     // Detect Address Size Fault by final output
87
     elsif AArch64.0AOutOfRange(walkstate, walkparams.ps, walkparams.tgx, va) then
88
       fault.statuscode = Fault_AddressSize;
89
90
95
```

96 **return** (fault, walkaddress, walkstate, descriptor);

#### 7.8.5 AArch64.S2Translate

```
1 (FaultRecord, AddressDescriptor) AArch64.S2Translate(FaultRecord fault_in,
       AddressDescriptor ipa, boolean s1aarch64, SecurityState ss, boolean s2fs1walk,
       AccType acctype, boolean aligned, boolean iswrite, boolean ispriv)
2
     walkparams = AArch64.GetS2TTWParams(ss, ipa.paddress.paspace, s1aarch64);
     FaultRecord fault = fault_in;
3
4
5
     // Prepare fault fields in case a fault is detected
6
     fault.statuscode = Fault_None; // Ignore any faults from stage 1
7
     fault.secondstage = TRUE;
8
     fault.s2fs1walk = s2fs1walk;
     fault.ipaddress = ipa.paddress;
9
10
11
     if walkparams.vm != '1' then - Check if in a two-stage regime
12
       // Stage 2 translation is disabled
13
       return (fault, ipa);
14
15
     if (AArch64.S2InvalidTxSZ(walkparams, s1aarch64) ||
         AArch64.IPAIsOutOfRange(ipa.paddress.address, walkparams)) then
18
       fault.statuscode = Fault_Translation;
19
20
       fault.level
                      = 0:
21
       return (fault, AddressDescriptor UNKNOWN);
22
23
     AddressDescriptor descaddress;
24
     TTWState walkstate;
25
     bits(64) descriptor;
     bits(64) new_desc;
26
27
     bits(64) mem_desc;
28
29
       (fault, descaddress, walkstate, descriptor) = AArch64.S2Walk (fault, ipa,
       walkparams, ss, acctype, iswrite, s1aarch64);
                                                                       \<u>Do translation table walk</u>
30
       if fault.statuscode != Fault_None then ← Check for stage 2 translation fault
31
32
          return (fault, AddressDescriptor UNKNOWN);
33
34
       elsif AArch64.S2HasPermissionsFault(s2fs1walk, walkstate, ss, walkparams, ispriv
48
       , acctype, iswrite) then ← Check for stage 2 permission fault
          fault.statuscode = Fault_Permission;
49
50
51
       new_desc = descriptor;
52
55
56
```

```
68
69
74
75
77
78
    ipa_64 = ZeroExtend(ipa.paddress.address, 64);
79
    // Output Address
80
    oa = StageOA(ipa_64, walkparams.tgx, walkstate); ← Compute final PA
81
    MemoryAttributes s2_memattrs;
82
92
      s2_memattrs = walkstate.memattrs;
93
94
98
99
    MemoryAttributes memattrs;
100
      101
102
104
105
    pa = CreateAddressDescriptor(ipa.vaddress, oa, memattrs);
    106
```

#### 7.8.6 AArch64.S2Walk

```
(FaultRecord, AddressDescriptor, TTWState, bits(64)) AArch64.S2Walk(
1
2
       FaultRecord fault_in, AddressDescriptor ipa, S2TTWParams walkparams,
      SecurityState ss, AccType acctype, boolean iswrite, boolean s1aarch64)
3
     FaultRecord fault = fault_in;
4
5
     ipa_64 = ZeroExtend(ipa.paddress.address, 64);
6
7
     TTWState walkstate;
8
11
       12
13
     // Detect Address Size Fault by TTB
     if AArch64.OAOutOfRange(walkstate, walkparams.ps, walkparams.tgx, ipa_64) then
14
       fault.statuscode = Fault_AddressSize;
15
       fault.level
                    = 0;
16
       return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64) UNKNOWN);
17
18
19
     bits(64) descriptor;
20
     AddressDescriptor walkaddress;
21
22
     walkaddress.vaddress = ipa.vaddress;
23
27
       walkaddress.memattrs = walkstate.memattrs;
28
29
30
31
     DescriptorType desctype;
     32
       fault.level = walkstate.level;
33
34
35
       FullAddress descaddress;
36
40
         ipa_64 = ZeroExtend(ipa.paddress.address, 64);
         descaddress = AArch64.TTEntryAddress(walkstate.level, walkparams.tgx,
41
      walkparams.txsz, ipa_64, walkstate.baseaddress);
42
                                                     Get PA of entry to read
43
       walkaddress.paddress = descaddress;
44
45
       46
       47
         return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64) UNKNOWN);
48
       desctype = AArch64.DecodeDescriptorType(descriptor, walkparams.ds, walkparams.
49
      tgx, walkstate.level);
50
51
       case desctype of
52
         when DescriptorType_Table
          walkstate = AArch64.S2NextWalkStateTable(walkstate, walkparams, descriptor);
Extract next level table address
53
54
55
           // Detect Address Size Fault by table descriptor
56
           if AArch64.OAOutOfRange(walkstate, walkparams.ps, walkparams.tgx, ipa_64)
      then
57
            fault.statuscode = Fault_AddressSize;
58
            return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64)
      UNKNOWN);
59
         when DescriptorType_Page, DescriptorType_Block
60
```

```
61
            walkstate = AArch64.S2NextWalkStateLast(walkstate, ss, walkparams, ipa,
       descriptor);
                                                      \Extract page start address
62
          when DescriptorType_Invalid
63
            fault.statuscode = Fault_Translation;
64
            return (fault, AddressDescriptor UNKNOWN, TTWState UNKNOWN, bits(64) UNKNOWN
65
       );
                      Return fault if invalid
66
67
          otherwise
68
            Unreachable();
69
70
71
     until desctype IN {DescriptorType_Page, DescriptorType_Block};
72
     elsif AArch64.0AOutOfRange (walkstate, walkparams.ps, walkparams.tgx, ipa_64) then
79
        fault.statuscode = Fault_AddressSize; Check output address is within bounds
80
85
     return (fault, walkaddress, walkstate, descriptor);
86
```

# 7.8.7 AArch64.FetchDescriptor

```
1 (FaultRecord, bits(N)) FetchDescriptor(bit ee, AddressDescriptor walkaddress,
       FaultRecord fault_in)
2
     // 64-bit descriptors for AArch64
4
6
     bits(N) descriptor;
7
     FaultRecord fault = fault_in;
     AccessDescriptor walkacc;
8
0
10
     walkacc.acctype = AccType_TTW;
11
13
14
21
22
     PhysMemRetStatus memstatus;
23
      (memstatus, descriptor) = PhysMemRead(walkaddress, N DIV 8, walkacc);
24
      if IsFault(memstatus) then
        fault = HandleExternalTTWAbort(memstatus, fault.write, walkaddress, walkacc, N
25
       DIV 8, fault);
        if IsFault(fault.statuscode) then
26
          return (fault, bits(N) UNKNOWN);
27
28
29
31
32
      return (fault, descriptor);
```

Chapter 8

# **Relaxed virtual memory**

This chapter is based, in part, on: Relaxed virtual memory in Armv8-A [54] by Ben Simner, Alasdair Armstrong, Jean
 Pichon-Pharabod, Christopher Pulte, Richard Grisenthwaite, and Peter Sewell. Published in the proceedings of the 31st
 European Symposium on Programming (ESOP, 2022).

Now we will introduce the main concurrency architecture design questions that arise for virtual memory in Arm. As usual, the architecture defines the *envelope* of behaviours which hardware must guarantee and on which software may rely. This envelope must be tight enough to give the guarantees software needs to function, but still loose enough to admit the range of existing and conceivable microarchitectures whose optimization techniques are necessary for performance.

This chapter therefore will discuss both the relevant microarchitecture as we understand it, and also the behaviours which it is believed software relies upon. The discussion will touch on points of several kinds: some which are clear in the current Arm prose documentation; some where Arm are in the process of architecting a change; some that are not documented but where the semantics is (perhaps, after discussion with Arm) clear or constrained by current hardware or software practice; and, some where their modelling raised questions for which the architecture is not yet well-defined, and Arm must make an architectural decision.

Ideally, we would be able to specify which points belong to which kind. It is, however, not so easy. There is no clean separation between aspects there are clearly defined in the architecture reference, and those that are not; instead, the manual has a shallow covering of many of the behaviours described here. In other places, the reference may have been updated or changed over the course of the work, clarifying parts of the architecture, and while this may have happened concurrently with discussing those and other points with Arm, the reference text itself is solely the responsibility of Arm. In §8.9 we will return to this question, and more directly address the kinds of each point discussed.

**Chapter overview** The body of this chapter will explore a sequence of key behaviours, some of which the architecture guarantee and some that it does not. Each contains a description of the behaviour, including whether software relies on it or known hardware guarantees it; a short discussion of the architectural intent as we understand it; and any associated litmus tests.

This chapter will discuss a variety of interesting behaviours. In an attempt to make this chapter more approachable,
it is broken down into a logical progression: slowly building up from the most simple and fundamental parts of
the architecture, to increasingly more complex cases.

We will first discuss (in §8.2) how translation affects the prior 'data memory' **TODO: PS: user-mode**? tests covered in previous work. Then, we shall see how the caching of translation entries is limited (§??) and the fundamental behaviours of the translation table walk (§8.4). Building upon that, we will see that these translation table walks may be cached and re-used in later translations, which is explored in detail in §8.5. Then (in §8.6), we will explore how the various kinds of TLB maintenance interact with those cached translations, and other translation table walks. Finally, we touch on how all of the above fit together with system registers and other context changing and synchronising operations in §8.7.

1783

1784

1821	§8.2 Aliased data memory	
1822	§8.2.1 Virtual coherence	
1823	§8.2.2 Aliasing different locations	
1824	§8.2.3 Might be same (physical) address	
1825	§8.3 What can be cached in TLBs	
1826	§8.4 Reads not from TLB	
1827	§8.4.1 Out-of-order execution	
1828	§8.4.2 Enforcing thread-local ordering	
1829	§8.4.3 Enhanced Translation Synchronization	
1830	<b>§</b> 8.4.4 Forwarding to the translation table walker	
1831	§8.4.5 Speculative execution	95
1832	§8.4.6 Single-copy atomicity	
1833	§8.4.7 Multi-copy atomicity	
1834	<b>§</b> 8.4.8 Translation-table-walk intra-walk ordering	
1835	<b>§</b> 8.4.9 Multiple translations within a single instruction	
1836	§8.5 Caching of translations in TLBs	
1837	§8.5.1 Cached translations	
1838	§8.5.2 TLB fills	105
1839	§8.5.3 micro-TLBs	
1840	§8.5.4 Partial caching of walks	107
1841	§8.5.5 Reachability	
1842	§8.6 TLB maintenance	
1843	§8.6.1 Recovering coherence	
1844	<b>§</b> 8.6.2 Thread-local ordering and TLBI	
1845	§8.6.3 Broadcast	
1846	§8.6.4 Virtualization	
1847	§8.6.5 Break-before-make	
1848	§8.6.6 ASIDs and VMIDs	
1849	§8.6.7 Access permissions	
1850	§8.7 Context synchronisation	
1851	§8.7.1 Relaxed system registers	
1852	§8.8 Details likely to change	
1853	§8.9 Contributions	

# 1854 8.1 Virtual memory litmus tests

As previously discussed, one fundamental idea to come out of the field of
 relaxed memory is the concept of litmus tests. Virtual memory is no different,
 and exploring the architectural intent is best done through the creation, dis cussion and evaluation of small programs which are representative examples
 of common patterns.

However, as we explore more of the system semantics more and more of the
system state plays an integral role in the behaviours we see. For this reason
we need a new language for describing the state of the system, with features
not supported by the language supported by the previous litmus, rmem, herd,
and diy tools [17, 28, 14, 21], in particular, the translation table state.

The litmus tests here are given in the isla-axiomatic test format. I describe the isla tool itself, and the extended test format syntax, in more detail in **TODO: ?REF?**.

**A virtual memory litmus test** To illustrate this isla test format, Figure 8.1 contains the test listing for a non-trivial virtual memory litmus test called **CoW** (or "Copy-on-Write").

This test is derived from sequence of operations the Linux kernel takes when 1871 performing copy-on-write. Thread 0 tries to write to a location (call it x) that 1872 is currently read-only (line 1 in the thread code), then when the fault is taken 1873 the Linux exception handler begins executing (line 1 in the handler), Linux 1874 performs some checks that it's okay to copy and that it hasn't already done 1875 so (not part of the test), and then copies the physical page (lines 3 and 4 in 1876 the handler, although the test here only copies one value as demonstration), 1877 before flushing the data caches (line 5) so that later reads will be guaranteed 1878 to see the copied values. Then Linux needs to swap over the pagetable entry 1879 for x from a read-only view on the original page to a writeable mapping on 1880 the freshly copied page. It does this by first 'breaking' the entry, making 1881 it invalid (line 7), then performing the necessary TLB maintenance (line 9), 1882 before writing a new mapping to the new page (line 11). Now, Linux can 1883 return from the handler (line 13) and re-try the store instruction, hopefully 1884 this time successfully writing to the new page. 1885

1886 The test format is split into 4 main parts:

1887  $\triangleright$  The initial state, comprised of:

1880

- 1888 the per-thread register state.
  - the global memory and pagetable state.
- $_{1890}$   $\triangleright$  The thread code and any exception-handler code.
- The interesting final state, as a predicate over the final register and
   memory state.
- And, optionally, whether the outcome is allowed or forbidden by the model.

Initial state The initial state has three virtual addresses (x, y and z), and two physical addresses (pa1 and pa2).
Initial register values are written like 0:R4=z, meaning register R4 on Thread 0 initially contains the value z (in this case, a virtual address). Helper functions like pte3, page and mkdesc3 are used to get the address of the leaf entry, the page offset and to create a new valid descriptor with the given OA, a more detailed description of the functions are given later.

Behind the scenes, isla creates a full instantiation of the Arm translation tables, but with some holes for symbolic
values where the test may modify the tables. There is a default translation table, where the code and the tables
themselves are mapped by default and everything else is invalid.

<sup>1003</sup> The pagetable setup is then defined in a small DSL which defines a delta to that default table, specifying that <sup>1004</sup> certain pages should be mapped or unmapped initially, as well as being able to specify the set of locations and

AArch64	CoW
Initial State	
0:R0=0x2	
0:R1=x	
0:R3=y	
0:R4=z	
0:R5=z	
0:R6=0b0	
0:R7=pte3(x)	
0:R8=page(x)	
0:R9=mkdesc3(oa=pa2)	
0:R10=pte3(x)	
0:R20=0b0	
0:VBAR_EL1=0x1000	
0:PSTATE.EL=0b00	
virtual x y z;	
physical pa1 pa2;	
$x \mapsto pa1$ with [AP = 0b11] and default;	
$x \mapsto invalid;$	
$x \mapsto pa2$ with [AP = 0b01] and default;	
$y \mapsto pal;$	
$z \mapsto pa2;$	
identity 0x1000 with code;	
*pai = 1;	
*paz = 0; Thread 0	
Thread 0 EL1 Handler	
02 CBNZ X20 evit	
03 LDB X2 [X3]	
04 STB X2 [X4]	
05. DC CTVAC. X5	
06. DSB SY	
07. STR X6. [X7]	
08. DSB SY	
09. TLBT VALE1TS, X8	
10. DSB SY	
11. STR X9. [X10]	
12. MOV X20.#1	
13. ERET	
14. exit:	
15. MRS X21,ELR_EL1	
16. ADD X21, X21, #4	
17. MSR ELR_EL1,X21	
18. ERET	
Final State	
pa1=1 & pa2=2	
Allow	

Figure 8.1: Test CoW: code listing

75



Figure 8.2: Test CoW: execution diagram

<sup>1905</sup> their initial memory values the test will need.

Fundamentally we categorise those locations as either virtual, intermediate, or physical. The line virtual x y 1906 z in the CoW test allocates 3 virtual contiguous pages, and labels their page-aligned addresses as x, y, and z. It 1907 then allocates two physical pages with addresses pa1 and pa2. Next, the setup defines the initial value of the 1908 translation tables, as well as specifying the set of potential translation tables that may be in use by the test (for 1909 isla to create symbolic 'holes' for those). Namely, the initial state starts with x mapped to pa1 with the access 1910 permissions bits set to 0b11 (read-only). The next two lines tell isla that during the test x may become unmapped 1911 (the descriptor may be invalid), or mapped to pa2 with AP=0b01. The test also defines two other variables, y and z 1912 as aliases to the two physical pages, to help with copying the data between them, just as Linux would. Since there 1913 is an exception handler in this test, we need to ensure that the code page of the handler is mapped executable at 1914 EL1, which is what the identity 0x1000 with code line does (note that the handler section starts within the 1915 0x1000 page). Finally, we say that the initial values of pa1 and pa2 are 1 and 0 respectively. 1916

**Register translation helpers** The initial register state can reference parts of the initial state related to pagetables through the use of helper functions. Here are the helpers used by CoW, and most of the tests in this section. The full description of this format is given in **TODO: ?REF**?if more information is needed.

- pte<N>(va): The (intermediate) physical address of the level N entry in the default translation tables that
   maps va.
- b desc<N>(va): The 64-bit descriptor from the initial state of the level N entry that maps va (the value of pte<N>(va) in the initial state).
- <sup>1924</sup>  $\triangleright$  page(va): The page number that va is in (equivalently: va  $\gg$  12).
- <sup>1925</sup> b mkdesc<N>(oa=pa): A fresh 64-bit descriptor for a valid leaf entry at level N where the output address is
   <sup>1926</sup> given by the oa parameter.
- <sup>1927</sup> b mkdesc<N>(table=pa): A fresh 64-bit descriptor for a valid table entry at level N where the next-level-table
   <sup>1928</sup> address is given by the table parameter.
- <sup>1929</sup> Entries listed as f<N> mean a family of functions f1, f2, f3 and so on.

**Execution diagrams** Figure 8.2 is the isla-generated execution diagram for the CoW test. It illustrates a candidate execution which isla found (with any symbolic holes filled with concrete values) which matched the final state of the execution, and was consistent with the axioms of the model (given in Chapter 9).

<sup>1933</sup> The execution is rendered as a diagram, with separate traces for each thread, with multiple columns per thread, <sup>1934</sup> for translations and explicit events. In the diagram, there is one thread (Thread 0), and all events belong to

its trace. There are two columns; the right-hand side are the explicit events rendered in program-order, and 1935 the left-hand side contains translation events alongside any explicit events from the same instruction. Not all 1936 events from the trace are displayed in the execution diagram; many uninteresting events, of register reads and 1937 writes, and translation reads of unchanged entries, are suppressed. The execution displayed here is one where the 1938 initial store's translation table walk (event a1) reads an valid entry from the initial state but which did not have 1939 permissions to do a write, and so generates a Fault event (a2). The execution continues, copying the memory over 1940 to a new page (events b-c), before updating the translation tables to point to the new page (d-h, see §8.6.5), before 1941 returning from the exception handler (j) and re-trying the store which succeeds in writing to the new page (k2), 1942 giving a final state consistent with the expected final state from the test listing in Figure 8.1. 1943

In general, while there could be multiple executions that correspond to the final execution, the tests are usually written in a way to ensure that there is only one consistent candidate execution which corresponds to the final state. In cases where the test is forbidden by the model, we still have isla induce a concrete candidate, and render a diagram of the interesting forbidden execution.

# 1948 8.2 Aliased data memory

<sup>1949</sup> Much of the previous work on relaxed memory has been concerned with what we shall call 'data memory': the <sup>1950</sup> weak behaviour of concurrent loads and stores to memory. For Arm, we shall see that these previous models were <sup>1951</sup> implicitly assuming that all locations in the test were virtual addresses, with well-formed, constant, and injective, <sup>1952</sup> address translation mappings, which mapped all locations as readable, writable, and executable, normal cacheable <sup>1953</sup> memory.

<sup>1954</sup> Consider a non-injective mapping. Such mappings give rise to *aliasing*: the situation where two distinct virtual <sup>1955</sup> addresses in the same address space map to the same output physical address. This section will explore how the <sup>1956</sup> behaviours of those data memory tests change in the presence of aliasing.

#### 1957 8.2.1 Virtual coherence

For data memory accesses, one of the most fundamental guarantee that architectures provide is *coherence*: in any execution, for each memory location, there is a total order of the accesses to that location, consistent with the program order of each thread, with reads reading from the most recent write in that order. Hardware implementations provide this, despite their elaborate cache hierarchies and out-of-order pipelines, by a combination of coherent cache protocols and pipeline hazard checking, identifying and restarting instructions when possible coherence violations are detected.

For Arm, coherence is with respect to physical addresses [1, B2.3.1 (p157)] [1, D5.11.1 (p4931)]. This means that if two virtual addresses alias to the same physical address, then:

<sup>1966</sup> b a load from one virtual address cannot ignore a program-order previous store to the other, as seen in the
 <sup>1967</sup> following CoWR.alias test [Figure 8.3]:

AArch64	CoWR.alias
	Initial State
0:R0=0x1	
0:R1=x	
0:R3=y	
<pre>physical pa1;</pre>	
x  -> pal;	
y  -> pal;	
*pa1 = 0;	
	Thread 0
STR X0,[X1]	
LDR X2,[X3]	
	Final State
0:X2=0	
	Forbid
	Thread 0
a: ri b:	W x/pa1 = 0x1 ↓ po R y/pa1 = 0x0

This test is a variation on the standard CoWR test, where the VA is replaced with two distinct VAs, which both alias to the same PA.

The initial state is a configuration with two virtual addresses, x and y, which are both mapped to the physical address pa1, whose initial value is 0. The thread then stores 1 to x, then loads y. It is then forbidden for this load to read 0.

While the Armv8-A architecture reference manual describes data caches as being physically-indexed [1, D5.11.1 (p4931)] and so accesses via the same PA are 'fully coherent', further discussions with Arm clarify that this implies not just this coherence test, but that all prior data memory behaviours previously examined still apply when subjected to aliasing.

#### Figure 8.3: CoWR.alias test

▷ a load from one virtual address cannot ignore the write that a program-order previous load of the other 1968 address saw (CoRR0.alias+po [Figure 8.4], CoRR2.alias+po [Figure 8.5]). 1969

▷ a load from one virtual address can have its value forwarded from a store to the other, and similarly on a speculative branch (MP.alias3+rfi-data+dmb [Figure 8.6], PPOCA.alias [Figure 8.6]). 1971

1970
AArch64	CoRR0.alias+po
Initial	State
0:R0=0b1	1:R1=x
0:Rl=x	1:R3=y
	1:PSTATE.SP=0b0
	1:PSTATE.EL=0b00
<pre>physical pa1;</pre>	
x  -> pa1;	
y  -> pa1;	
*pa1 = 0;	
Thread 0	Thread 1
CTTD V0 [V1]	LDR X0, [X1]
SIR XU, [XI]	LDR X2, [X3]
Final	State
1:X0=1 & 1:X2=0	
For	·bid
Thread 0	Thread 1
a: W x/pa1 = 0x1	b: R x/pa1 = 0x1 rf po c: R y/pa1 = 0x0

This test is a variation of the data memory CoRR0 test, where one of the loads has been replaced with a load of a distinct virtual address which aliases to the same underlying physical address. Note that, like the original test, it is forbidden to read from the initial state in the later load, as this would violate coherence: exactly what the earlier text from the manual explicitly forbade.

Figure 8.4: CoRR0.alias+po test

AArch64			CoRR2.alias+p
	Initia	al State	
0:R0=0b01	1:R0=0b10	2:R1=w	3:R1=y
0:R1=u	1:R1=v	2:R3=x	3:R3=z
		2:PSTATE.SP=0b0	3:PSTATE.SP=0b0
		2:PSTATE.EL=0b00	3:PSTATE.EL=0b00
<pre>physical pa1;</pre>			
u  -> pal;			
v  -> pa1;			
w  -> pal;			
x  -> pa1;			
y  -> pa1;			
z  -> pa1;			
*pa1 = 0;			
Thread 0	Thread 1	Thread 2	Thread 3
STTR X0 [X1]	STR X0 [X1]	LDR X0, [X1]	LDR X0, [X1]
SIR X0, [XI]	Sik kö, [ki]	LDR X2, [X3]	LDR X2, [X3]
	Fina	1 State	
2:X0=1 & 2:X2=2 & 3:X	0=2 & 3:X2=1		
	Fo	orbid	
		r <u> </u>	
Thread 0	Thread 1	Thread 2	Thread 3
a: W u/pa1 = 0x1	b: W v/pa1 = 0x2	c: R w/pa1 = 0x1	e: R y/pa1 = 0x2
		rf po	DO
	$\sim$		
	$\sim$	d: K X/pal = 0x2	I: K Z/pal = 0x1
		rf	

This test is a variation of the data memory CoRR2 test. Here there are many options for adding aliasing, so we choose the maximally aliased version where each individual store and load uses a distinct virtual address but where all those virtual addresses alias to the same physical one.

This gives us a classic coherence shape, where it is forbidden for different threads to observe writes to the same physical location in different orders.

Figure 8.5: CoRR2.alias+po test



These tests are variations of the standard PPOCA and MP+rfi-data+dmb tests, but with some aliasing. Both are examples of *forwarding*: a thread-local read of a write before that write has been propagated to memory. These two tests, determined to be allowed architecturally from our discussions with Arm, show that the processor can forward from a write even if the read was for a different virtual address so long as the physical addresses match, even down a speculative path.

Figure 8.6: PPOCA.alias and MP.alias3+rfi-data+dmb tests.

# 1972 8.2.2 Aliasing different locations

In the previous section, we explored taking tests over a single location, and rewriting the test to use many locations, which all alias to the same address. One can also take a test that has multiple locations and make some of them alias to the same address.

<sup>1976</sup> Multi-location data memory tests, which are architecturally allowed, may become forbidden in the presence of

aliasing. For example, taking the traditional MP+pos test, when the two locations are aliased to the same physical

address then we get the forbidden MP.alias+pos test [Figure 8.7]. This new test is, essentially, equivalent to the

1979 old CoRR0 test: coherence with two writes and two reads to the same location; just using different aliases.

Arch64	MP.alias+pc
I	Initial State
0:R0=0x1	1:R1=y
0:R1=x	1:R3=x
0:R2=0x1	
0:R3=y	
<pre>physical pal;</pre>	
x  -> pa1;	
y  -> pa1;	
*pa1 = 0;	
Thread 0	Thread 1
STR X0, [X1]	LDR X0, [X1]
<b>STR</b> X2, [X3]	LDR X2, [X3]
	Final State
1:X0=1 & 1:X2=0	
	Forbid

Thread 0	Thread 1
a: W x/pa1 = 0x1	c: R y/pa1 = 0x1
co 🖡 po	rf po
b: W y/pa1 = 0x1	d: R x/pa1 = 0x0

Because x and y alias to the same physical address pa1, the two loads (c and d) read the same location, and so cannot read different writes out-of-order.

Figure 8.7: Test MP.alias+pos

## 1980 8.2.3 Might be same (physical) address

There is a corner case that we now should consider. For load and store instructions, when the last register used in the calculation of the address is read, the address becomes known. This allows, in the flat model, for program-order later instructions to begin execution (or at least, know they will not be restarted) at that point.

With the introduction of address translation, however, this point happens much later, after the whole translation
table walk is performed. Between the read of the register and the completion of the translation table walk, other
instructions may perform some part of their functionality. This may include reading from a different virtual
address, before the physical address of a program-order previous instruction is known, but after the virtual address
is known.

One might expect that, when deciding whether to propagate a store, if the page offset of the virtual address is different to that of the in-flight program-order earlier instructions, then the write could go ahead early, knowing that the access could not be to the same physical address as any of those instructions. However, this is not the case. Although the accesses definitely will not access the same physical address, the program-order earlier access may still fault, meaning the write will not be reached. This means that writes must wait for program-order earlier translations to finish (or at least, be known to not fault) before they can be propagated to other threads.

# <sup>1995</sup> 8.3 What can be cached in TLBs

As was described in §7.7, Arm hardware can have TLBs, caching previously seen translations. But, there are some restrictions to this; both in what information a TLB must cache when it does so, but also in what kind of information it is not permitted to cache at all.

# 1999 8.3.1 Microarchitectural TLBs

Here we must make a clear distinction between the actual *microarchitectural* translation caching one may encounter
 inspecting hardware, and the architectural model being discussed here.

While there are possibly many different ways to describe the same architectural intent, here we carefully choose one which will make building tooling, extending the model, discussions with architects, and explaining individual tests easier. We will first look at a specific example to pin down terminology and gain some intuition for hardware, before giving a model MMU and TLB that abstracts away from the details.

Microarchitectural MMU – A53 Let us explore more closely how the actual hardware fill and walk works on 2006 a modern microprocessor. The Arm Cortex A53 is an Arm-designed application class processor. Previous relaxed 2007 memory work included exercising this core design extensively during litmus testing validation of the models, 2008 finding it to be relaxed, exhibiting many relaxed behaviours, but not aggressively so. This makes the A53 a good 2009 candidate as a demonstrator of an average relaxed processor design. While other processors by Arm are more 2010 aggressive in their optimisations, the MMU and TLB layout of the A53 seems typical: other cores, such as the 2011 A57 TODO: ?CITE?, A72 TODO: ?CITE?, A76 TODO: ?CITE?, A78 TODO: ?CITE? and A715 TODO: ?CITE? 2012 all have comparable, or simpler, TLB configurations. 2013

The Arm A53 Technical Reference Manual (TRM) describes, in detail, the structure of the Memory Management Unit [57, 5-2] of the A53, and its constituent parts. Figure 8.8 shows a hand-written block diagram representing the key information from the TRM.

<sup>2017</sup> We see that each core has its own MMU, and that each MMU contains a unit that will perform the translation <sup>2018</sup> table walk, in addition to a selection of translation caching structures:

- $_{2019}$   $\triangleright$  one instruction micro-TLB;
- 2020 ▷ one data micro-TLB;
- $_{2021}$   $\triangleright$  one unified TLB;
- $_{2022}$   $\triangleright$  one walk cache; and,
- $_{2023}$   $\triangleright$  one IPA cache.

The microarchitectural TLBs store whole translations: virtual to physical mappings, plus permissions and so-on, tagged with their context. The TLBs are arranged hierarchically. With small, 10-entry, 'micro' TLBs for instruction

and data streams separately, and one large 512-entry unified TLB.



Figure 8.8: A53 Memory Management Block Diagram.

<sup>2027</sup> On a TLB miss, the MMU performs a translation table walk using the walker, computing the Arm translation <sup>2028</sup> table walk ASL code which we previously explored in §7.6.

When it begins this walk, the MMU first checks the walk cache for a matching entry. Walk cache entires are mappings from virtual address to the physical address of the last level translation table. If an entry is present the MMU can skip most of the walk entirely, performing just the very last read to read the leaf entry.

If a second stage of translation is required during the walk, the IPA cache is used (and may be, or not, used many times during the same walk). The IPA cache stores mappings from intermediate physical to physical memory – with no associated virtual address – which can be used during both the final stage 2 walk and any intermediate stage 2 walks during a stage 1 walk.

#### 2036 TODO: PS: walk cache s1 only? BS: that is one of thibaut's questions to RG

The MMU is free to save the result of any translation table walk into these structures, including for walks due to speculation, prefetching, or architectural execution. This, essentially, allows the MMU to perform a walk for any arbitrary VA or IPA, at any point in time.

#### 2040 8.3.2 Model MMU

To abstract away from any specific microarchitecture, we will model the MMU as if it were a separate asynchronous unit, one for each thread, each with an overapproximate 'TLB'.

Later, we will see tests that justify and ground this particular choice of abstraction, and we will explore this model and the mathematics which corresponds to it in more rigorous detail. But for now, we can imagine this model MMU as a set of (concurrently) executing translation table walks and a 'model TLB' cache of translation table entries.

Model TLB entries In general, the architecture permits hardware to cache whatever information from the translation process the hardware sees fit, this may include the output of whole translation table walks (complete virtual to physical mappings) or individual translation table entries, or even the result of partial walks (the address of the last-level table, for example).

It would not be feasible to even attempt to enumerate all the possible shapes of TLBs and the kinds of information they can cache. Instead, we will define a *model* TLB. This model will treat the TLB as a cache of writes of translation table entries, each tagged with some context. This allows the model to cache any combination of entries read from a translation table walk, making it weak enough to allow all known TLB implementations, but strong enough to not break any of the guarantees Arm require of those TLB implementations. These guarantees are explored, in detail, in §8.4 and §8.5.  $\begin{array}{l} \mbox{TranslationTableEntry} \equiv u64 \\ \mbox{Context} \equiv ArchContext \times Stage \times option VA \times option IPA \times PA \times Level \\ \mbox{ArchContext} \equiv VMID \times ASID \times Regime \\ \mbox{CachedTranslationTableEntry} \equiv PA \times TranslationTableEntry \times Context \\ \mbox{TLB} \equiv set CachedTranslationTableEntry \end{array}$ 

#### Figure 8.9: Model TLB type definitions.

Each entry in the model TLB contains the information about the write itself: the physical address of the entry, and the cached 64-bit entry. But it must also be tagged with some contextual information, some used during TLB lookup and some used to identify cached entries during TLB invalidation. Figure 8.9 gives a consise summary of the model TLB definition in some psuedo-type-definitions.

2061 This contextual information includes:

2065

2084

- <sup>2062</sup> b the architectural context information of the translation: the VMID, ASID (or a "global indicator"), and the translation regime;
- <sup>2064</sup> > some *extended context* information, required for implementing TLB maintenance:
  - the virtual address, intermediate physical address, and/or physical address of the translation;
- the translation stage and level at which the write was used;
- <sup>2067</sup> the system register values used in the translation (those which can be cached); and,
- <sup>2068</sup> for an entry used for a Stage 1 translation, whether it has been invalidated at both stages.

The model MMU then performs all translations by doing a full translation table walk, but being able to optionally satisfy any read during that walk from a matching entry in the model TLB which matches the architectural context and input address.

We imagine that any behaviour exhibited by a specific micro-architectural MMU and TLB configuration would also be explainable in this model.

**TLB fills** Hardware has a variety of mechanisms which may lead to a translation table walk: direct architectural execution of instructions, pre-fetching of data or instructions, and speculation down branches. These translation table walks may result in TLB misses, and those misses then result in reads from memory and the MMU 'filling' the TLB with a copy of the information it can use in future.

Arm do not wish to enumerate all the possible speculation machinery or prefetchers so instead opt for a model that is weaker: at any point in time, any thread's MMU can spontaneously perform a translation table walk for any virtual or intermediate-physical address for the current architectural context (VMID, ASID, etc, as in §8.3.2), and any reads that the translation table walk performs can either read from other TLB entries, or perform a non-TLB read of memory and then potentially cache a copy of the write it reads from in the TLB tagged with the extended context information from the walk. The behaviour of those non-TLB reads are explored more in §8.4.

# 8.3.3 Invalid entries

It is architecturally forbidden to cache information from attempted translations which result in translation faults, access flag faults, or address size faults (Note that a translation table walk may give rise to other faults as well, as discussed in §7.3.2, such as permission faults and alignment faults, which do not impose restrictions on TLB caching). More specifically, a TLB entry cannot be a write of a translation table entry which is the *direct* cause of such a fault. In particular, the TLB cannot cache translation table entries whose valid bit is not set.

This is important, as it gives software a mechanism in which it can safely update a mapping without potentially having multiple entries in the TLB for the same virtual address. These problems are described in more detail during the exploration of break-before-make in §8.6.5.

#### 2093 TODO: PS: no forward refs to tests?

84



Figure 8.10: Test CoWTf.inv+po

# 2094 8.4 Reads not from TLB

The requirement that invalid entries are not cached in the TLB gives us a way to directly observe non-TLB reads: translation table reads which result in a translation fault *must* have come from a non-TLB read.

<sup>2097</sup> We will see that these reads have some important properties that software can rely on, but that some of those <sup>2098</sup> properties will depend on certain architecture features being enabled (namely FEAT\_ETS).

In this section will we explore the properties these reads have, and the guarantees software can rely on. We shall see that these reads are affected by thread-local re-ordering, even to a greater extent than data memory reads, and the synchronization that recovers the sequential semantics. We will see how these reads from the translation table walk relate to data memory reads, with respect to coherence, multi-copy atomicity, write forwarding and so on. Finally, we will see how the FEAT\_ETS architectural feature can change the required synchronization software needs to perform.

# 2105 8.4.1 Out-of-order execution

<sup>2106</sup> First, let us consider whether reads that do not come from the TLB preserve the original program order.

po-previous writes One of the simplest questions one might ask is whether a translation-table-walk non-TLB
 read can ignore a program-order previous store.

<sup>2109</sup> This scenario is captured by the CoWTf.inv+po test [Figure 8.10]. Starting with a VA x initially invalid at level 3,

and so cannot have its level 3 entry cached in any TLB (directly or indirectly), the test then overwrites the invalid

entry with a new valid entry pointing to the physical address pa1. Program-order later, the thread then attempts to read x.

We see that the thread can take a translation fault. This fault is caused by reading an invalid entry, which was read from a stale entry in memory, ignoring the program-order previous store to the translation table entry's location.

One explanation that suffices to allow this outcome is that the instructions can be locally re-ordered; the translation table walk of the later load instruction can happen much earlier than the program-order previous store, and satisfy its read from memory first.

**po-previous reads** Similarly, the reads of a translation table walk can be locally re-ordered with respect to program-order earlier loads of the translation table entry, as demonstrated in the CoRpteTf.inv+po test [Figure 8.11].

Arch64	CoRpteTf.inv+
Ini	tial State
0:R0=desc3(y)	1:R1= <b>pte3</b> (x)
0:R1= <b>pte3</b> (x)	1:R3=x
	1:VBAR_EL1=0x1000
	1:PSTATE.SP=0b0
	1:PSTATE.EL=0b00
option default_table	es = true;
<pre>physical pa1;</pre>	
<pre>intermediate ipa1;</pre>	
x  -> invalid;	
$x \mapsto pa1;$	
y  -> pa1;	
identity 0x1000 with	n code;
*pa1 = 1;	
Thread 0	Thread 1
CMD V0 [V1]	LDR X0, [X1]
STR XU, [X1]	LDR X2, [X3]
	Thread 1 EL1 Handle
	0x1400:
	<b>MOV</b> X2,#0
	MRS X13,ELR_EL1
	ADD X13,X13,#4
	MSR ELR_EL1,X13
	ERET
Fi	nal State
1:X0=desc3(y) & 1:X2	2=0
	A 11

The translation read (event c1) can be re-ordered with respect to the program-order previous load of 13pte(x) (b), even though the load read the new translation table entry, for the same location the translation reads from.



Figure 8.11: Test CoRpteTf.inv+po

AICH04	LB.TT.inv+pos	
Initial	State	
0:R1=x	1:R1=y	
0:R2=mkdesc3(oa=pal)	1:R2=mkdesc3(oa=pa1)	
0:R3= <b>pte3</b> (y)	1:R3= <b>pte3</b> (x)	
0:VBAR_EL1=0x1000	1:VBAR_EL1=0x2000	
0:PSTATE.SP=0b0	1:PSTATE.SP=0b0	
0:PSTATE.EL=0b00	1:PSTATE.EL=0b00	
physical pal;		
x  -> invalid;		
y  -> invalid;		
$x \mapsto pa1;$		-
$y \mapsto pa1;$		The writes to
*pa1 = 1;		from propaga
identity 0x1000 with co	ode;	
identity 0x2000 with co	ode;	order earlier t
Thread 0	Thread 1	ding them fro
<b>MOV</b> X0,#0	MOV X0,#0	8
LDR X0, [X1]	LDR X0, [X1]	
<b>STR</b> X2, [X3]	<b>STR</b> X2, [X3]	
Thread 0 EL1 Handler	Thread 1 EL1 Handler	
0x1400:	0x2400:	
MRS X13,ELR_EL1	MRS X13,ELR_EL1	
ADD X13,X13,#4	ADD X13, X13, #4	
MSR ELR_EL1,X13	MSR ELR_EL1,X13	
ERET	ERET	
Final	State	
0:X0=1 & 1:X0=1		

The writes to the translation tables (b and d) are forbidden from propagating to other threads before the programorder earlier translations (a1 and c1) are satisfied, forbidding them from reading from each other's writes.



Figure 8.12: Test LB.TT.inv+pos

**po-future writes** A translation table walk read may not, in general, be re-ordered with program-order later stores.

<sup>2122</sup> This is consistent with the description in §8.2.3, as the program-order later store might not architecturally happen

if the translation table walk read were to fault. So, the later writes are speculative until the translation has finished, preventing the write from propagating until then.

This forbids both the general re-ordering of the propagation of the write to other threads (LB.TT.inv+pos [Figure 8.12]) with program-order earlier translation table walks, and, translations reading from program-order later writes (CoTW1.inv [Figure 8.13]).

### 2128 8.4.2 Enforcing thread-local ordering

Since non-TLB reads do not necessarily preserve the program order, it appears that there are no coherence guarantees one can make about them. However, by introducing some thread-local ordering constructs, we can recover some of the strong guarantees we are used to.

To force a non-TLB read to happen after some program-order earlier event we can insert the two-instruction sequence DSB SY ; ISB between them. The DSB ("Data Synchronization Barrier") waits for all loads to satisfy and for all stores to have finished and be visible to translation table walkers, before the ISB ("Instruction Synchronization Barrier") flushes the pipeline and restarts any program-order later instructions, including any translation table walks they perform.

<sup>2137</sup> **Locally-ordered-previous writes** If we introduce this sequence into the previous CoWTf.inv+po test we obtain <sup>2138</sup> the CoWTf.inv+dsb-isb test [Figure 8.14], which is forbidden by Arm. This is because the non-TLB reads, in the

AArch64 CoTW1.	inv
Initial State	
0:R1=x	
0:R2=desc3(y)	
0:R3= <b>pte3</b> (x)	
0:VBAR_EL1=0x1000	
0:PSTATE.EL=0b00	
0:PSTATE.SP=0b0	
<pre>physical pa1;</pre>	
x  -> invalid;	
$x \mapsto pa1;$	
y  -> pa1;	a1: T s1:13pte(x)
*pa1 = 1;	
identity 0x1000 with code;	b: W.0x303000/s1:
Thread 0	
LDR X0, [X1]	The store to the trans
<b>STR</b> X2, [X3]	with the program-ord
Thread 0 EL1 Handler	
0x1400:	preventing that walk
MOV X0,#0	
MRS X13, ELR_EL1	
ADD X13, X13, #4	
MSR ELR_EL1,X13	
ERET	
Final State	
0:X0=1	
Forbid	



The store to the translation table (b) cannot be re-ordered with the program-order earlier translation table walk (a), preventing that walk from reading from the store.

Figure 8.13: Test CoTW1.inv

absence of non-coherent TLB caching structures (discussed more in §8.6.1), will read from the coherent storage subsystem, and so will be required to see the new write, or something coherence after it.

Locally-ordered-previous reads If a program-order previous load has already seen some other-thread write, either through a translation (CoTTf.inv+dsb-isb [Figure 8.15]), or through a normal data load of the translation table (CoRpteTf.inv+dsb-isb [Figure 8.16]), then translation table non-TLB reads which are ordered after that read must also see that write, or a write coherence after it. These tests use the DSB; ISB sequence previously described, but any ordering to the translation table walk (described in §8.4.3) will suffice.

Microarchitecturally this is because translation table walkers are 'separate observers'. The idea is that the MMU performs reads of memory the same way any of the other observers (threads) do, meaning that those reads behave almost exactly like normal data memory reads.

This 'separate observers' principle is a reasonable model, however, we will see later on in §8.4.4 where it begins to break down.

Instruction synchronization barrier and control dependencies The ISB instruction naturally orders all translation table walks of program-order later instructions with the ISB itself. This is because the ISB effectively restarts all program-order later instructions, including any translations they do.

However, an ISB is not naturally ordered with respect to program-order *earlier* instructions. That is why in the previous tests we introduced a DSB. But a control-dependency would also work (CoTTf.inv+ctrl-isb [Figure 8.17]).

Address dependencies In previous work, address dependencies were assumed fundamental, but now we can define what an address dependency is: a register dataflow dependency into the translation table walk reads.

Address dependencies remain a strong way to order events. Arm, here and in general, avoid speculation of the values and addresses of the explicit reads and writes to memory. This means that a translation table walk will not start until after its address dataflow dependent registers are fully determined. Note, that this does not mean that pre-fetching and caching of the walk cannot happen, it's just that the architectural translation table walk must retrieve any cached values after it is known what the address will be, see **§TODO**: **?REF?**.

For non-TLB translation reads this means that a non-TLB read is locally ordered after any read whose value flows into the non-TLB read, as in CoRpteTf.inv+addr [Figure 8.18].

Memory barriers Much of the earlier work in relaxed-memory concurrency was dedicated to the behaviour of *barriers*. The Arm data memory barrier (DMB) creates ordering between memory events program-order earlier than the barrier, and memory events program-order after the barrier.

AArch	64 CoWTf.inv+dsb-isb	
	Initial State	
0:R0=	desc3(y)	
0:R1=	pte3(x)	
0:R3=	x	
0:VBA	R_EL1=0x1000	
0:PST	ATE.SP=0b0	
physi	cal pa1;	
x  ->	invalid;	
$x \mapsto j$	pal;	
у  ->	pal;	
*pal	= 1;	
ident	ity 0x1000 with code;	
Thread 0		
STR )	<0,[X1]	
DSB S	SY	
ISB		
LDR >	<2,[X3]	
	Thread 0 EL1 Handler	
0x140	00:	
MOV	<2,#0	
MRS )	<20, ELR_EL1	
ADD >	<20, X20, #4	
MSR H	ELR_EL1,X20	
ERET		
	Final State	
0:X2=	0	
	Forbid	



non-TLB read of the entry (d1) because of the intervening DSB; ISB sequence, creating local order. This ordering ensures that the non-TLB read respects the coherence order up to the point of the write a, preventing the non-TLB read from reading from a write coherence-before a.

Figure 8.14: Test CoWTf.inv+dsb-isb

AArch64	CoTTf.inv+dsb-isb
Ini	itial State
0:R0=desc3(y)	1:R1=x
0:R1=pte3(x)	1:R3=x
	1:VBAR_EL1=0x1000
	1:PSTATE.SP=0b0
	1:PSTATE.EL=0b00
<pre>physical pa1;</pre>	
<pre>x  -&gt; invalid;</pre>	
$x \mapsto pa1;$	
y  -> pa1;	
*pa1 = 1;	
identity 0x1000 with	h code;
Thread 0	Thread 1
	LDR X2, [X1]
	MOV X0, X2
STR X0, [X1]	DSB SY
	ISB
	LDR X2, [X3]
	Thread 1 EL1 Handler
	0x1400:
	<b>MOV</b> X2,#0
	MRS X13, ELR_EL1
	<b>ADD</b> X13, X13, #4
	MSR ELR_EL1,X13
	ERET
Fi	inal State
1:X0=1 & 1:X2=0	
	Forbid

The second translation-table non-TLB read of x (e1) is locally ordered after the first translation table walk (b1) because of the intervening dsb; isb sequence, and so cannot see a write coherence-before the write the earlier (b1) translation-read read from.



Figure 8.15: Test CoTTf.inv+dsb-isb

AArch64	CoRpteTf.inv+dsb-isb	
Initial State		
0:R0=desc3(y)	1:R1=pte3(x)	
0:R1=pte3(x)	1:R3=x	
	1:VBAR_EL1=0x1000	
	1:PSTATE.SP=0b0	
	1:PSTATE.EL=0b00	
option default_tables	= true;	
physical pal;		
<pre>intermediate ipa1;</pre>		
<pre>x  -&gt; invalid;</pre>		
$x \mapsto pa1;$		
y  -> pal;		
identity 0x1000 with code;		
*pa1 = 1;		
Thread 0	Thread 1	
	LDR X0, [X1]	
STTR X0 [X1]	DSB SY	
DIR AU, [AI]	ISB	
	LDR X2, [X3]	
	Thread 1 EL1 Handler	
	0x1400:	
	<b>MOV</b> X2,#0	
	MRS X13,ELR_EL1	
	MRS X13,ELR_EL1 ADD X13,X13,#4	
	MRS X13,ELR_EL1 ADD X13,X13,#4 MSR ELR_EL1,X13	
	<pre>MRS X13,ELR_EL1 ADD X13,X13,#4 MSR ELR_EL1,X13 ERET</pre>	
Fina	MRS X13, ELR_EL1 ADD X13, X13, #4 MSR ELR_EL1, X13 ERET al State	
Fina 1:x0= <b>desc3</b> (y) & 1:x2=	MRS         X13, ELR_EL1           ADD         X13, X13, #4           MSR         ELR_EL1, X13           ERET         1           State         0	

The final translation table walk of x (e1) cannot be reordered with the program-order previous load of pte3(x)(b), because of the intervening DSB; ISB sequence. The non-TLB translation read of pte3(x) (e1) therefore must read from the same write as the earlier load, or something coherence-after it.



Figure 8.16: Test CoRpteTf.inv+dsb-isb

AArch64	CoTTf.inv+ctrl-isb	
Ini	tial State	
0:R0=desc3(y)	1:R1=x	
0:R1= <b>pte3</b> (x)	1:R3=x	
	1:VBAR_EL1=0x1000	
	1:PSTATE.SP=0b0	
	1:PSTATE.EL=0b00	
<pre>physical pa1;</pre>		
x  -> invalid;		
$x \mapsto pa1;$		
y  -> pa1;		
*pa1 = 1;		
identity 0x1000 with code;		
Thread 0	Thread 1	
	<b>MOV</b> X0,#0	
	LDR X0, [X1]	
	<b>EOR</b> X4, X0, X0	
STTR X0 [X1]	CBNZ X4,LC00	
DIR RO, [RI]	LC00:	
	ISB	
	<b>MOV</b> X2,#0	
	LDR X2, [X3]	
	Thread 1 EL1 Handler	
	0x1400:	
	MRS X13,ELR_EL1	
	ADD X13,X13,#4	
	MSR ELR_EL1,X13	
	ERET	
Fi	nal State	
1:X0=1 & 1:X2=0		
Forbid		

Control-ISB locally-orders the later translation table walk (d1) after the resolution of the control flow, which happens only after the satisfaction of the read b2.



Figure 8.17: Test CoTTf.inv+ctrl-isb

AArch64	CoRpteTf.inv+addr		
Initial	Initial State		
0:R0=desc3(y)	1:R1=pte3(x)		
0:R1= <b>pte3</b> (x)	1:R3=x		
	1:VBAR_EL1=0x1000		
	1:PSTATE.SP=0b0		
	1:PSTATE.EL=0b00		
option default_tables	= true;		
<pre>physical pa1;</pre>			
<pre>intermediate ipa1;</pre>			
<pre>x  -&gt; invalid;</pre>			
$x \mapsto pa1;$			
y  -> pal;			
identity 0x1000 with c	identity 0x1000 with code;		
*pa1 = 1;			
Thread 0	Thread 1		
	LDR X0, [X1]		
STR X0, [X1]	EOR X4,X0,X0		
	LDR X2, [X3, X4]		
	Thread 1 EL1 Handler		
	0x1400:		
	MOV X2,#0		
	MRS X13,ELR_EL1		
	ADD X13, X13, #4		
	MSR ELR_EL1,X13		
	ERET		
Final	State		
1:X0= <b>desc3</b> (y) & 1:X2=0			
For	Forbid		

The address dependency from the load b to the second load, orders the reads due to the translation table walk of that load (c1) after b. Since c1 is a non-TLB read, it cannot read from a write coherence-before the write b read from.



Figure 8.18: Test CoRpteTf.inv+addr



Figure 8.19: Test CoWTf.inv+dmb

<sup>2168</sup> We will see that this applies to *explicit* memory events only: the principal reads and writes that load and store

<sup>2169</sup> instructions perform, not the implicit reads and writes they do during translations (or instruction fetching, TODO: <sup>2170</sup> ref: ifetch chapter).

Ordering of the explicit memory events does not, automatically, induce ordering between those explicit events and any reads due to translation table walks performed by those instructions. In the next subsection, we will

see how FEAT\_ETS (§8.4.3) extends the architecture to include more orderings between translations and other memory events in the same thread.

Figure 8.19 shows a simple coherence test, with a data memory barrier between a store to the translation tables and a load whose translation table walk might read from that. We can see that the barrier does not enforce that the translation table walk sees the update to the translation tables. From the previous tests, we know this means that the translation table walk happened (microarchitecturally) before the store was propagated to memory.

#### 2179 The arm DMB vs DSB instructions TODO: PS: discuss DMB v DSB

The architectural intent for DMB's ordering with respect to translation table walkers in the absence of FEAT\_ETS is still tentative, so we shall focus on the fragment with FEAT\_ETS**TODO**: ... and continue.

## 2182 8.4.3 Enhanced Translation Synchronization

#### 2183 TODO: PS: litmus tests?

Recent versions of the Arm architecture require support for FEAT\_ETS: Enhanced Translation Synchronization.
 This feature does not change the ISA, but instead, requires implementations to enforce extra ordering.

<sup>2186</sup> The Arm Architecture Reference Manual says the following [1, D5.2.5 (p4802)] :

If FEAT\_ETS is implemented, and a memory access RW1 is Ordered-before a second memory access RW2, then RW1 is also Ordered-before any translation table walk generated by RW2 that generates any of the following:

- $\triangleright$  A Translation fault.
- ▷ An Address size fault
- $\triangleright$  An Access flag fault.

2187

This prose description is a little ambiguous, and we feel, needs some clarification: The scenario being described here is a case with two instructions,  $I_1$  and  $I_2$ , each either a load or store. Imagine  $I_1$  and  $I_2$  both executing to

completion, without generating any translation, address size, or access flag faults. Then, each instruction would

93



**Figure 8.20:** ETS Ghost events example: A load instruction (I<sub>1</sub>) followed followed (in program order) by a store instruction (I<sub>2</sub>), which faults. The address dependency means that the read event  $E_{10}$  is syntactically ordered-before the (ghost) write event  $E_{20}$ , and so the read event is ordered before the reads of the translation table walk for I<sub>2</sub> read from the TLB or memory (represented by the dashed ob line).

have generated one or more explicit memory events. For example, a store might generate up to 8 separate write events (one for each byte). Call those events  $E_{ij}$  for the *j*th explicit event of instruction  $I_i$ .

Each explicit event  $E_{ij}$  would have required a translation table walk, generating translation read events which we can call  $T_{ijk}$  for the *k*th translation-table-walk read for the *j*th explicit memory event for instruction  $I_i$ .

Then, if  $I_2$  generates a translation, address size, or access flag fault, and  $E_{1n}$  would have been locally-ordered-before  $E_{2m}$  in the imagined execution without the fault (and which we can consider a kind of *ghost* event in the real

execution), and FEAT\_ETS is enabled, then,  $E_{1n}$  is locally ordered before any translation table read  $T_{2m}$  in the execution with the fault. This scenario is described pictographically in Figure 8.20.

<sup>2199</sup> The intuition here is that, microarchitecturally, on implementations that support FEAT\_ETS, when an instruction

takes an exception, the access that caused it is re-tried once the prefix of instructions is non-restartable. This

reduces *spurious aborts*: faults that come from an out-of-order read of a (what is now) stale value from memory.

**Other effects** The architecturally desired effect of FEAT\_ETS seems to be that no additional context-synchronisation should be required to prevent spurious aborts, and that simple local orderings (barriers, dependencies) should be enough. To make this so, ETS must implicitly enforce more than just the aforementioned ordering constraints.

Specifically, TLBI instructions must have stronger thread-local orderings to translation-table walks (described in more detail later); translation table walks must be (other) multi-copy atomic; and, translation table walk reads must be coherent and single-copy atomic.

non-ETS fragment There is a question here as to whether we should consider the non-ETS behaviours of the
 architecture. On the one hand, hardware in use today is from a pre-ETS version of the architecture and so we
 cannot assume that the behaviours of those devices are consistent with ETS. On the other hand, ETS is a feature
 that is widely assumed by software, even if not present on hardware.

Linux, for example, assumes implementations are ETS compatible even when they are not. Building models that capture the full extent of the non-ETS fragment would have questionable benefits as one would have to assume an ETS model when verifying software. Additionally, as ETS is becoming a mandatory feature, the concerns over non-ETS hardware will diminish over time, perhaps even by the publication of this thesis, they will be questions of the past. Finally, the semantics of this non-ETS fragment is still unclear; there are numerous questions, especially around forwarding and multi-copy atomicity generally, which are grey areas in the non-ETS fragment which Arm have yet to explicitly decide one way or another.

<sup>2219</sup> For these reasons we will assume FEAT\_ETS is present and enabled, unless explicitly stated otherwise.

**Ordering to the translation table walk** We can now define which constructs give rise to local ordering into a translation table walk. Address dependencies, and locally-ordered context-synchronisation (in particular, the DSB; ISB sequence) always give rise to ordering to the translation table walks. Control dependencies, on their own, never give rise to such ordering. If using FEAT\_ETS, then a plain DSB orders translation table walks of program-order later instructions after it. **TODO: BS: even if there's no fault?** Other barriers may give ordering to the translation table walker, if using FEAT\_ETS and the translation results in a translation fault, and those barriers would have ordered the event that would have happened otherwise.

AArch64		R.TR.inv+dmb+t	
	Initial State		
0:R0=0x2	1:R0=mkdesc3(oa=pa1)	2:R1= <b>pte3</b> (w)	
0:R1=x	1:R1= <b>pte3</b> (w)		
0:R2=0x2	1:R3=w		
0:R3= <b>pte3</b> (w)	1:VBAR_EL1=0x1000		
	1:PSTATE.SP=0b0		
	1:PSTATE.EL=0b00		
<pre>physical pa1;</pre>			
w  -> invalid;			
$w \mapsto pa1;$			
$w \mapsto raw(2);$			
x  -> pa1;			
*pa1 = 0;			
identity 0x1000 wit	h code;		
Thread 0	Thread 1	Thread 2	
<b>STR</b> X0, [X1]	<b>STR</b> X0, [X1]	TOP VO (V1)	
DMB SY	MOV X2,#1	IDR X0, [X1]	
<b>STR</b> X2, [X3]	LDR X2, [X3]		
	Thread 1 EL1 Handler		
	0x1400:		
	MRS X13,ELR_EL1		
	ADD X13,X13,#4		
	MSR ELR_EL1, X13		
	ERET		
	Final State	·	
1:X2=0 & 2:X0=2 & 2	:X2=mkdesc3(oa=pa1)		
-	Allow		

The write of the new valid entry (d) can be forwarded locally to the translation of w (e1) allowing the read of w (e2) to satisfy early. **TODO: PS: Thread2 needs explaining** 





## 2227 8.4.4 Forwarding to the translation table walker

Writes take time to propagate out to memory to other cores. One common performance optimization is *gathering*: collecting multiple writes together in a store buffer and propagating them all out together.

To maintain uniprocessor semantics, the core can read from its own store buffer, in effect, allowing it to read from writes before they've been propagated out to other cores. This behaviour is referred to as *write forwarding*.

Although the translation table walker is described as a 'separate' observer, it is also part of the core that hosts it, and is allowed to read from that core's store buffer, effectively allowing writes to be 'forwarded' to the walker, as shown in the R.T.R.inv+dmb+trfi test [Figure 8.21].

<sup>2235</sup> The simplest model here is one where non-TLB translation reads behave as a normal data memory read, reading <sup>2236</sup> either from forwarding from the store buffer, or from the coherence-latest write in the storage subsystem.

# 2237 8.4.5 Speculative execution

To facilitate fast out-of-order pipelines the machine has to begin fetching and executing the next instruction before the earlier instructions are finished. But, those instructions might control the flow of execution through the program. Executing later instructions before they are finished means that those later instructions are being executed *speculatively*: they may, if the predicted flow turns out to be incorrect, need to be discarded, **TODO: PS:** what about restarting on coherence violations? to avoid the need for rollback across threads.

When executing down a speculative path like this, there are additional restrictions that the core must adhere to. For example, stores should not be propagated out to memory, although they can still be read from by program-order later reads in the same thread.

<sup>2246</sup> Since we know reads and writes can be performed speculatively, their associated translations must also be allowed

AArch64	MP.RTf.inv+dmb+ctrl	
Initial State		
0:R0=desc3(z)	1:R1=y	
0:R1=pte3(x)	1:R3=x	
0:R2=0b1	1:VBAR_EL1=0x1000	
0:R3=y	1:PSTATE.SP=0b0	
	1:PSTATE.EL=0b00	
physical pal pa2;		
x  -> invalid;		
$x \mapsto pa1;$		
z  -> pa1;		
*pa1 = 1;		
y  -> pa2;		
identity 0x1000 with c	:ode;	
Thread 0	Thread 1	
STR X0. [X1]	LDR X0, [X1]	
DMB SY	CBNZ X0,L0	
STR X2. [X3]	L0:	
	LDR X2, [X3]	
	Thread 1 EL1 Handler	
	0x1400:	
	<b>MOV</b> X2, #0	
	MRS X13,ELR_EL1	
	ADD X13, X13, #4	
	MSR ELR_EL1,X13	
	ERET	
Final State		
1:X0=1 & 1:X2=0		
Allow		

The non-TLB read in Thread 1 (e1) is not locally ordered after the earlier load (d), despite the control dependency. This is because the processor can speculatively perform the translation table walk, before the earlier read is satisfied.



Figure 8.22: Test MP.RTf.inv+dmb+ctrl

to have been performed speculatively. This is what allows the MP.RTf.inv+dmb+ctrl test [Figure 8.22] to see an old
value for the translation table entry, as the translation can be performed speculatively. TODO: PS: If this were a
"user" test, I'd say that e1 was satisfied out-of-order w.r.t. d, not that e1 was "performed speculatively".
Or I'd expect to see a test with control-flow speculation, or argument that the second instruction is
speculative until the first is known not to fault. Are you not distinguishing between out-of-order and
speculative execution any more? TODO: BS: but speculation implies OoO?

However, forwarding from a speculative write to the translation table walker is disallowed. Since reads to read-sensitive locations (such as devices) can have side-effects, software can protect those locations by marking them as device memory in the translation tables, or leaving them unmapped altogether. A speculative write could update the translation tables arbitrarily, including allowing reads to read-sensitive locations, so it must be forbidden for a translation read to read from a still speculative write. The MP.RT.inv+dmb+ctrl-trfi test [Figure 8.23] demonstrates this, requiring that the translation table walk on the speculative path cannot read from the still-speculative store to the translation tables.

Instruction restarts A related, but separate, concept, is that of instruction restarts. In the TODO: PS: usermode? base memory model a read might be satisfied early, out-of-order with respect to program-order previous instructions, even before those instructions' accesses addresses are known. If such an earlier access turned out to be to the same address, and the later access is not a read of the same write, then the later access must be restarted to avoid coherence violations.

Translation table walk reads, while they are reads, do not do this hazard checking, and so are not required to be restarted to recover coherence. See §8.2 for more discussion on this. TODO: PS: 8.2 has a lot of stuff, point to specifics?

Arch64	MP.RT.inv+dmb+ctrl-	
	trfi	
Ini	itial State	
0:R0=0b1	1:R1=y	
0:R1=x	1:R2=mkdesc3(oa=pa1)	
0:R2=0b1	1:R3= <b>pte3</b> (w)	
0:R3=y	1:R5=w	
	1:VBAR_EL1=0x1000	
	1:PSTATE.SP=0b0	Thread 0
	1:PSTATE.EL=0b00	
<pre>physical pa1 pa2;</pre>		a: $W \times /pa1 = 0 \times 1$
w  -> invalid;		
$w \mapsto pa1;$		↓ po
x  -> pal;		b: dmb sy
*pa1 = 0;		po
y  -> pa2;		$c \cdot W w/m^2 = 0 w 1 f 1$
identity 0x1000 with	h code;	c. wy/puz owir
Thread 0	Thread 1	
	LDR X0, [X1]	
<b>STR</b> X0, [X1]	CBZ X0,LC00	
DMB SY	LC00:	The non-TLB rea
<b>STR</b> X2, [X3]	<b>STR</b> X2, [X3]	not read from a
	LDR X4, [X5]	1
	Thread 1 EL1 Handler	when on a specul
	0x1400:	after d. TODO: F
	<b>MOV</b> X4, #2	
	MRS X13,ELR_EL1	
	<b>ADD</b> X13, X13, #4	
	MSR ELR_EL1,X13	
	ERET	
Fi	nal State	
1:X0=1 & 1:X4=0		
	Forbid	



The non-TLB read of the translation table entry (f1) cannot read from a forwarded thread-local write (event e) when on a speculative path, requiring that f1 be ordered after d. **TODO: PS: manual layout this** 

Figure 8.23: Test MP.RT.inv+dmb+ctrl-trfi

# 2268 8.4.6 Single-copy atomicity

In the base memory model, there are two key guarantees on the *atomicity* of reads and writes: single-copy and multi-copy atomicity.

Recall that, single-copy atomic reads always read the maximum it can from another single-copy atomic write; in particular a 64-bit atomic never partially reads from another 64-bit atomic write.

Translation table walk reads are 64-bit single-copy-atomic reads of memory. This means that each of the reads
generated by a translation table walk will read the entire descriptor in one shot. This causes the CoWroW.inv+dsbisb test [Figure 8.24] to be forbidden, disallowing reading the output address obtained from one write, and access
permissions from another.

# 2277 8.4.7 Multi-copy atomicity

Multi-copy atomicity is a guarantee that requires any update to memory to propagate to all other threads simultaneously. This is one of the core guarantees Armv8 and RISC-V give, but earlier versions of Arm and IBM's current Power architectures do not. This has a caveat for Armv8, which is described as *other*-multi-copy atomic: threads can observe their own writes early (through write forwarding).

Microarchitecturally, a thread can only read another thread's write by reading from a global coherent storage
 subsystem. This ensures that after the thread reads from that write, any other thread must also see that write, or
 something coherence after it. While this is a property that the base model seems to have, whether it is true for
 accesses during translation table walks is a separate question.

The non-TLB reads during a translation table walk, in fact, do seem to respect this property: if one other thread has observed a write through a translation table walk then future translation table walk non-TLB reads by other threads will also observe that write (or something newer). Axiomatically, if one thread translation-reads-from a write, then all translation-table-walk reads locally-ordered after another memory event, which is itself ordered after the other thread's translation-table-walk read, will be ordered after that translation-table-walk read.

<sup>2291</sup> There are three combinations of multi-thread reads of interest, where a weaker architecture (with separate <sup>2292</sup> pagetable and data memory storage) might have mixed non-multi-copy atomic behaviours. The first of these is

AArch64 CoWroW.inv+dsb-isb		
Initial State		
0:R0=mkdesc3(oa=pa1, AP=0b11)		
0:R1=pte3(x)		
0:R2=0x1		
0:R3=x		
0:VBAR_EL1=0x1000		
0:PSTATE.SP=0b0		
<pre>physical pa1;</pre>		
<pre>x  -&gt; invalid;</pre>		
$x \mapsto$ pal with [AP = 0b11] and default;		
*pa1 = 0;		
<pre>identity 0x1000 with code;</pre>		
Thread 0		
<b>STR</b> X0, [X1]		
DSB SY		
ISB		
<b>STR</b> X2, [X3]		
Thread 0 EL1 Handler		
0x1400:		
MRS X20,ELR_EL1		
<b>ADD</b> X20, X20, #4		
MSR ELR_EL1, X20		
ERET		
Final State		
pal=1		
Forbid		

The translation table walk of the second store must read from the entire write from the earlier store, or not at all, forbidding its translation walk from reading a mix of both the initial state and the earlier write. This means there should be no way the final store can happen, as it must either be invalid or read-only.

Note that, isla does not generate candidates with non-atomic reads which are supposed to be singlecopy atomic, and so the diagram is hand-drawn **TODO: Draw it**.

Figure 8.24: Test CoWroW.inv+dsb-isb

the most basic; translation-read to translation-read, that is, the pagetable accesses are multi-copy atomic, and this is what forbids reading the old translation table value in Thread 2 in the WRC.TRTf.inv+po+dsb-isb test [Figure 8.25]. The other two are combinations of read-to-translation-read and translation-read-to-read, these show us that the translation accesses and explicit data accesses are architecturally unified: information about the memory state learned through one kind of access apply to accesses of the other. This is what forbids the following WRC.RRTf.inv+dmb+dsb-isb [Figure 8.26] and WRC.TRR.inv+po+dsb [Figure 8.27] tests, from reading the old value from memory at the end of Thread 2.

#### 2300 TODO: PS: these all need text captions

## 2301 8.4.8 Translation-table-walk intra-walk ordering

All the tests so far have been concerned with changes to at most one of the translation table entries during a single walk, however, as we saw in §7 a translation table walk may perform many reads for a single translation.

The ASL for the translation table walker performs each translation, in order, starting with the root, and ending with the leaf entry.

While reads in a thread can be re-ordered, translation-reads within a translation table walk cannot, as this would require the hardware to do value speculation on the next-level table address, and as discussed in §8.4.5 reading from speculative values in a translation table walk is generally forbidden.

Requiring the translation reads from a translation table walk to be satisfied in translation walk order has an observable effect, for example in the following ROT.inv+dsb test [Figure 8.28] the translation table walk of the read in Thread 1 must see the writes to the translation table done by Thread 0 in the order they were propagated out to memory, and so reading from the old level 3 entry is forbidden.

## 2313 8.4.9 Multiple translations within a single instruction

Some instructions generate multiple explicit memory events, such as for the load pair and store pair instructions, or misaligned accesses, or potentially some read-modify-writes. When there are multiple explicit memory events, there will be a dedicated translation for each of them, with its own translation table walk.

Here the architecture as it is written today is overly sequentialised: the ASL for these cases performs each translation (and the respective access) in some order, but the architectural intent is that the separate translations should be unordered with respect to each other.

Arch64	W	RC.TRTf.inv+po+dsb-is
	Initial State	
0:R0= <b>desc3</b> (z)	1:R1=x	2:R1=y
0:R1=pte3(x)	1:R2=0b1	2:R3=x
0:PSTATE.EL=0b00	1:R3=y	2:VBAR_EL1=0x2000
0:PSTATE.SP=0b0	1:VBAR_EL1=0x1000	2:PSTATE.SP=0b0
	1:PSTATE.EL=0b00	2:PSTATE.EL=0b00
	1:PSTATE.SP=0b0	
physical pal pa2;		
x  -> invalid;		
$x \mapsto pa1;$		
z  -> pa1;		
*pa1 = 1;		
y  -> pa2;		
identity 0x1000 with	n code;	
identity 0x2000 with	n code;	
Thread 0	Thread 1	Thread 2
		LDR X0, [X1]
CMD VO (V1)	LDR X0, [X1]	DSB SY
SIR AU, [AI]	<b>STR</b> X2, [X3]	ISB
		LDR X2, [X3]
	Thread 1 EL1 Handler	Thread 2 EL1 Handle
	0x1400:	0x2400:
	<b>MOV</b> X0,#0	<b>MOV</b> X2,#0
	MRS X13,ELR_EL1	MRS X13, ELR_EL1
	ADD X13, X13, #4	
	MSR ELR_EL1, X13	MSR ELR_EL1, X13
ERET		ERET
	Final State	
1:X0=1 & 2:X0=1 & 2:	x2=0	
	Forbid	



Figure 8.25: Test WRC.TRTf.inv+po+dsb-isb





d: W y/pa2 = 0x1

g: isb

h1: T s1:13pte(x)  $\rightarrow$  h2: Fault (R)

• po

↓po i: eret

 $\operatorname{trf}$ 

in in child i		wne.rnn.mv+po+usi
	Initial State	
0:R0=mkdesc3(oa=pal)	1:R0=0b0	2:R1=y
0:R1= <b>pte3</b> (x)	1:R1=x	2:R3= <b>pte3</b> (x)
0:PSTATE.EL=0b00	1:R2=0b1	2:PSTATE.SP=0b0
0:PSTATE.SP=0b0	1:R3=y	2:PSTATE.EL=0b00
	1:VBAR_EL1=0x1000	
	1:PSTATE.EL=0b00	
	1:PSTATE.SP=0b0	
physical pal pa2;		
<pre>x  -&gt; invalid;</pre>		
$x \mapsto pa1;$		
y  -> pa2;		
*pa1 = 1;		TODO: DS: 10
identity 0x1000 with c	ode;	10D0: F3: w
Thread 0	Thread 1	Thread 2
		LDR X0, [X1]
<b>STR</b> X0, [X1]	LDR X0, [X1]	DSB SY
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3]	DSB SY LDR X2, [X3]
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler	DSB SY LDR X2, [X3]
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler 0x1400:	DSB SY LDR X2, [X3]
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler 0x1400: MRS X13, ELR_EL1	DSB SY LDR X2, [X3]
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler 0x1400: MRS X13, ELR_EL1 ADD X13, X13, #4	DSB SY LDR X2, [X3]
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler 0x1400: MRS X13, ELR_EL1 ADD X13, X13, #4 MSR ELR_EL1, X13	DSB SY LDR X2, [X3]
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler 0x1400: MRS X13, ELR_EL1 ADD X13, X13, #4 MSR ELR_EL1, X13 ERET	DSB SY LDR X2, [X3]
<b>STR</b> X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler 0x1400: MRS X13, ELR_EL1 ADD X13, X13, #4 MSR ELR_EL1, X13 ERET Final State	DSB SY LDR X2, [X3]
STR X0, [X1]	LDR X0, [X1] STR X2, [X3] Thread 1 EL1 Handler 0x1400: MRS X13, ELR_EL1 ADD X13, X13, #4 MSR ELR_EL1, X13 ERET Final State 2=0	DSB SY LDR X2, [X3]



Figure 8.27: Test WRC.TRR.inv+po+dsb

Arch64	ROT.inv+ds
Initia	l State
0:R0=mkdesc3(oa=ipal)	1:R1=x
0:R1= <b>pte3</b> (x, new_table	) 1:VBAR_EL1=0x1000
0:R2=mkdesc2(table=0x2	.83000)
0:R3= <b>pte2</b> (x)	
0:PSTATE.EL=0b01	
physical pal;	
<pre>intermediate ipa1;</pre>	
assert pal == ipal;	
ipal  -> pal;	
x  -> invalid at level	2;
$x \mapsto table(0x283000) a$	t level 2;
<pre>sltable new_table 0x28</pre>	0000 {
<pre>x  -&gt; invalid;</pre>	
$x \mapsto ipal;$	
};	
identity 0x1000 with c	ode;
Thread 0	Thread 1
<b>STR</b> X0, [X1]	
DSB SY	LDR X0, [X1]
<b>STR</b> X2, [X3]	
	Thread 1 EL1 Handler
	0x1400:
	// read ESR_EL1.ISS, to s
	MRS X14,ESR_EL1
	<b>AND</b> X14, X14, #0b111
	CMP X14,#0b111
	MOV X17,#1
	MOV X18,#2
	// if ESR_EL1.ISS.DFSC =
	CSEL X0,X17,X18,eq
	// advance ELR
	MRS X13,ELR_EL1
	ADD X13,X13,#4
	MSR ELR_EL1, X13
	// return
	ERET
Final	State
1. VO-1	
1:X0=1	

The translation-table walk from the read of x in Thread 1 must perform its translation non-TLB reads in the order they appear in the walk, forbidding reading from the new level 2 table entry in d1, but then reading the stale initial value for that entry from memory.

The test listing contains some concrete values to make it executable in isla, namely fixing the location of the new table at 0x280000 so it's not symbolic, and the exact location of the level 3 entry within the new table will be at 0x283000 (known from the fixed isla configuration). Whether the exception comes from the level 2 or the level 3 entry can be determined by reading the ISS field of the ESR\_EL1 register, which the exception handler does.



Figure 8.28: Test ROT.inv+dsb

Misaligned accesses, and the load and store pair instructions, should generate explicit memory events and associated translations which are unordered with respect to each other.

2322 TODO: PS: litmus test with misalgined?

AArch64	CoWinvT+po	
I	initial State	
0:R0=0b0		
0:R1= <b>pte3</b> (x)		
0:R3=x		
0:VBAR_EL1=0x10	000	
0:PSTATE.SP=0b0	)	
physical pal pa	12;	Thread 0
x  -> pa1;		
$x \ \mapsto \ \text{invalid};$		a: $W 0x303000/s1:13pte(x) = 0x0$
identity 0x1000	with code;	, trf
	Thread 0	
STR X0, [X1]		b1: T s1:13pte(x) $\rightarrow$ b2: R x/pa1 = 0x0
LDR X2, [X3]		The translation read (b1) of the last level entry for y con
Threa	d 0 EL1 Handler	The translation read (D1) of the last-level entry for x can
0x1400:		be re-ordered with respect to the program-order earlier
MOV X2,#1		store (a) to $nte3(x)$
MRS X13, ELR_E	L1	store (a) to pees(x).
ADD X13, X13, #	4	
MSR ELR_EL1,X	13	
ERET		
	Final State	
0:X2=0		
	Allow	

Figure 8.29: Test CoWinvT+po

# **8.5** Caching of translations in TLBs

We have seen in §8.4 that, while non-TLB reads do not necessarily preserve the program-order without additional synchronisation due to the out-of-order execution of instructions, those translation table reads get satisfied from the coherent storage subsystem or from forwarding from earlier stores, much like the normal explicit data reads do. This section will explore what happens when translation table walk reads may instead be satisfied from the TLB.

<sup>2329</sup> Unfortunately for the programmer, the TLB need not be coherent with memory: it can have stale values. This <sup>2330</sup> section explores the behaviours that arise from this caching of stale values.

# 2331 8.5.1 Cached translations

In the previous section we carefully constructed tests which began with an initially invalid translation, to avoid
 TLB caching issues. Here, we will generally start with entries that are valid, and so might be present in the TLB.

The following CoWinvT+po test [Figure 8.29] begins with an *initially valid* (and therefore potentially initially *cached in the TLB*) translation for the virtual address x. It then updates the last-level translation table entry for x, setting it to 0, making it invalid (and thus unmapping x). Then, program order later, the same thread tries to read X.

The read can succeed, as its translation can read from the old value from memory. We saw earlier that translation table walks can be re-ordered with respect to program order, but even inserting thread-local ordering to the translation, such as in test CoWinvT+dsb-isb [Figure 8.30], does not forbid it.

AArch64	CoWinvT+dsb-isb	
Initi	al State	
0:R0=0b0		
0:R1= <b>pte3</b> (x)		
0:R3=x		
0:VBAR_EL1=0x1000		
0:PSTATE.SP=0b0		
<pre>physical pa1 pa2;</pre>		
x  -> pal;		
$x \mapsto invalid;$		
identity 0x1000 wi	th code;	
Th	read 0	
STR X0, [X1]		
DSB SY		
ISB		
LDR X2, [X3]		
Thread 0	EL1 Handler	
0x1400:		
MOV X2,#1		
MRS X20,ELR_EL1		
<b>ADD</b> X20, X20, #4		
MSR ELR_EL1, X20		
ERET		
Final State		
0:X2=0		
А	allow	

	Thread 0
	a: W 0x303000/s1:l3pte(x) = 0x0
	po
	b: dsb sy
	po
trf	c: isb
d1: T s1:13	d2: R x/pa1 = 0x0

The translation read (d1) of the last-level entry for x is required to be satisfied after the earlier store (a) to the entry's location because of the intervening dsb sy; isb sequence, but can be satisfied from a cached value in the TLB, allowing d1 to read from a stale value.



•

2341 8.5.2 TLB fills

Translation table walks can be requested by the core in two different ways: (1) through the architectural execution of an instruction; or, (2) from a spontaneous translation table walk (for example, due to speculation and prefetching of data or instructions). In either case, the result of that walk can be cached in the TLB and recalled for other translation table walks.

Architecturally a TLB fill is no different to a normal translation table walk; each fill originates from a non-TLB
 read, with all the behaviours described in the previous sections. Later translation table walks are allowed, however,
 to recall an earlier value and then reuse that rather than doing a fresh read.

Spontaneous walks The hardware may, at any time, try to prefetch or speculatively read some address.
Architecturally these appear as spontaneous translation table walks. Those spontaneous walks may be cached.
We can see this occuring in the following MP.RT.inv+poloc-dmb+ctrl-isb test [Figure 8.31], where a spontaneous translation and the resulting TLB fill allows a future translation table walk to see a stale value.

Speculative paths Since translation table walks, and therefore TLB fills from the result of those walks, can
 happen at any point, there is no need to consider TLB fills of architectural translation table walks down speculative
 paths as any such behaviour is subsumed by a spontaneous fill.

However, as described earlier, we saw that writes cannot be forwarded to translation table walks when down
speculative paths (§8.4.5), as this would lead to security violations. This naturally excludes TLB fills of still
speculative writes; since a speculative write cannot be used in the result of a translation table walk, it cannot end
up cached in a TLB.

# 2360 8.5.3 micro-TLBs

So far we have spoken as if entries are, at any particular moment in time, either present in the TLB or not. Hardware, however, may have multiple *micro*-TLBs for the same thread, each with their own potential cached entry for the same address.

In effect, these micro-TLBs behave as if they were a larger non-deterministic TLB with potentially many values for each entry. The presence of these smaller caching structures in a superscalar machine means that different instructions may be accessing different TLBs at the same time. This allows later instructions to 'skip' over a previously seen cached entry, and then see it again later.

<sup>2368</sup> These effects can be seen in the CoTfT+dsb-isb test [Figure 8.32], where the presence of these micro-TLBs (or <sup>2369</sup> other distributed caching structures) allow later events (even locally-ordered later) to see old cached entries after

AArch64	MP.RT.inv+poloc
	dmb+ctrl-isl
Initia	d State
0:R0=mkdesc3(oa=pal)	1:R1=y
0:R1= <b>pte3</b> (x)	1:R3=x
0:R2=0b0	1:VBAR_EL1=0x1000
0:R3= <b>pte3</b> (x)	1:PSTATE.SP=0b0
0:R4=0b1	1:PSTATE.EL=0b00
0:R5=y	
<pre>physical pa1 pa2;</pre>	
x  -> invalid;	
$x \mapsto pa1;$	
y  -> pa2;	
*pa1 = 0;	
*pa2 = 0;	
identity 0x1000 with c	code;
Thread 0	Thread 1
	LDR X0, [X1]
STR X0, [X1]	CBNZ X0,L0
<b>STR</b> X2, [X3]	LO:
DMB SY	ISB
<b>STR</b> X4, [X5]	MOV X2,#1
	LDR X2, [X3]
	Thread 1 EL1 Handler
	0x1400:
	MRS X13,ELR_EL1
	ADD X13, X13, #4
	MSR ELR_EL1,X13
	ERET
Fina	l State
1:X0=1 & 1:X2=0	

A spontaneous walk and fill can happen on Thread 1 after the write of the valid entry to pte3(x) (a), but before the immediate re-invalidation of that entry (b), allowing the later translation table walk to see the old cached entry (g1), even though the architectural translation table walk could not have happened while the valid entry was visible.



Figure 8.31: Test MP.RT.inv+poloc-dmb+ctrl-isb

AArch64	CoTfT+dsb-isb	
Initial State		
0:R0=0b0	1:R1=x	
0:R1= <b>pte3</b> (x)	1:R3=x	
	1:VBAR_EL1=0x1000	
	1:PSTATE.SP=0b0	
	1:PSTATE.EL=0b00	
<pre>physical pa1;</pre>		
x  -> pa1;		
$x \mapsto \text{invalid};$		
y  -> pal;		
*pa1 = 0;		
identity 0x1000 with	code;	
Thread 0	Thread 1	
	LDR X2, [X1]	
	MOV X0, X2	
STR X0, [X1]	DSB SY	
	ISB	
	LDR X2, [X3]	
	Thread 1 EL1 Handler	
	0x1400:	
	MOV X2,#1	
	MRS X13,ELR_EL1	
	<b>ADD</b> X13,X13,#4	
	MSR ELR_EL1,X13	
	ERET	
Final State		
1:X0=1 & 1:X2=0		
	Allow	



The earlier translation read (b1) reads from the new invalid entry, reading from memory (as it cannot have been in the TLB), but a later translation read (f1) of the same location can still potentially see a stale cached entry.

Figure 8.32: Test CoTfT+dsb-isb

2370 earlier events witnessed a TLB miss.

**Break-before-make and restrictions** We will see later that the ability to have multiple cached entries for a single address causes problems for software managing coherence, and imposes extra restrictions on software practice. This means that, in general, the effects of the micro-TLBs are restricted to only those combinations that do not cause *break-before-make violations* (see §8.6.5).

## 2375 8.5.4 Partial caching of walks

TLBs need not cache entire virtual to physical translations. Instead, they are free to cache any subset of the reads from the walk separately.

Caching up to last-level table The most common kind of caching structure we are aware of in microarchitecture is the walk cache (see §8.3.1). Traditionally a TLB would store entire virtual to physical mappings, making it fast to lookup the translation (often a single cycle), but there was limited space and induced heavy burden on a TLB miss or TLB invalidation. Walk caches store the last-level table entry, allowing TLB invalidation of leaf entries and TLB misses to re-use a prefix of the walk and perform a minimal number of accesses. This can be seen in the MP.RTT.inv3+dmb-dmb+dsb-isb test [Figure 8.33], where a walk cache could allow the table entry to be cached separately from the last-level entry, allowing the last translation read to read from a much newer value.

Caching of whole translation A common configuration for the TLB is to cache whole translation walks, from virtual to physical. This kind of caching has an important caveat: there is no requirement for the TLB to remember the intermediate physical address of any stage 2 translations that were done during the walk, including the final stage 2 walk of the access address itself. This means that TLB invalidations by IPA might not remove all the cached data associated with a cached entry for that IPA, if there is a whole cached translation which used that entry. TODO: ?REF?.

Independent caching of IPAs In a two-stage regime, the virtual addresses are first translated into intermediate
 physical address. The secondary translations based on the intermediate physical addresses, either of the final
 output address or of any of the intermediate table addresses, may be cached in the TLB without remembering the
 originating virtual address.

<sup>2395</sup> This means these cached translations may be recalled for translations of different virtual addresses.

AArch64	MP.RTT.inv3+dmb	
	dmb+dsb-ish	
Initial State		
0:R0=0b0	1:R1=y	
0:R1= <b>pte2</b> (x)	1:R3=x	
0:R2=mkdesc3(oa=pal)	1:VBAR_EL1=0x1000	
0:R3= <b>pte3</b> (x)	1:PSTATE.SP=0b0	
0:R4=0b1	1:PSTATE.EL=0b00	
0:R5=y		
virtual x y;		
<pre>physical pa1 pa2;</pre>		
assert x[4821] = $y[4821]$	4821];!	
x  -> invalid;		
$x \mapsto pa1;$		
$x \ \mapsto \ \text{invalid} \ \text{at level}$	2;	
y  -> pa2;		
*pa1 = 0;		
*pa2 = 0;		
identity 0x1000 with a	code;	
Thread 0	Thread 1	
<b>STR</b> X0, [X1]	LDR X0, [X1]	
DMB SY	DSB SY	
<b>STR</b> X2, [X3]	ISB	
DMB SY	MOV X2,#1	
<b>STR</b> X4, [X5]	LDR X2, [X3]	
	Thread 1 EL1 Handler	
	0x1400:	
	MRS X13,ELR_EL1	
	ADD X13, X13, #4	
	MSR ELR_EL1, X13	
	ERET	
Fina	l State	
1:X0=1 & 1:X2=0		
Allow		

The translation-read of the level 2 entry for x (i1) can read from stale writes from a translation that the subsequent level 3 translation-read (i2) does not read from, as the level 2 entry could have been cached in the 'TLB' (in this case, a co-located 'walk cache' structure), while the level 3 entry gets read from memory. **TODO: PS: explain the magic numbers and tfr edge a bit.** 



Figure 8.33: Test MP.RTT.inv3+dmb-dmb+dsb-isb

In addition, pre-fetching may perform translations of arbitrary IPAs. This means that any cached translations
 might not correspond to any valid whole translation table walk, but may still be used during such walks.

This is most clear in ROT.invs1+dmb2 [Figure 8.34], where, although the IPA was never reachable from the stage 1 translations, the old IPA to PA mapping was cached and used later.

Caching of individual entries Architecturally, Arm wish to allow many more implementations of TLBs and
 translation caching structures than currently known hardware contains.

The weakest variation on this is allowing each individual translation table entry to be cached separately and independently.

One could construct litmus tests for each of the possible combination of translation table entries, but there would
be overwhelmingly many of these, or even a 'most relaxed' version where every translation table entry comes
from different previous translations, but this would be too large to show here. So, for simplicity I show just one of
them here, ROT.inv2+dmb [Figure 8.35]; where the last-level entry came from a newer value than the previous
levels.

## 2409 8.5.5 Reachability

One key property that the TLB must have is that it may only cache translation table entries which are *reachable* from a translation in the current context. That is, it can only cache an entry which is the result of a valid translation table walk, either using values from memory or other valid translation table entries from the TLB, using the current system register state.

This means that writes coherence-before the most recent write at the time a translation table entry location becomes reachable are not visible to the walker, and cannot have been cached in any TLB.

This is captured in the RUE+isb [Figure 8.36] ("Read-unreachable-entry") test, which is forbidden as the write to the translation table from before the time the location becomes reachable by translation table walkers cannot have been cached in any TLBs, or read from by any spontaneous walks.

<sup>2419</sup> This area is currently under discussion with Arm.

# 2420 8.6 TLB maintenance

Recovering coherence for translation reads in the presence of TLB caching can be achieved through the use of TLB *maintenance* instructions: namely the TLBI ("TLB invalidate" instruction).

TLB maintenance generally causes two microarchitectural effects: erasing stale entries from the TLB, ensuring future TLB fills (for example, due to a translation read) will see the coherent value from memory; and, discarding any partially executed instructions, on other cores, which had already begun execution using a stale entry but had not yet finished executing. We will explore both of these effects and the subtle interaction with other parts of the virtual memory systems architecture in more detail throughout this section.

### 2428 8.6.1 Recovering coherence

We saw earlier, in Section 8.5.1, that stale values cached in the TLB can cause coherence violations in the translation, for example, in the CoWinvT+dsb-isb test [Figure 8.30]. By inserting the correct TLBI sequence into that test, we produce a new test, CoWinvT.EL1+dsb-tlbi-dsb-isb [Figure 8.37], which is now forbidden.

<sup>2432</sup> There are many flavours of TLBI that could have been inserted into this test, the one in the figure is TLBI VAE1, or,

TLB invalidation by virtual address, for the EL1&0 translation regime. Using a TLBI-by-VA means the programmer has to provide the virtual page to invalidate, and the TLBI only affects addresses for that specific invalidated entry, not all of them.

<sup>2436</sup> Using the incorrect TLBI leads to insufficient invalidation occuring. For example, if in the aforementioned <sup>2437</sup> CoWinvT.EL1+dsb-tlbi-dsb-isb the TLBI had the wrong page, then it would have no effect and the test would

2438 remain allowed.

AArch64	ROT.invs1+dmb
Initia	l State
0:R0=mkdesc3(oa=pal)	1:R1=x
0:R1= <b>pte3</b> (x, s2_table)	1:VBAR_EL1=0x1000
0:R2=0b0	1:VBAR_EL2=0x2000
0:R3= <b>pte3</b> (x, s2_table)	
0:R4=mkdesc3(oa=ipal)	
0:R5= <b>pte3</b> (x)	
0:PSTATE.EL=0b01	
<pre>physical pa1;</pre>	
<pre>intermediate ipa1;</pre>	
x $\mid - >$ invalid at level	2;
$x \mapsto ipa1;$	
ipal  -> pal;	
ipal $\mapsto$ invalid;	
*pa1 = 1;	
identity 0x1000 with c	ode;
identity 0x2000 with c	ode;
Thread 0	Thread 1
<b>STR</b> X0, [X1]	
DMB SY	MOV X0, #0
<b>STR</b> X2, [X3]	LDR X0, [X1]
DMB SY	
<b>STR</b> X4, [X5]	
	Thread I ELI Handler
	0x1400:
	MRS X13, ELR_EL1
	<b>ADD</b> X13, X13, #4
	MSR ELR_EL1,X13
	ERET
	Inread I EL2 Handler
	Ux2400:
	MRS X13, ELR_EL2
	ADD X13, X13, #4
	MSK ELR_EL2,X13
E:	ERET Ctata
Final State	
1:XU=1	1
Al	low

The translation read of the stage 2 leaf entry for x (f2) can read from an old cached version, from the write (a) even though it was not reachable by any translation table walk for any VA, as the IPA it maps was not mapped by any stage 1 tables before it was overwritten by (b). This test relies on translation table walks being naturally ordered (by iio), see §8.4.8.



Figure 8.34: Test ROT.invs1+dmb2

AArch64	ROT.inv2+dmb	
Initial	State	
0:R0=0b0	1:R1=x	
0:R1= <b>pte3</b> (x, new_table) 1:VBAR_EL1=0x1000		
0:R2=mkdesc2(table=0x283000)		
0:R3= <b>pte2</b> (x)		
0:PSTATE.EL=0b01		
physical pal;		
<pre>intermediate ipal;</pre>		
assert pal == ipal;		
ipal  -> pal;		
<pre>x  -&gt; invalid at level 2;</pre>		
$x \mapsto table(0x283000)$ at level 2;		
<pre>sltable new_table 0x280000 {</pre>		
x  -> ipal;		
$x \mapsto invalid;$		
};		
identity 0x1000 with c	ode;	
Thread 0	Thread 1	
STR X0, [X1]	MOV X0, #1	
DMB SY	LDR X0, [X1]	
<b>STR</b> X2, [X3]		
	Thread I ELI Handler	
	0x1400:	
	MRS X13,ELR_EL1	
	ADD X13, X13, #4	
	MSR ELR_EL1,X13	
T_ 1	ERET	
Final State		
1:X0=0		
Allow		

The translation-read of the level 3 entry (d2) can read from a stale cached translation, which was cached before the write to the level 2 entry (c). Note that this test assumes that the original new\_table was reachable (and therefore could be cached) before the write c. See §8.5.5 for a discussion on this.



Figure 8.35: Test ROT.inv2+dmb

AArch64	RUE+isb	
Initial State		
0:R0=mkdesc3(oa=pa1)		
0:R1=0x0		
0:R2= <b>pte3</b> (x, new_table)		
0:R3=ttbr(asid=0x01, base=r	ew_table)	
0:R4=x		
0:VBAR_EL1=0x1000		
0:PSTATE.EL=0b01		
0:PSTATE.SP=0b1		
<pre>intermediate ipa1;</pre>		
physical pal;		
*pa1 = 0;		
	,	
sitable new_table 0x2C0000	{	
identity 0x1000 with code;		
x  -> invalid;		
1,		
identity 0x1000 with code;		
Thread 0		
01. STR X0, [X2]		
02. STR X1, [X2]		
03. MSR TTBR0_EL1,X3		
04. ISB		
05. MOV X1,#1		
06. LDR X1, [X4]		
Thread 0 EL1 Handle	er	
01. 0x1200:		
02. MRS X20,ELR_EL1		
03. ADD X20,X20,#4		
04. MSR ELR_EL1,X20		
05. ERET		
Final State		
0:X1=0		



The write to the new\_table translation table entry for x (a) is not visible at the point of the change of TTBR (c), and so the later translation table walk (e1) cannot read from it.

Note that isla currently does not do any kind of reachability analysis, and so does not forbid this test.

Figure 8.36: Test RUE+isb

AArch64	CoWinvT.EL1+dsb-	
	tlbi-dsb-isb	
Initial	State	
0:R0=0b0		
0:R1= <b>pte3</b> (x)		
0:R3=x		
0:R5= <b>page</b> (x)		
0:VBAR_EL1=0x1000		
0:PSTATE.EL=0b01		
<pre>physical pa1 pa2;</pre>		
x  -> pal;		
$x \mapsto invalid;$		
<pre>identity 0x1000 with code;</pre>		
Three	ad 0	
STR X0, [X1]		
DSB SY		
TLBI VAE1,X5		
DSB SY		
ISB		
LDR X2, [X3]		
Thread 0 EL1 Handler		
0x1000:		
MOV X2,#1		
MRS X13,ELR_EL1		
ADD X13, X13, #4		
MSR ELR_EL1, X13		
ERET	~	
Final State		
0:X2=0		
Fort	oid	

The read of the translation table entry for x (f1) is required to happen after the earlier store (a), because of the intervening dsb sy; isb sequence (d and e), and cannot be satisfied from the TLB, because of the TLBI (c), forbidding it from still seeing a stale value. Note that TLBI instructions can only be executed from EL1, so this test starts execution at EL1 rather than the usual default of EL0.

Figure 8.37: Test CoWinvT.EL1+dsb-tlbi-dsb-isb





The TLBI (b) can be re-ordered with program-order earlier events, due to the lack of DSBs ordering it after them, allowing the store (a) to happen later, letting the final translation read (e1) still see the old stale translation.

Figure 8.38: Test CoWinvT.EL1+tlbi-dsb-isb

#### 2439 FEAT\_nTLBPA

Armv8.4-A introduced a new optional Arm feature, FEAT\_nTLBPA [1, A2.2.1 (p79)].

This feature adds a field to the memory model feature register (AA64MMFR1\_EL1) which can identify whether the current processor's TLB (and related microarchitectural caching structures) may contain non-coherent copies of stage 1 entries indexed by those entries intermediate physical address. Microarchitecturally, this corresponds to there being non-coherent caches associated with the TLB, which must be flushed on a TLBI.

These caches would allow TLB misses to read from a non-coherent cache, thus not seeing the most up-to-date value from the coherent storage subsystem like described in §8.4.

Note that the text in the reference manual is a little ambiguous, the entry in A2.2.1 describes it as a "mechanism to identify if [TLB caching] does not include non-coherent caches [of old translation entries] since the last completed TLBI". This change adds a field to the register, whose reserved value in Armv8.0 corresponds to the non-coherent caches existing. This implies that implementation of the feature is not only the existence of the runtime identification register's field, but additionally that its value is 0b0001 (that is, that non-coherent caches do not exist). This further implies that in processors without FEAT\_nTLBPA one should assume that TLBs may contain non-coherent caching structures.

## 2454 8.6.2 Thread-local ordering and TLBI

TLB maintenance instructions are not naturally locally ordered with respect to other instructions in the instruction stream, this means that they can be re-ordered with other instructions. To ensure they are synchronized with other instructions, the programmer can use the DSB barrier instruction to order instructions before and after it.

Leaving out one, or both, of the DSBs around the TLBI leads to insufficient ordering around the TLBI and allows the invalidation to occur at the wrong time. For example, the CoWinvT.EL1+tlbi-dsb-isb test [Figure 8.38] is allowed as the initial write and TLBI may be re-ordered, negating the architectural effect of the TLBI.

#### 2461 TODO: talk about FEAT\_ETS

# 2462 8.6.3 Broadcast

<sup>2463</sup> Arm provide broadcast variants of the TLBI instructions. These are generally suffixed with the letters IS (for <sup>2464</sup> "Inner-shareable").

AArch64	MP.RT.EL1+dsb-tlbiis-	
	dsb+dsb-isb	
Initial State		
0:R0=0b0	1:R1=y	
0:R1= <b>pte3</b> (x)	1:R3=x	
0:R2=0b1	1:VBAR_EL1=0x1000	
0:R3=y	1:PSTATE.SP=0b0	
0:R4= <b>page</b> (x)	1:PSTATE.EL=0b00	
0:PSTATE.EL=0b01		
<pre>physical pa1 pa2;</pre>		
x  -> pa1;		
$x \mapsto invalid;$		
y  -> pa2;		
<pre>identity 0x1000 with code;</pre>		
Thread 0	Thread 1	
STR X0, [X1]	TOR X0 [X1]	
DSB SY	DSB SV	
TLBI VAE1IS,X4	TSB	
DSB SY	LDB X2. [X3]	
<b>STR</b> X2, [X3]	inc, [inc]	
	Thread 1 EL1 Handler	
	0x1400:	
	MOV X2,#1	
	MRS X13,ELR_EL1	
	ADD X13,X13,#4	
	MSR ELR_EL1,X13	
	ERET	
Final State		
1:X0=1 & 1:X2=0		
Forbid		



The broadcast TLBI on Thread 0 (c) ensures that the earlier unmapping (a) is seen by the ordered later translation read on Thread 1 (i1), by ensuring Thread 1's local TLB is cleaned of any stale entries for x.

Figure 8.39: Test MP.RT.EL1+dsb-tlbiis-dsb+dsb-isb

<sup>2465</sup> Broadcast TLBIs, sometimes referred to as TLB *shootdowns*, allow one processor to perform maintenance on <sup>2466</sup> another core's TLB.

This is in contrast to other systems, such as for IBM's Power architecture, where maintenance of other cores must be achieved in software through the use of only thread-local invalidation instructions.

**TLB invalidation on another core** One of the simplest examples is a message passing invalidation pattern, where the old entry is removed and a message is sent to another core. This can be seen in the MP.RT.EL1+dsbtlbiis-dsb+dsb-isb test [Figure 8.39].

Instruction restarts Broadcast TLBIs must do more than touch the other thread's TLB. If the other processor
 had already performed translation, using the old stale value, but has not yet finished execution, then that instruction
 must be restarted.

This ensures that Arm broadcast TLBIs have the same behaviour as the traditional software IPI-based shootdown
 (With context synchronization); but also provides a needed security guarantee.

If a mapping is taken away from a process, then future writes to the physical location it used to map to, should
 not be visible to that process anymore.

This guarantee is captured in the RBS+dsb-tlbiis-dsb [Figure 8.40] (Read-broken-secret) test. Once a mapping has been *broken*, and sufficient TLB maintenance performed, any future reads or writes to the original physical location will not be visible through that mapping anymore. Note, however, that this does not mean that instructions which have already completed their execution will be restarted, even if they occur after an earlier restarted instruction. This can be seen in the RBS+dsb-tlbiis-dsb+poloc test [Figure 8.41], where the program-order later load can see the old value, even after the first faults.

While here I describe things in terms of instruction restarting, these behaviours can be (and presumably are)
implemented in terms of waiting: instead of the TLBI forcibly restarting instructions that already started but
haven't finished, the TLBI can simply wait for them to complete. This phrasing of waiting for completion is how
this process is described in the Arm ARM [1, D5.10.2 (p4928)].

Atomic TLBIs In the previous RBS-shaped tests, I describe the behaviour in terms of writes that occur 'before' the TLBI.
AArch64	RBS+dsb-tlbiis-dsb
Init	tial State
0:R0=0b0	1:R1=x
0:R1= <b>pte3</b> (x)	1:VBAR_EL1=0x1000
0:R5= <b>page</b> (x)	
0:R2=0x2	
0:R3=y	
0:PSTATE.EL=0b01	
<pre>physical pal;</pre>	
x  -> pa1;	
$x \mapsto invalid;$	
y  -> pa1;	
*pa1 = 0;	
identity 0x1000 with code;	
Thread 0	Thread 1
STR X0, [X1]	
DSB SY	
TLBI VAE1IS,X5	LDR X0, [X1]
DSB SY	
<b>STR</b> X2, [X3]	
	Thread 1 EL1 Handler
	0x1400:
	MOV X0,#1
	MRS X13,ELR_EL1
	ADD X13,X13,#4
	MSR ELR_EL1,X13
	ERET
Final State	
1:X0=2	
Forbid	



The broadcast TLB1 of x (c) ensures that the execution of the load of x in Thread 1 either entirely executes using the old translation and finishes before the TLB1 does, or begins execution after the TLB1 finishes.



Arch64	RBS+dsb-tlbiis-	
	dsb+poloo	
Initial State		
0:R0=0b0	1:R1=x	
0:R1=pte3(x)	1:R3=x	
0:R5= <b>page</b> (x)	1:VBAR_EL1=0x1000	
0:R2=0x2		
0:R3=y		
0:PSTATE.EL=0b01		
physical pal;		
x  -> pa1;		
$x \mapsto invalid;$		
y  -> pa1;		
*pa1 = 0;		
identity 0x1000 with	a code;	
Thread 0	Thread 1	
<b>STR</b> X0, [X1]	MOV X0 #1	
DSB SY		
TLBI VAE1IS,X5	MOV X2, #1	
DSB SY	LDR X2. [X3]	
<b>STR</b> X2, [X3]		
	Thread 1 EL1 Handler	
	0x1400:	
	MRS X13,ELR_EL1	
	<b>ADD</b> X13,X13,#4	
	MSR ELR_EL1, X13	
	ERET	
Final State		
1:X0=1 & 1:X2=0		



Even though the broadcast TLBI on Thread 0 (c) ensures that not-yet-completed instructions using the old mapping are restarted, it does not require that the second load of x in Thread 1 (h) be restarted if it has already satisfied its value, as that value must have come from a write before the TLBI.

Figure 8.41: Test RBS+dsb-tlbiis-dsb+poloc

Microarchitecturally a TLBI instruction is very non-atomic: it sends messages to all other cores, performs some action, and sends messages back to the originating core. The program-order earlier DSB ensures that program-order earlier instructions are complete before sending the messages. The program-order later DSB ensures that all program-order later instructions wait for those messages to return.

The presence of these DSBs ensure that the TLBI's effect happens entirely at that point in the instruction stream, and cannot be broken up and re-ordered amongst the other instructions in the stream. This, coupled with the fact that these messages *strengthen* and never weaken the behaviour of other cores, means that you cannot observe a partial TLBI effect. So long as the programmer takes care to maintain the required thread-local ordering.

Because of this, we can think of the TLBI as executing either before an instruction or after an instruction, but
do not need to consider a TLBI executing in the middle of another instruction. This allows us to simplify things,
fitting TLBIs into a (generalised) coherence order, with other writes occurring either before or after.

### 2502 8.6.4 Virtualization

Throughout this section we have considered tests for stage 1 translation with virtual mappings. But many of these questions and behaviours also apply to the stage 2 intermediate physical mappings, with some key differences.

Virtual to physical and IPA caches The existence of TLBs that cache virtual to physical mappings (§8.5.4)
 complicates the TLB maintenance sequence required for changes to the intermediate physical mappings.

When invalidating stale second stage entries from the TLB, it is required for the programmer to do *two* sets of invalidations: first one TLB invalidation to remove any of the old entries for the old IPA to PA, then, perhaps surprisingly, a second TLB invalidation is needed to remove any stale whole translation, VA to PA mappings or any combination thereof, as these could have indirectly cached the result of a second stage translation without

<sup>2511</sup> remembering the IPA.

This can be seen in MP.RT.EL2+dsb-tlbiipais-dsb+dsb-isb [Figure 8.42], where invalidation of *just* the IPA is not enough. Adding an invalidation of the VA (or all VAs), like in MP.RT.EL2+dsb-tlbiipais-dsb-tlbiis-dsb+dsb-isb

<sup>2514</sup> [Figure 8.43], ensures that later translations cannot see the stale value anymore.

AArch64	MP.RT.EL2+dsb-
	tlbiipais-dsb+dsb-isb
Initia	al State
0:R0=0b0	1:R1=y
0:R1=pte3(ipa1, s2_table):R3=x	
0:R2=0b1	1:VBAR_EL2=0x2000
0:R3=z	1:PSTATE.EL=0b00
0:R4= <b>page</b> (ipal)	
0:PSTATE.EL=0b10	
<pre>physical pa1 pa2;</pre>	
intermediate ipal ipal	2;
x  -> ipa1;	
ipal  -> pal;	
ipal → invalid;	
y  -> ipa2;	
ipa2  -> pa2;	
z  -> pa2;	
identity 0x2000 with	code;
*pa1 = 0;	
*pa2 = 0;	
Thread 0	Thread 1
<b>STR</b> X0, [X1]	LDR X0, [X1]
DSB SY	DSB SY
TLBI IPAS2E1IS,X4	ISB
DSB SY	MOV X2,#1
<b>STR</b> X2, [X3]	LDR X2, [X3]
	Thread 1 EL2 Handler
	0x2400:
	MRS X13, ELR_EL2
	<b>ADD</b> X13, X13, #4
	MSR ELR_EL2,X13
	ERET
Final State	
1:X0=1 & 1:X2=0	
Allow (if not ETS)	

Despite the TLB invalidation of the stale IPA (c), a later stage 2 translation-read of that IPA (i1) can still see the old stale value.





AArch64	MP.RT.EL2+dsb	
	tibiipais-dsb-tibiis-	
Initial State		
0:R0=0b0	1:R1=v	
0.R1= <b>nte3</b> (inal. s2 table).R3=x		
0:R2=0b1	1:PSTATE.EL=0b00	
0:R3=z	1:PSTATE.SP=0b0	
0:R4=page(ipal)	1:VBAR EL2=0x2000	
0:PSTATE.EL=0b10		
physical pal pa2;		
intermediate ipal ipal	2;	
x  -> ipa1;		
ipal  -> pal;		
$ipa1 \mapsto invalid;$		
y  -> ipa2;		
ipa2  -> pa2;		
z  -> pa2;		
identity 0x2000 with code;		
*pa1 = 0;		
*pa2 = 0;		
Thread 0	Thread 1	
STR X0, [X1]		
DSB SY	TDB VO (V1)	
TLBI IPAS2E1IS,X4	LDR X0, [X1]	
DSB SY	DSB 51	
TLBI VMALLE1IS	15D	
DSB SY	LDR X2, [X3]	
<b>STR</b> X2, [X3]		
	Thread 1 EL2 Handler	
	0x2400:	
	MOV X2,#1	
	MRS X13, ELR_EL2	
	ADD X13, X13, #4	
	MSR ELR_EL2,X13	
	ERET	
Final State		
1:X0=1 & 1:X2=0		
Forbid		

By performing TLB invalidation of the stage 1 entries (e) after invalidating the stage 2 ones (c1), it is guaranteed that the later translation-read (k1) cannot see the old stale value anymore.



Figure 8.43: Test MP.RT.EL2+dsb-tlbiipais-dsb-tlbiis-dsb+dsb-isb

### 2515 8.6.5 Break-before-make

TLBs are not required to store only a single cached translation for a given address. There may, in general, be multiple valid translations cached in the TLB.

To avoid this possibility, the architecture provides a *break-before-make* sequence, which will ensure that there cannot be two cached translations existing in the TLB at the same time.

The architecture requires break-before-make when writing to the translation tables to update an already valid entry with a new valid entry, and the change involves any of the following<sup>1</sup>:

- <sup>2522</sup> > A change in output address, if the new or old entry is writeable.
- <sup>2523</sup> ▷ A change in output address, if the new and old locations have different contents.
- <sup>2524</sup> ▷ A change in memory type.
- <sup>2525</sup> ▷ A change in block size (e.g. replacing a page of 4KiB leaf with a 2MiB block mapping).
- <sup>2526</sup> For those cases where break-before-make is required, the programmer must:
- <sup>2527</sup> (1) write an invalid entry to overwrite the currently valid translation table entry in memory;
- 2528 (2) perform a dsb sy (or equivalent);
- (3) perform any TLB maintenance required to sufficiently invalidate the old entry from any TLB(s) required;
- <sup>2530</sup> (4) perform a dsb sy (or equivalent);
- (5) write the new valid translation table entry, overwriting the old invalid entry.

Litmus test For completeness, the BBM+dsb-tlbiis-dsb [Figure 8.44] gives possibly the simplest valid to valid concurrent update test,

#### 2534 Break-before-make violations

Architecturally, there is no hard requirement to perform break-before-make. Failure to do so simply leads to a
 degraded state, defined by CONSTRAINEDUNPREDICTABLE behaviour.

The Arm reference manual does make it clear that failure to perform break-before-make when required can lead to failure of single-copy atomicity, coherence or even the full breakdown of uniprocessor semantics. While the reference manual does not give motivation for this, we can speculate that this is to allow hardware to perform multiple translations during execution of the instruction, for example, during hazard checking. As such, we do not try to give a full picture of CONSTRAINEDUNPREDICTABLE behaviour arising from break-before-make not being followed.

<sup>2543</sup> Understanding CONSTRAINEDUNPREDICTABLE in full is future work, but a quick summary might be 'any behaviour that this program could have performed if it wanted to'. That is, an instantenous change in the state to a random new state that would have been reachable by executing arbitrary code at that same exception level, security state and translation regime.

### 2547 8.6.6 ASIDs and VMIDs

In an effort to reduce the expense of TLB maintenance the architecture provides a mechanism to separate out the address spaces by tagging translations with *address space identifiers* (or ASIDs). These ASIDs allow TLB entries to be tagged with only the address space they are used with, and allow TLB maintenance operations to selectively target only the address space being updated.

<sup>2552</sup> Crucially, this allows software to switch between address spaces without having to invalidate the TLB.

<sup>2553</sup> This idea is extended not just to address spaces at EL1 (used primary for the operating system and its processes),

<sup>2554</sup> but to EL2 with *virtual machine identifiers* (or VMIDs). These VMIDs serve the same function as ASIDs, giving IDs

to address spaces, except in this case IDs to second-stage IPA to PA address spaces.

<sup>&</sup>lt;sup>1</sup>See the Arm ARM "TLB maintenance requirements and the TLB maintenance instructions" [1, D5.10.1 (p4913)] . for the full list of conditions.

AArch64	BBM+dsb-tlbiis-dsb
Initia	al State
0:R0=0b0	1:R1=x
0:R1=pte3(x)	1:VBAR_EL1=0x1000
0:R2=mkdesc3(oa=pa2)	1:PSTATE.SP=0b0
0:R4=0b1	1:PSTATE.EL=0b00
0:R6= <b>page</b> (x)	
0:PSTATE.EL=0b01	
<pre>physical pa1 pa2;</pre>	
x  -> pal;	
$x \mapsto invalid;$	
$x \mapsto pa2;$	
identity 0x1000 with code;	
*pa2 = 2;	
Thread 0	Thread 1
STR X0, [X1]	
DSB SY	
TLBI VAE1IS,X6	LDR X0, [X1]
DSB SY	
<b>STR</b> X2, [X1]	
	Thread 1 EL1 Handler
	0x1400:
	MOV X0,#1
	MRS X13,ELR_EL1
	ADD X13, X13, #4
	MSR ELR_EL1, X13
	ERET
Final State	
1:X0=0	
Allow	

The update of the translation table entry for x in Thread 0 follows the break-before-make sequence, first *breaking* x (a), then performing the necessary TLBI sequence (b-c-d), before *making* x be the new mapping (e). This ensures the concurrent access in Thread 1 is guaranteed to see either the old value, the intermediate broken page (and so a page fault), or the new value. This test is the variant whose final state asserts that the old value was read.



Figure 8.44: Test BBM+dsb-tlbiis-dsb

### 2556 8.6.7 Access permissions

Accesses which result in permission faults can have been satisfied from the TLB, and writes which update translation table entries AP field can be cached in the TLB.

Translations can give rise to permission faults. These are unlike translation faults, in that, they are based not just upon the descriptor read, but also on the *kind* of access requested: whether a read, or a write.

Accesses which result in permission faults result in exceptions, much like translation faults do, but may have been read from the TLB. This can clearly be seen in the CoWinvTp.ro+dsb-isb test [Figure 8.45], where ordered after a write to the translation tables a permission failure is experienced, whose descriptor must have come from the TLB.

Multiple cached entries The changing of access permissions not necessarily being break-before-make vi olations allows us to observe multiple cached entries within the TLB. It is permitted for these entries to exist
 simultaneously.

<sup>2567</sup> When reading from the TLB, and there existing multiple entries for the same input address, it is allowed for the <sup>2568</sup> hardware to generate a *TLB conflict abort*. These aborts are reported as data aborts.

<sup>2569</sup> If the hardware does not generate a conflict abort, then translation reads of that address are ConstrainedUN-<sup>2570</sup> PREDICTABLE, nondeterministically able to read one or the other or an "amalgamation" of the values [1, K1.2.3 <sup>2571</sup> (p11243)].

Here there seems a contradiction: it is not required to perform break-before-make, but there is no requirement that only one entry be cached in the TLB. We can side step this issue by constructing a test that only changes a single bit of the descriptor, in a way that is not a break-before-make violation, and therefore avoiding any questions about how 'amalgamation' of entries happens. This can be seen with the MP.RTpT.ro+dmb-dmb+dsb-isb-dsb-isb test [Figure 8.46], where the existence of multiple cached entries in the TLB allows multiple translation-reads to read from different stale writes.

**Atomic TLB reads** Existence of multiple cached translation table entries in the TLB, without break-beforemake violations, introduces the question of whether those TLB fills and subsequence TLB reads must read from entire single-copy atomic writes of the original translation table entries (much like a read of memory would) or whether a translation read can read from a mix of different writes. RMD+dmb [Figure 8.47] ("Read-mixeddescriptor") shows that translation reads cannot read partially read from a write, it must read from the entire write or none of it.

AArch64 CoWin	vTp.ro+dsb-isb	
Initial State		
0:R0=0x0		
0:R1=pte3(x)		
0:R2=0x1		
0:R3=x		
0:VBAR_EL1=0x1000		
0:PSTATE.SP=0b0		
physical pal;		
x  -> pa1 with [AP = 0b11]	and default;	
$x \mapsto invalid;$		
*pa1 = 0;		
identity 0x1000 with code;		
Thread 0		
STR X0, [X1]		
DSB SY		
ISB		
<b>MOV</b> X13,#0		
STR X2, [X3]	11	
Thread 0 EL1 Han	dler	
0x1400:		
// read ESR_ELI.ISS to see if F	Permission or Tra	
MRS X14, ESK_ELI		
AND X14, X14, #0b1111		
CMP X14, #Ubilii		
MOV A15, $\#1$ // Permission MOV X16 $\#2$ // Translation		
MOV A10, #2 // Iranslation		
MBS X20 FLR FL1		
ADD X20, X20, #4		
MSR ELR EL1, X20		
ERET		
Final State		
0:X13=1		
Allow		



The translation-read (d1) of x, which happens after the program-order earlier store to the translation tables (a) because of the intervening dsb; isb sequence (b-c), can read from a stale value and result in a permission fault, as the read-only entry from the initial state may be cached in the TLB.

Figure 8.45: Test CoWinvTp.ro+dsb-isb

AArch64	MP.RTpT.ro+dmb	
Init	tial State	
0.B0=mkdesc3(oa=pal	0.P0=mkdogg2(op=pp10P=0kHDby	
0:R1=pte3(x)	1:R4=x	
0:B2=0b0	1:VBAR_EL1=0x1000	
0:R3=pte3(x)	1:PSTATE.SP=0b0	
0:R4=0b1		
0:R5=y		
physical pal pa2;		
x  -> pal with [AP =	0b11] and default;	
$x \mapsto pa1$ with [AP =	0b10] and default;	
$x \mapsto invalid;$		
y  -> pa2;		
*pa1 = 0;		
identity 0x1000 with code;		
Thread 0	Thread 1	
	LDR X0, [X1]	
	DSB SY	
STR X0, [X1]	ISB	
DMB SY	LDR X13, [X4]	
<b>STR</b> X2, [X3]	<b>MOV</b> X2, X13	
DMB SY	DSB SY	
<b>STR</b> X4, [X5]	ISB	
	LDR X13, [X4]	
	MOV X3, X13	
	Inread I ELI Handler	
	UX1400:	
	MRS X14 ESR FL1	
	AND X14, X14, #0b	
	<b>CMP</b> X14, #0b1111	
	MOV X15, #1 // Per	
	MOV X16, #2 // Tra	
	CSEL X13, X15, X1	
	MRS X20,ELR_EL1	
	<b>ADD</b> X20, X20, #4	
	MSR ELR_EL1,X20	
	ERET	
Fir	nal State	
1:X0=1 & 1:X2=1 & 1:X3=0		
Allow		

The first translation-read of x (i1) reads from the write that removes read permissions (a) and this write must have come from the TLB because of the intervening invalidation (c), message pass (e-f), and dsb; isb sequence (g-h). The later translation-read of x (m1) can still see an even older value with read permissions, from the initial state, as it may *also* have been cached in the TLB.



Figure 8.46: Test MP.RTpT.ro+dmb-dmb+dsb-isb-dsb-isb

AArch64	RMD+dmb
Initia	l State
0:R0=mkdesc3(oa=pa2, AP=DbRD)x	
0:R1= <b>pte3</b> (x)	1:VBAR_EL1=0x1000
0:R2=0x1	1:PSTATE.SP=0b0
0:R3=y	
<pre>physical pa1 pa2;</pre>	
x $  \rightarrow$ pal with [AP = 0	b11] and default;
$x \mapsto pa2$ with [AP = 0b10] and default;	
y  -> pa2;	
*pa1 = 0;	
*pa2 = 1;	
identity 0x1000 with c	:ode;
Thread 0	Thread 1
<b>STR</b> X0, [X1]	MOV X0 #0
DMB SY	
<b>STR</b> X2, [X3]	HDR X0, [X1]
	Thread 1 EL1 Handler
	0x1400:
	MRS X20,ELR_EL1
	<b>ADD</b> X20, X20, #4
	MSR ELR_EL1,X20
	ERET
Final State	
1:X0=1	
Forbid	

The translation-read of x (d1) cannot read from both the 64-bit single-copy atomic write a as well as from the initial state. Note that this test does not, as far as we can see, violate the break-before-make requirements, as currently prescribed by the Arm manual, as the contents in memory of both locations pa1 and pa2 are the same at the time of the write to the translation tables.

This diagram was generated by hand, as isla does not generate a candidate execution of this shape.



Figure 8.47: Test RMD+dmb

## 2584 8.7 Context synchronisation

There are many operations which change the current context the system is in. We will focus in on two of these: taking and returning from exceptions, and writing to system registers.

<sup>2587</sup> These actions can change the context that the system is executing in: the current exception level, the translation <sup>2588</sup> regime, the translation table base, the ASID or VMID, and a variety of other system configuration state.

### 2589 8.7.1 Relaxed system registers

So far, in this and previous work, register reads and writes have been completely coherent: instructions programorder after a write to a register will always read from that write (or an intervening write) when it reads that register. System registers break this guarantee.

Arm System registers may require the programmer to insert explicit synchronization, as stated clearly in the Arm reference manual [1, D13.1.2 (p5235)] :

Reads of the System registers can occur out of order with respect to earlier instructions executed on the same PE, provided that both:

- Any data dependencies between the instructions, including read-after-read dependencies, are respected.
- ▷ The reads to the register do not occur earlier than the most recent Context synchronization event to its architectural position in the instruction stream.

2595

<sup>2596</sup> This means a read of a system register might not read from the most recent write to that system register.

<sup>2597</sup> To ensure that writes to system registers are seen by program-order later reads, the programmer can ensure

that a *Context synchronization* event occurs. These are typically things which flush the pipeline causing future

<sup>2599</sup> instructions to restart: The ISB instruction and taking and returning from exceptions.

There are two important caveats: (1) this does not apply to non-System registers, such as special-purpose or general-purpose registers, and they never require synchronization; and (2), the synchronization required for System registers depends on the *kind* of accesses.

There are typically two kinds of accesses to System registers: direct, and indirect. Direct accesses are the way we think of registers: instructions which specifically read or write to those registers. Indirect accesses happen when an instruction which does not explicitly mention the register by name performs an access, a read or a write, to that register, during the execution of its behaviour.

Because of the out-of-order nature of the pipeline, these indirect register reads and writes may occur out-of-order with respect to any program-order earlier direct reads or writes of that register. This means that before any direct read, and after any direct write, the programmer must perform a context-synchronizing event to ensure that these direct accesses occur in-order with respect to other indirect accesses. The programmer does not have to insert context-synchronization *after* any direct read, as it is guaranteed that register reads or writes cannot be affected by program-order later accesses.

**System register ASL** In the previous chapter we explored the Arm ASL code for the translation table walk and for one of the store instructions. We saw that this ASL code reads from system registers (as indirect reads).

A naive attempt at a first interpretation of the relaxed semantics is to allow these reads to read-from the most recent indirect write and any program-order later direct writes since the last context synchronization event. However, this would not give the correct behaviour.

The Arm ASL is not written to accommodate relaxed system register behaviours. It leaves questions open about whether these registers can be redundant re-read during execution, whether the instruction reads the entire register at once or piecemeal over the course of execution, and whether repeated accesses to the same register within an instruction are able to read-from different writes. These questions, and others, are still under discussion with Arm.

We will see later in **\$TODO: ?REF?** that we give a simple incomplete (and possibly unsound) interpretation in our model in the *pointed set* semantics of system registers, which allows the model to observe some of the known

<sup>2625</sup> behaviours in this area, without yet fully exploring the architecture.

Caching of system registers in TLBs In addition to being out-of-order due to pipeline effects, some system
 registers may be indirectly cached within the TLB.

We have already seen one of these: the MAIR register. Direct writes to the MAIR may not be seen by program-order later translations, even after context-synchronization, as the translations may get their value from the TLB and the TLB may have stored a result which depended on the previous value of the MAIR, effectively causing a stale read of it at that point in the instruction stream.

To ensure that an update to the MAIR is observed by later translations therefore requires both TLB maintenance and context synchronization, in that order.

The registers which can be cached in this way, and the behaviours that arise from this caching, are still under current investigation with Arm.

## 2636 8.8 Details likely to change

There have been a few places so far I've added words to the effect of 'this is currently under discussion with Arm'. In this section I will summarise those things which we know some things, but also know that it is is likely to change and the ways in which it will.

Caching of entries in the TLB The biggest change we are aware of to the model is a strengthening in caching
of entries in the TLB. We have assumed that the TLB can cache any combination of translation table entries for a
walk, and recall any cached combination as well. We are aware that Arm wish to strengthen this, to make the
model TLB more in-line with the hardware TLB: essentially requiring the model TLB to behave as a walk cache,
caching whole walks (or prefixes of walks) rather than individual entries.

**TLB Invalidation** We have primarily considered TLB invalidations of cached level 3 (that is, last level) entries. When invalidating entries higher in the table, they affect more of the address space (as described in §7.3.1), and so the TLB invalidation must affect addresses outside of just the page referenced. The model currently does not support this, this is a simple oversight and we believe not hard to update the model to handle this case.

More complex invalidation patterns, for example, zeroing and invalidating table entries, is still under discussion
 with Arm.

**ETS** We are aware of changes to the architecture regarding ETS. Every attempt has been made to try incorporate those changes into the model as we become aware, however, often they are changes driven by others and we only become aware as they are publically released.

It is very likely the parts of the model dedicated to ETS will gain new strength over the coming weeks and months, it is unfortunately not possible at this time to give a detailed description of what the final state of ETS will look like, partially for confidentiality reasons and partially because Arm have not yet decided.

**System registers** As previously described, the current state of relaxed system register reads and writes is unclear. We are in discussion with Arm on this aspect. It is not possible at this time to describe what the final model will look like, or what changes will need to be done to the model presented here.

**Exceptions and context-synchronisation** We are in discussion with Arm about the nature of exceptions and their context-synchronising nature and how this interacts with the memory model.

We believe the changes required to the model presented here will be minor, although they will probably be neither a relaxation nor a strengthening of the current model, but rather an incomparable change.

## 2664 8.9 Contributions

We have now covered all the relaxed memory behaviours, and will, in the next chapter, move on to discuss the extant models created to capture those behaviours. But before that, it may at this point be unclear what the *contribution* of this chapter is. They come in three forms: (1) the attempt at some systematic coverage of the kinds of behaviours which systems software must account for; (2) the precise, formal description (in prose, and as litmus

tests) of those behaviours; and, (3) the clarification of the architecture where such behaviours were otherwise 2669 unclear. 2670

**Coverage of behaviours** While this chapter attempts to systematically cover the behaviours we imagine 2671 software may try to rely on, starting from the basics of translation table walks and exploring the effects of 2672 out-of-order pipelines, caching, and barriers, we cannot claim it is exhaustive. As this is a manually compiled and 2673 curated list of behaviours, from reading the text and talking with architects, there are surely corner cases missed 2674 and software patterns overlooked. However, we believe we have covered those patterns known and required for 2675 the features we cover enough for software verification efforts of microkernels and hypervisors. 2676

**Clarification of architecture** Attempts to clarify the architecture come primarily from the confidential discus-2677 sions with architects. The behaviours discussed usually fell into one of three categories, whether they were clear 2678 already, needed further exploration or are, still, under invesitgation by Arm. 2679

The first major category are those behaviours which were already clear and potentially covered in the architecture 2680 text. As alluded to right at the start of this chapter, these are not whole sections or sub-sections or even necessarily 2681 whole tests. The most obvious cases are §8.3.3 ('Invalid entries'), §8.2.1 ('Virtual coherence'), and §8.6.5 ('Break-2682 before-make'). These are fundamental behaviours to the correctness of all modern systems software, and for 2683 which the architecture reference manual has clear words (at least, enough to cover the basic sequences software 2684 rely upon). 2685

Most of the subsections fall into a more general category, of things that either had some associated reference 2686 materials, or was otherwise clear from discussion with architects, but for which further invesitgation was needed. 2687 This includes: forwarding (\$8.4.4) and speculation (\$8.4.5) for translation table walks; multi-copy atomic translation 2688 table walks (§8.4.7); intra-instruction ordering (§8.4.8,§8.4.9); micro-TLBs (§8.5.3) and partial walk caching (§8.5.4); 2689 a variety of TLBI questions (§8.6); and, system register accesses (§8.7.1).

2690

Despite the work conducted here, from reading the architecture reference text, discussions with architects, and 2691 the testing of existing hardware, there are still many questions which are under current investigation with Arm. 2692

These include further questions about the scope of TLBIs, interaction with exceptions and interrupts, changes in 2693

cacheability, translations for instruction fetching, and relaxed system register accesses. Those areas will require 2694

more work before giving a concrete semantics. 2695

## **Chapter 9**

# An axiomatic VMSA model

This chapter is based, in part, on: Relaxed virtual memory in Armv8-A [54] by Ben Simner, Alasdair Armstrong, Jean
 Pichon-Pharabod, Christopher Pulte, Richard Grisenthwaite, and Peter Sewell. Published in the proceedings of the 31st

We now define a semantic model for Armv8-A relaxed virtual memory that, to the best of our knowledge, captures the Arm architectural intent for the questions discussed in Chapter 8, including Stage 1 and Stage 2 translation-table walks and the required TLB maintenance.

In §8 we described the design issues in microarchitectural terms, discussing the behaviour of translation table walks and TLB caching, along with the needs of system software. We now abstract from microarchitecture: constructing a model based on ordering between translation-read events and others, avoiding modelling TLBs and out-of-order pipelines directly.

<sup>2708</sup> This chapter will present this model, as an extention to the 'user-mode' Armv8-A axiomatic model presented in <sup>2709</sup> **§TODO: ?REF?**.

### **9.1** Extended candidate executions

The base Armv8-A axiomatic model is defined as a predicate over *candidate executions*, each of which is a graph with various events (reads, writes, barriers) and relations over them, notably the per-thread program order po, the per-location coherence order co, the reads-from relation rf from writes to reads, the addr, data, and ctrl-dependency subsets of po, and others.

<sup>2715</sup> We extend these candidates with both new events and new relations over those events, as well as modifying some <sup>2716</sup> of the original ones.

### 2717 9.1.1 Candidate events

- <sup>2718</sup> In addition to the events of the original model, we add the following events to the candidates:
- <sup>2719</sup> For reads originating from architected translation-table walks.
- <sup>2720</sup> These roughly correspond to the actual satisfaction from memory which with TLBs may happen very early.
- TLBI events for each TLBI instruction, with a single such event per TLBI instruction, corresponding to the
   TLBI being completed on all relevant cores.
- <sup>2723</sup> > TE and ERET events for taking and returning from an exception (these might not correspond to changes in exception level).
- ▷ MSR events for writes to relevant system registers, such as the TTBR.
- <sup>2726</sup> ▷ DSB events for DSB instructions.
- <sup>2727</sup> ▷ Fault events for translation and permission faults.

**Translation-reads** During execution of the ASL TranslateAddress function (§7.6) there will be many reads, which would usually generate R events. When those reads happen during the TranslateAddress call, they instead generate T events. This means that each translation table walk may generate up to 24 T events, before the instruction generates the (explicit) R | W event.

2696

2697

<sup>2700</sup> European Symposium on Programming (ESOP, 2022).

Alternative representations were explored, including leaving them as R events or collecting all reads into a single large translation event. But these options did not give the clarity and fine granularity we desired in the model, and would require more relations and axioms than presented here.

We also choose not to include TLB hits and misses in the model directly, but instead model the TLB as a relaxation of the values the walk can read from, much like normal R data memory read events and modelling load buffering, write gathering and caches.

We add a helper set, T\_f, for translation reads which read-from a write whose value is even. That is, an entry whose invalid bit is 0. If a translation read results in a fault, either because it was an invalid entry and we get a translation fault, or because the access permissions of the resulting translation do not permit the kind of requested access and so result in a permission fault, the candidate will contain a Fault event (partitioned into Fault\_t and Fault\_p for translation and permission faults) in po order where the explicit memory event would have been. See text on obETS for more discussion of these 'ghost' fault events.

We partition the T set into two subsets: Stage1 and Stage2 for translation read events from a stage 1 or stage 2 walk respectively (stage 2 reads during a stage 1 walk are marked as stage 2, not stage 1).

<sup>2746</sup> Finally, we leave the M set unchanged, which contains only the explicit reads and writes performed by instructions.

**TLBIs** As described in §7.7 Arm have a variety of TLBI instructions, with varying arguments. All of these TLBIs generate a single TLBI event.

<sup>2749</sup> To aide in modelling, there are a set of subsets of TLBI for various kinds of TLBI:

- <sup>2750</sup> ▷ TLBI-S1 for invalidations of Stage 1 entries.
- <sup>2751</sup> ▷ TLBI-S2 for invalidations of Stage 2 entries.
- <sup>2752</sup> ▷ TLBI-IPA for invalidations by intermediate physical address.
- <sup>2753</sup> ▷ TLBI-VA for invalidations by virtual address.
- $_{\rm 2754}$   $\,\triangleright\,$  TLBI-ASID for invalidations by ASID.
- 2755 ▷ TLBI-VMID for invalidations by VMID.
- $_{\rm 2756}$   $\,\triangleright\,$  TLBI-ALL for the TLBI ALL instructions.
- $_{\rm 2757}$   $\,\triangleright\,$  TLBI-IS for broadcast TLBIs.
- <sup>2758</sup> ▷ TLBI-EL1 for invalidations of the EL1&0 regime.
- 2759 ▷ TLBI-EL2 for invalidations of the EL2 regime.

<sup>2760</sup> These events do not *cut* the TLBI set into partitions, but rather any TLBI event may belong to multiple. For

- example, a TLBI VAE1IS event would belong to TLBI-VA, TLBI-VMID, TLBI-EL1, and TLBI-IS.
- <sup>2762</sup> We also include all TLBIs in a general C ("Cache maintenance") set.

**Exceptions** Despite not modelling exceptions in general in this work, we do need to include some exception machinery in the model to capture the minimal ordering requirements arising from both their context synchronisation

<sup>2765</sup> effects and also behaviours from crossing exception level boundaries.

<sup>2766</sup> To support this we add two new events to capture taking and returning from exceptions: TE ("Take-exception") <sup>2767</sup> and ERET.

**Barriers** The Arm DSB ("Data synchronization barrier") instruction is required for TLB maintenace, as was seen in the previous chapter. We include DSB events, one for each kind of DSB instruction:

- <sup>2770</sup> ▷ DSBSY and DSBISH (here, equivalent as we do not model shareability domains)
- $_{2771}$   $\triangleright$  DSBNSH, for thread-local effects.
- <sup>2772</sup> ▷ DSBST, DSBLD, for DSBs affecting only loads or stores.
- ▷ DSBISHST, DSBISHLD, and so on, for all combinations of DSB instruction domain and access types.

Arm define a hierarchy of barriers where, for example: DMB.LD < DMB.SY < DSB.SY That is, any ordering imposed by a DMB.LD is also imposed by a DMB.SY, and therefore also a DSB.SY.

<sup>2776</sup> To help capture this, and reduce the explosion in the number of relations in the model later on, we simplify and <sup>2777</sup> update the barrier story in the Arm model and include the helper sets given in Figure 9.1.

```
let dsbsy = DSBISH | DSBSY | DSBNSH
let dsbst =
                              DSBISHST
                                          DSBNSHST
                     DSBST
            dsbsv
                            let dsbld = dsbsy
                                         DSBNSHLD
                     DSBLD
                           1
                              DSBISHLD
let dsbnsh = DSBNSH
let dmbsy = dsbsy
                     DMBSY
          =
                              DMBST
                                      DSBST | DSBISHST | DSBNSHST
let dmbst
            dmbsy
                     dsbst
                                    dmbsy
                                    i
let dmbld
          =
                     dsbld
                              DMBLD
                                      DSBISHLD | DSBNSHLD
let dmb = dmbsy
                   dmbst
                            dmbld
let dsb = dsbsy
                   dsbst
                           dsblo
```

Figure 9.1: Barrier helper sets.

Context changing and synchronisation Finally, we add events for context-changing and context-synchronising
 operations. Context changes involve updates to system registers which change the current translation regime,
 which generate MSR events. We add a general context-synchronisation event set CSE which includes ISB, TE and
 ERET.

Changes to system registers may have relaxed behaviours, as described in §8.7.1, but full relaxation of the system 2782 register reads done by the Arm psueocode is unlikely to be valid, consistent or meaningful. Instead, we introduce 2783 a pointed-set semantics: when generating a candidate, we keep a per-system-register set of writes to that register, 2784 remembering which one is the most recent. On a write to that system register, we add it to the set. On a read of that 2785 system register, we generate one candidate for each value in the set, and then 'lock' the remainder of the execution 2786 of that instruction to that value so repeated reads will see the new value. When a context-synchronization event 2787 is generated (that is, an event that will be in the CSE set) all the sets are reduced to singleton sets containing only 2788 the most recent write. 2789

<sup>2790</sup> This gives us some relaxed behaviours, enough to see relaxed behaviours around changes to the TTBR, but we <sup>2791</sup> note that this is unlikely to be the full story for relaxed system registers.

### 2792 9.1.2 Candidate relations

<sup>2793</sup> In addition to those new events, we introduce new relations over those (and other) events:

- <sup>2794</sup> ▷ trf and tfr as analogues to rf and fr but for translation-read (T) events.
- i i o ("intra-instruction order") which relates events of the same instruction in the order they occur during
   execution of that instruction's intra-instruction semantics as defined by the Arm ASL.
- same-va, same-ipa, same-pa relations which relate events whose virtual, intermediate physical or physical
   address of the associated explicit memory access are the same.
- same-va-page, same-ipa-page, same-pa-page which relate events whose associated explicit memory
   events are in the same page (e.g. 4KiB chunk) of the virtual, intermediate physical or physical address space.
- same-asid, same-vmid relates events for which translations for the associated memory event are using the
   same ASID or VMID.
- > wco, a generalised coherence order which includes TLBIs.

Addresses, ASIDs and VMIDs Each translation table walk will read from registers and system registers and get a value for the (input) address, the current ASID and current VMID. We then relate each T with any other T where the translation associated with it is for the same virutal address (with same-va), the same intermediatephysical address (with same-ipa), or the same resulting physical address (same-pa). This means that all T events within a translation have the same same-\* relations. We also include relations which match translation's virtual, intermediate physical and physical addresses if they are in the same page rather than exactly, with the same rules, but as a same-\*-page relation.

<sup>2811</sup> If two translations are for the same ASID, their translation reads are related by same-asid. If two translations are <sup>2812</sup> for the same VMID, their translation reads are related by same-vmid.

To use these relations we also include TLBI events. A TLBI-X is related to T by same-X if the parameter to the TLBI instruction (the page, vmid, or asid) either passed by register, immediate or through the current context, if the T event's associated translation matches X. For example, a TLBI-IPA event would be same-ipa-page related to a T whose translation was for an intermediate physical address in the page provided as the parameter to the TLBI IPA instruction. **Generalised coherence order** We create an extended coherence order wco, which is the usual co (a per-location total order of writes to that location) as well as their relative ordering to all TLBI events.

One might be concerned at the validity of doing this, on two fronts. Generally, extending coherence to a total order over all locations is sound [9, §10.5 p174], and so there is no issue in doing this. Secondly, for broadcast TLBIs, microarchitecture will implement these with message passing to and from each core separately, and so there is no single moment the TLBI 'happens'. However, as described in §8.6.7 we seem to be able to consider TLBI instructions as executing 'atomically' so long as there are no break-before-make violations. This is a similar justification as to including DC and IC events in a similar generalised coherence order for instruction fetching [58, §5 p29].

**Dependencies** A candidate execution consists not only of events, and reads-from relations but also a set of dependencies: addr, data, ctrl, po and loc. We add iio and tdata to these.

The intra-instruction ordering iio relation relates two events in the same instruction in the order the Arm pseudocode generated the events. This relation therefore captures a total order over all events within an instruction, regardless of the intra-instruction dependencies (control, data) or unordered accesses (for example, for misaligned accesses). We are currently invesitaging a relaxation of this ordering, and associated changes in the underlying Arm pseudocode definitions, to enable a more relaxed definition of the ordering within an instruction to handle these cases.

We make loc relate events with the same physical address (for T events, this is the physical location of the translation table entry).

Program order (po) is restricted to explicit events: R, W, F, C, CSE and MSR. Implicit translation reads (T) and any
 indirect reads or writes of registers are not included in po.

Address dependencies were once fundamental, but now we can define address dependencies in the presence of address translation as dependencies into the translation table walk. To do this, we include a new relation tdata that relates reads with the translation read events of a translation which reads from the register written by that read to compute the address. The traditional addr can be recovered as tdata ; iio\* ; [M].

2843 9.2 Cat model

The base Arm axiomatic model had three axioms: internal, external, and atomic. These were acyclicity and emptyness checks of unions of set of relations: obs, dob, aob and bob. We will slightly modify three of these relations obs, bob and dob, and add 5 new ones (tob, obtlbi, ctxob, obfault, obETS) to handle the ordering between translations and TLBIs, and include them in the external acyclicity check. Then we will introduce one final new axiom translation-internal.

Figure 9.2 contains the axioms and relations for the updated Armv8-A relaxed virtual memory axiomatic model (RVM). Unchanged parts from the original are greyed out. Note that some helper relations are elided here, and will be described in more detail later.

## 2852 **9.3** Axioms

The RVM model axioms are, mostly, a syntactic extension to the original Armv8-A axiomatic model presented in **\$TODO: ref intro.** This is deliberate. Although there may be other, perhaps even nicer or more succinct, ways of phrasing the given model, the variation presented here is designed to be syntactically as close as possible to the original. This helps with readability for those familiar with the original; it allows us to present the differences to the original in an easier form; it makes recovery of the original model easier; and, it makes it easier to prove equivalance of the axiomatic models in the presence of constant address translation, increasing the confidence we have in the model.

The model has 3 kinds of axioms: internal ones for per-location guarantees, an external axiom for the global happens-before ordering, and the atomic axiom for RMWs (untouched in this work).

<sup>2862</sup> Internal axioms The new model has two per-location axioms: internal and translation-internal.

```
FaultFromR)]
| [R|W]; po; [Fault_T &
FaultFromReleaseW]
let speculative =
     ctrl
| addr; po
| [T]; instruction-order
                                                                                                           (* ETS-ordered-before *)
                                                                                                         let obETS =
   (obfault; [Fault_T]); iio^-1; [T_f]
   | ([TLBI]; po; [dsb]; instruction-order
   ; [T]) & tlb-affects
   * translation-ordered-before *)
(* transform
let tob =
    [T_f]; tfre
    [T]; iio; [R|W]; po; [W]
    | speculative; trfi
                                                                                                             * dependency-ordered-before *)
                                                                                                         (* observed by *)
let obs = rfe | fr | wco
| trfe
     ordered-before TLBI and translate *)
(* atomic-ordered-before *)
let aob = rmw
    [range(rmw)]; rfi; [A | Q]
     -S2])

& (same-translation ; [T & Stage1] ;

trf^-1 ; wco^-1)

(([T & Stage2] ; tlb_barriered ; [

TLBI-S2]) ; wco? ; [TLBI-S1])

& (same-translation ; [T & Stage1] ;

movbe TLB cached)
                                                                                                           (* barrier-ordered-before *)
                                                                                                         (* barrier ordered -before :>
let bob = [R]; po; [dmbst]
  [ dmbst]; po; [W]
  [ dmbst]; po; [W]
  [ dmbld]; po; [R|W]
  [ L]; po; [A]
  [ A | Q]; po; [R | W]
  [ F | C]; po; [dsbsy]
  [ Gdb]; po; [dsbsy]
     maybe_TLB_cached)
   * ordered-before TLBI *)
(a ofdefice before tell a)
let obtlbi =
    obtlbi_translate
    [R|W|Fault_T]; iio^-1; (
    obtlbi_translate & ext); [TLBI]
                                                                                                               | [dsb]; po
                                                                                                         (* Ordered-before *)
let ob = (obs | dob | aob | bob
| iio | tob | obtlbi | ctxob | obfault
| obETS)^+
 (* context-change ordered-before *)
(* context-change ordered-berore *)
let ctxob =
    speculative; [MSR]
    [CSE]; instruction-order
    [ContextChange]; po; [CSE]
    speculative; [CSE]
    po; [ERET]; instruction-order; [T]
                                                                                                          (* Internal visibility requirement *)
acyclic po-loc | fr | co | rf as internal
(* External visibility requirement *)
irreflexive ob as external
(* tetrainenent *)
    : ordered-before a translation fault *)
:t obfault =
                                                                                                         (* Atomic requirement *)
empty rmw & (fre; coe) as atomic
(* Writes cannot forward to po-future
translates *)
acyclic (po-pa | trfi) as translation-
internal
(* ordered-before a translation fault *)
let obfault =
    data; [Fault_T & FaultFromW]
    [speculative; [Fault_T & FaultFromW]
    [dmbst]; po; [Fault_T & FaultFromW]
    [dmbld]; po; [Fault_T & (FaultFromW]
    FaultFromR)]
    [A|Q]; po; [Fault_T & (FaultFromW |
```



```
(* Internal visibility requirement *)
acyclic po-loc | fr | co | rf as internal
(* Writes cannot forward to po-future translates *)
acyclic (po-pa | trfi) as translation-internal
```

2863

Unchanged from the original, the internal axiom captures the SC-per-location guarantee briefly discussed in **\$TODO: ?REF?.** Translations, however, do not have the same per-location guarantees. To account for this, we introduce a second axiom, translation-internal, which captures the weaker per-location guarantee for translation table walks. Since translation reads, in the presence of TLB caching and out-of-order pipelines, do not guarantee even coherence, the only behaviour this axiom ends up preventing is translation reads reading from program-order later stores.

External axiom The external axiom asserts acyclicity of the global happens-before ordering for Arm. The happens-before (called ob, 'ordered-before', in Arm) relation is the union of all the ordering relations, given in \$9.4.

```
(* Ordered-before *)
let ob = (obs | dob | aob | bob | iio | tob | obtlbi | ctxob | obfault |
        obETS)^+
(* External visibility requirement *)
irreflexive ob as external
```

2873

We choose to include all the pipeline and TLB effects as ordering requirements, rather than introducing new
ordering axioms just for translation and TLB invalidation. This produces a model that is more consistent with
the previous Arm memory models, and ensures ordering information gained through observing translation table
walks are respected by non-translation-table accesses.

Atomic axiom The atomic axiom remains unchanged. In this work we do not consider the interaction of translation with atomic accesses.

```
(* Atomic requirement *)
empty rmw & (fre; coe) as atomic
```

2880

## 2881 9.4 Relations

The RVM model modifies some of the original, and introduces some new, ordering relations. This section goes
 through each in detail, describing the mechanism and justifying the existence or non-existence of particular
 clauses.

2885 9.4.1 obs

```
(* observed by *)
let obs = rfe | fr | wco | trfe
```

2886

The 'observed-by' relation. It includes the original rf and fr (over physical locations), the *generalised coherence order* wco (§9.1.2), and the trfe (translation-reads-from-external) relation.

**Generalised coherence** Including wco, which is existentially quantified over the candidates, fixes some global order the writes and TLBIs happen in. Consider, informally, some microarchitectural execution. It would propagate writes to the coherent storage subsystem, and would complete TLBI instructions, and these events would be interleaved in some whole-machine trace. The generalised wco relation captures the relative ordering of these events in the axiomatic model, as they would have happened in the traces of machine executions. The model is then quantified over all such orderings, accounting for any interleaving of these events. **External translation reads** Inclusion of trfe enforces that translation-table-walk translation reads, which could not come from forwarding, must have originally come from the coherent storage subsystem and so the write must have been globally propagated before the translation read happened (§8.4.2, §8.4.7).

However, the translation read might have happened much later, either due to extreme out-of-order (§8.4.1) or TLB caching (§8.5.1), and so we do not include tfre (translation-from-reads-external) in ob.

Additionally, writes may be propagated to that thread's translation table walker before they are propagated to the coherent storage subsystem ( $\S$ 8.4.4), in other words, they can be forwarded. Therefore we do not include trfi

<sup>2902</sup> (translation-reads-from-internal) in obs.

2903 9.4.2 dob

```
let dob =
    addr | data
    speculative; [W]
    addr; po; [W]
    (addr | data); rfi
    (addr | data); trfi
```

2904

The dependency-ordered-before relation is mostly unchanged, we add a single (addr | data); trfi clause to the end to forbid thin-air creation of values (§8.4.1, §8.4.2, TODO: need dedicated thin air paragraph/test in prev chapter) similarly to the original model for data memory reads.

### 2908 9.4.3 bob

```
let bob =
    [R]; po; [dmbld]
    [W]; po; [dmbst]
    [dmbst]; po; [W]
    [dmbld]; po; [R|W]
    [L]; po; [A]
    [A | Q]; po; [R | W]
    [R | W]; po; [L]
    [F | C]; po; [dsbsy]
    [dsb]; po
```

2909

We rewrite the original barrier-ordered-before relation to use the barrier helpers defined in Figure 9.1. This does not change the underlying model for DMB instructions, but allows those same clauses to capture the barrier hierarchy imposing the same ordering when using stronger barriers (namely, DSBs).

The Arm DSB instruction has some extra ordering however. Firstly that a DSB SY orders TLBI instructions (§8.6.2) and so we include [F|C]; po; [dsbsy]. Secondly, all program-order later events must wait for an earlier DSB to finish before performing its explicit memory events, so we also include [dsb]; po in ob.

```
<sup>2916</sup> 9.4.4 tob
```

```
let tob =
    [T_f]; tfre
    [T]; iio; [R|W]; po; [W]
    speculative; trfi
```

2917

<sup>2918</sup> Translation table walks themselves impose ordering on the surrounding events.

Invalid writes The first of these is one of the key behaviours described in 8.3.3, that reads of invalid entries must not have come from the TLB. So we add the [T\_f]; tfre edge to capture this, that any translation-reads which read an invalid entry must happen before any writes coherence after the one it read from.

There is a major caveat here: write forwarding to the translation table walker. We cannot simply have  $[T_f]$ ; tfr as a thread-local write may be forwarded to the translation table walker before it's propagated to memory (§8.4.4). However, it should not be the case that the write is forwarded from a write that is too old or behind a DSB if FEAT\_ETS, except it may be the case that there might be other intervening writes in between. For now, we are unable to give a precise bound on the ordering for thread-local  $[T_f]$ ; tfr, and this area is still currently under investigation with Arm.

**Speculation** As we saw earlier, speculation interacts with translation in two ways: first, it is forbidden to read-from a still speculative write (§8.4.5), and, secondly, events program-order after an instruction which does a translation table walk are speculative until the translation table walk completes (§8.4.1).

<sup>2931</sup> To capture these we first define when one event is considered speculative until another event happens, with a <sup>2932</sup> new speculative relation, defined as following:

let speculative = ctrl | addr; po | [T]; instruction-order

2933

This captures all the control-flow dependencies that we model here, the classic ctrl and addr; po, as well as a new general [T]; instruction-order which says that all events ordered (iio|po)+ after a translation read are speculative until the translation read satisfies. We can then include speculative ; trfi to succinctly forbid any forwarding of still-speculative writes to translation table walks.

Finally, we include [T]; iio; [M]; po; [W] which captures that writes cannot propagate until program-order earlier instructions have their physical address (so, do not fault). Although, this edge is subsumed by the speculative; [W] edge in dob, it is kept here for clarity.

```
2941 9.4.5 ctxob
```

```
2942
```

NOTE: The model for exceptions and context-synchronising events is currently under revision, and what is presented here is likely to change.

```
let ctxob =
    speculative; [MSR]
    [CSE]; instruction-order
    [ContextChange]; po; [CSE]
    speculative; [CSE]
    po; [ERET]; instruction-order; [T]
```

2943

The ctxob relation captures the orderings required from context changing and synchronising operations, without trying to capture the full extent of the relaxed behaviours. As such, these orderings are likely to be incomparable to the real semantics.

Speculation The first guarantee we see is that context changes and synchronisation should not happen speculatively. Speculative context changes may end up creating translation table roots and therefore translation table walks using unreachable writes (§8.5.5). To prevent this we ensure that context changing operations only happen once they are non-speculative, by enforcing speculative; [MSR] in ob. Forbidding speculative execution of context synchronisation is done through the inclusion of speculative; [CSE] in ob.

**Context synchronising** A context synchronisation event (such as an ISB or ERET instruction) should ensure that program-order earlier context-changing events are seen by program-order later instructions. Microarchitecturally this is achieved by having context-synchronisation events flushing the pipeline, restarting all program-order later instructions. For now this effect seems fixed in the architecture (§8.7), and so we get [CSE]; instruction-order in ob subsuming the earlier ISB orderings.

To ensure that context changes are seen after the synchronisation we include [ContextChange]; po; [CSE], and the union of these two relations ensures the context change is ordered before any program-order later events.

**Exceptions** Taking and returning from exceptions are context synchronising (§8.7), and so those are captured by the previous clauses. However, translation reads of a lower exception level should not satisfy during execution at a higher exception level. We over approximate this with po; [ERET]; instruction-order; [T] ensuring all translation reads after an ERET wait. 2963

### 9.4.6 obfault and obETS

```
(* ordered-before a translation fault *)
let obfault =
    data ; [Fault_T & FaultFromW]
    speculative ; [Fault_T & FaultFromW]
    [dmbst] ; po ; [Fault_T & FaultFromW]
    [dmbld] ; po ; [Fault_T & (FaultFromW | FaultFromR)]
    [A|Q] ; po ; [Fault_T & (FaultFromW & FaultFromR)]
    [R|W] ; po ; [Fault_T & FaultFromW & FaultFromReleaseW]
(* ETS-ordered-before *)
let obETS =
    (obfault; [Fault_T]); iio^-1; [T_f]
    | ([TLBI]; po; [dsb]; instruction-order; [T]) & tlb-affects
```

2964

To capture the specific guarantees described by FEAT\_ETS ( $\S8.4.3$ ,  $\S8.6.2$ ), we include 'ghost' Fault events in the candidate executions. These events sit in the execution (in po order) where the explicit memory event would have been if there was no fault, and tags the fault with the kind of fault it was (translation or permission).

**Ordering to a fault** To fully capture the strength of FEAT\_ETS we keep track of syntactic dependencies *into* the instruction which faulted, and apply those dependencies to the Fault event itself. obfault then the syntactic subset of bob and dob where the right-hand side of each clause is substituted with a Fault\_T (a translation fault).

Using obfault we can then keep track of the (syntactic) subset of ob that would have ordered the explicit event after, and associate those relations with the Fault\_T event instead. obETS's first clause then adds to ob this ordering, but attached to the translation read of the invalid entry itself, as architected by FEAT\_ETS.

Note that dependencies and orderings *from* a faulting instruction seem not respected, and so we do not induce
 orderings out of a Fault\_T.

FEAT\_ETS and TLBI The second clause of obETS captures the second architected behaviour of FEAT\_ETS (§TODO: TLBI ordering needs ETS explained), that faults after a thread-local TLBIs do not need context synchronisation to be ordered after the TLBI. Note that one still needs a DSB to complete the TLBI in that case.

2979 9.4.7 obtlbi

```
(* ordered-before TLBI *)
let obtlbi =
    obtlbi_translate
    | [R|W|Fault_T]; iio^-1; (obtlbi_translate & ext); [TLBI]
(* translate ordered-before TLBI *)
let obtlbi_translate =
    [T & Stage1]; tlb_barriered; [TLBI-S1]
| (([T & Stage2]; tlb_barriered; [TLBI-S2]); wco?; [TLBI-S1])
& (same-translation; [T & Stage1]; maybe_TLB_cached)
| ([T & Stage2]; tlb_barriered; [TLBI-S2])
    & (same-translation; [T & Stage1]; trf^-1; wco^-1)
```

2980

Finally, there is the obtlbi relation which captures the ordering from translations (and their explicit memory events) and the TLB invalidations which affect them. The relation is split in two: the obtlbi\_translate clause enforces order between stale translations and the TLBIs they are invalidated by, the second clause covers broadcast TLBIs.

**Capturing stale TLB entries** When a translation read happens, it is allowed for it to read from a stale write (§8.5.1). That is, the translation need not be ordered before writes which come after the write it actually reads from. Consequently the tfre relation is not included in ob.

We strengthen this, by including some edges from translations to TLBIs, when there is an interposing newer write. The general shape of this ordering is illustrated in Figure 9.3.

136



Figure 9.3: General obtlbi\_translate shape.

This shape is succinctly captured by the tlb\_barriered auxiliary relation, which relates any translate-read that reads from a write which is wco before another write which is wco before a TLBI which targets the address, ASID or VMID of the translation:

2993

We cannot simply include tlb\_barriered in ob, however. Instead, we must consider the orderings for stage 1 and stage 2 translation reads separately.

**Stale stage 1 reads** For stage 1 translation reads, either in single-stage regimes or as part of a two-stage regime, we can include a variant of tlb\_barried specialised to stage 1 translation-reads and TLBIs which affect stage 1 entries.

**Stale stage 2 reads** Stage 2 walks are more subtle. The requirement to perform stage 1 invalidation (§8.6.4) means that, in those instances, we do not get tlb\_barriered directly.

Instead, we have to case split on the execution: either, (1) the translation table walk does a stage 1 translation read which reads-from an older write, in which case there may have been a whole cached translation that must be invalidated; or, (2) one of the stage 1 translation reads of the translation table walk reads from a write that is newer than the stage 2 TLBI and so there cannot have been any cached whole translation entries in the TLB and so, logically, we only need the stage 2 invalidation. These cases are illustrated in Figure 9.4, and correspond to the two clauses of obtlbi\_translate which match on stage 2 translation reads.



Figure 9.4: obtlbi stage 2 scenarios.

We capture the general shape of (1), where a translation-read may have been cached in the TLB, with the following maybe\_TLB\_cached relation:

let maybe\_TLB\_cached = ([T]; trf^-1; wco; [TLBI]) & tlb-affects^-1

3009

We then use this relation to add ordering from a stage 2 translation-read to the stage 1 TLBI, wco-after a stage 2 TLBI that removed any stale IPA mappings, which would remove any cached whole-translation any stage 1 translation-read might have read from, and after which any fresh translation table walk would be required to not see the stale stage 2 entry the translation-read read from.

**Broadcast TLBIs** Recall that broadcast TLBIs impose restrictions on other threads (§8.6.3). When a broadcast TLBI's invalidation affects a translation on another core, then it must also affect the explicit memory effect associated with it. This shape is illustrated in Figure 9.5, and corresponds to the final clause of obtlbi.



Figure 9.5: obtlbi broadcast TLBI shape.

3017 Connecting TLB invalidations to translation reads The final part of the puzzle is how to relate TLBI events
 3018 with translations which may be affected by the invalidation. Recall that the TLBIs are grouped into subsets of
 3019 TLBI-S1, TLBI-VA, and so on. We define a tlb\_might\_affect that is the cross-product of these with the same-\*
 3020 relations:

```
let tlb_might_affect =
    TLBI-SI & ~TLBI-S2 &
                            TLBI-VA
                                       &
                                          TLBI-ASID & TLBI-VMID] ; (same-va-
  Г
    page & same-asid & same-vmid) ; [T & Stage1]
    TLBI-S1 & ~TLBI-S2 & ~TLBI-VA
                                          TLBI-ASID & TLBI-VMID] ; (same-
ΙΓ
                                       &
    asid & same-vmid) ; [T & Stage1]
TLBI-S1 & ~TLBI-S2 & ~TLBI-VA & ~TLBI-ASID & TLBI-VMID] ; same-vmid
       [T & Stage1]
   ~TLBI-S1 & TLBI-S2 &
                            TLBI-IPA & ~TLBI-ASID & TLBI-VMID] ; (same-ipa
    -page & same-vmid) ; [T & Stage2]
TLBI-S1 & TLBI-S2 & ~TLBI-IPA &
                           ~TLBI-IPA & ~TLBI-ASID & TLBI-VMID] ; same-vmid
       [T & Stage2]
                TLBI-S2 & ~TLBI-IPA & ~TLBI-ASID & TLBI-VMID] ; same-vmid
    TLBI-S1 &
  Г
       ΓT]
    TLBI-S1 &
                            ~TLBI-IPA & ~TLBI-ASID & ~TLBI-VMID) * (T &
    Stage1)
                 TLBI-S2 & ~TLBI-IPA & ~TLBI-ASID & ~TLBI-VMID) * (T &
  (
    Stage2)
```

3021

<sup>3022</sup> Finally, we get tlb-affects by attaching tlb\_might\_affect to events in the same thread, and if a TLBI-IS, to <sup>3023</sup> ones in other threads too:

```
let tlb-affects =
   [TLBI-IS]; tlb_might_affect
| ([~TLBI-IS]; tlb_might_affect) & int
```

3024

3025

## 9.5 Interface

To support the new Armv9 ISA and the new concurrency interface, we produce architecture-specific definitions using the new isla-cat language features in isla-axiomatic.

<sup>3028</sup> Barriers are instances of the sail\_barrier outcome. For Arm we instantiate these with the Arm Barrier union.

Figure 9.6 contains the isla-cat definitions for the Arm barriers. The sail Arm Barrier union is reproduced here for the reader's benefit but is not required (nor present) in the source cat file. Similar unions, structs, enums and

```
(* F for all fences *)
accessor F: bool = is sail_barrier
(* from the Arm ASL1, compiled to sail *)
union Barrier = {
  Barrier_DSB : DxB,
  Barrier_DMB : DxB, // The nXS field is ignored from DMBs
  Barrier_ISB : unit,
Barrier_SSBB : unit,
Barrier_SSBB : unit,
Barrier_SSBB : unit,
}
(* accessors for each relevant constructor *)
accessor z_dmb: bool =
     .match {
         Barrier_DMB => true,
         _ => false
     }
accessor z_dsb: bool =
     .match {
         Barrier_DSB => true,
         _ => false
     }
accessor z_isb: bool =
     .match {
        Barrier_ISB => true,
         _ => false
     }
(* cat event sets for the different barriers *)
define DMB(ev: Event): bool =
     F(ev) & z_dmb(ev)
define DSB(ev: Event): bool =
    F(ev) & z_dsb(ev)
define ISB(ev: Event): bool =
     F(ev) & z_isb(ev)
```

Figure 9.6: isla-cat accessors for Arm barriers.

corresponding accessors and definitions exist for the Arm barrier domains (NSH, ISH, OSH) and access types (ST,

<sup>3032</sup> LD, SY), elided here for brevity.

We make use of **accessors** to access fields of the sail structs and unions, both here for barriers, and also for exceptions (faults), and TLBIs, as well as defining the trf and wco (found in Figure 9.7) relations. The full isla-cat-defined interface can be found in the Appendix**TODO:** make appendix.

```
declare wco(Event, Event): bool
(* wco has domain and range of W,CacheOp *)
assert forall ev1: Event, ev2: Event =>
     wco(ev1, ev2) -->
(W(ev1) | C(ev1) | (ev1 == IW)) & (W(ev2) | C(ev2))
(* wco is transitive *)
assert forall ev1: Event, ev2: Event, ev3: Event =>
    wco(ev1, ev2) & wco(ev2, ev3) --> wco(ev1, ev3)
(* wco is total *)
assert forall ev1: Event, ev2: Event, ev3: Event =>
    wco(ev1, ev3) & wco(ev2, ev3) & ~(ev1 == ev2) -->
    wco(ev1, ev2) | wco(ev2, ev1)
(* wco is irreflexive *)
assert forall ev1: Event, ev2: Event, ev3: Event =>
    wco(ev1, ev2) --> ~(ev1 == ev2)
(* wco is antisymmetric *)
assert forall ev1: Event, ev2: Event =>
    wco(ev1, ev2) --> ~wco(ev2, ev1)
(* all write/cache-op pairs are wco related *)
assert forall ev1: Event, ev2: Event =>
    W(ev1) & C(ev2) -->
     wco(ev1, ev2) | wco(ev2, ev1)
(* wco is consistent with co *)
assert forall ev1: Event, ev2: Event =>
co(ev1, ev2) --> wco(ev1, ev2)
(* all C are wco after IW
 \star n.b. all W are wco after IW, because all W are co after IW and co =>
    WCO
 *)
assert forall ev: Event =>
     C(ev) \longrightarrow wco(IW, ev)
```

Figure 9.7: wco.cat: isla-cat definition of wco.