

Arm<sup>®</sup> Server Base System Architecture 7.1  
**Platform Design Document**  
Non-confidential



## Contents

Release information	4
Arm Non-Confidential Document Licence (“Licence”)	7
<b>About this document</b>	<b>9</b>
Terms and abbreviations	9
References	10
Rules-based writing	11
Identifiers	11
Examples	12
Feedback	12
<b>1 Server Base System Architecture</b>	<b>13</b>
1.1 Arm Base System Architecture	13
1.2 Introduction	13
1.2.1 SBSA levels and minimum component versions	13
1.3 Level 3	14
1.3.1 PE architecture	14
1.3.2 Memory map	14
1.3.3 Interrupt controller	14
1.3.4 PPI assignments	15
1.3.5 System MMU and device assignment	15
1.3.6 Watchdog	15
1.4 Level 4	15
1.4.1 PE architecture	15
1.4.2 System MMU and device assignment	15
1.4.3 Peripheral subsystems	15
1.5 Level 5	16
1.5.1 PE architecture	16
1.5.2 Interrupt controller	16
1.5.3 System MMU and device assignment	17
1.5.4 Clock and Timer subsystem	17
1.5.5 Watchdog	17
1.5.6 PPI assignments	17
1.6 Level 6	18
1.6.1 PE architecture	18
1.6.2 System MMU and device assignment	18
1.6.3 Watchdog	19
1.6.4 Armv8 RAS extension requirements	19
1.7 Level 7	19
1.7.1 PE architecture	19
1.7.2 MPAM	20
1.7.3 Entropy source	21
1.7.4 SMMU and device assignment	21
1.7.5 Performance Monitoring Unit	22
1.7.6 System RAS	22
1.7.7 PCIe	23
1.8 Future requirements	24
1.8.1 PE architecture	24
1.8.2 Self-hosted debug	25
1.8.3 System MMU and device assignment	25
1.8.4 System RAS	25
1.8.5 PCIe	25
1.8.6 GPU accelerated compute	27

1.9	SBSA checklist	27
1.9.1	SBSA Level 3 checklist	28
1.9.2	SBSA Level 4 checklist	31
1.9.3	SBSA Level 5 checklist	32
1.9.4	SBSA Level 6 checklist	32
1.9.5	SBSA Level 7 checklist	33
1.9.6	SBSA Future Level checklist	35
<b>A</b>	<b>Level 3 - firmware</b>	<b>37</b>
A.1	Memory map	37
A.2	Clock and Timer Subsystem	37
A.3	Watchdog	38
<b>B</b>	<b>Performance Monitoring Unit</b>	<b>39</b>
B.1	Sampling	39
B.2	PE PMU	39
B.3	PE PMU events	40
B.3.1	Fractional cycle accounting	43
B.3.2	PE utilization	43
B.3.3	Top-down accounting	44
B.3.4	Workload events (SVE)	44
B.3.5	Branch predictor effectiveness	45
B.3.6	Latency	47
B.3.7	Memory workload	48
B.4	System performance monitors	51
B.4.1	Recommendations	53
<b>C</b>	<b>Server RAS</b>	<b>55</b>
C.1	Justifications and impact	55
C.1.1	PE architecture	55
C.1.2	System architecture components	55
C.1.3	Software faults	57
C.2	Server RAS architecture requirements	58
C.2.1	Software faults	59
C.2.2	Recommended RAS features	60
<b>D</b>	<b>Self-hosted Debug for Armv9-A</b>	<b>62</b>
D.1	Goals	62
D.2	Levels of functionality	62
D.3	Self-hosted debug capabilities	62
D.4	External debug capabilities	62
D.5	PE Trace	63
D.5.1	Introduction	63
D.5.2	Embedded Trace Extension (ETE)	63
D.5.3	ETE Level 1	63
D.6	System Trace Macrocell	64
D.6.1	Introduction	64
D.6.2	Level 1	64
<b>E</b>	<b>GPU accelerated compute</b>	<b>68</b>

Copyright © 2016-2022 Arm Limited. All rights reserved.

## Release information

### Version 7.1 (06 Oct 2022)

- Major additions are Errata for 7.0, Future requirements, and Appendix D and E.
- SBSA 7.1 version (542).
- Errata for 7.0 - PMU events issues for cache accesses (PMU\_EV\_02, PMU\_EV\_03, PMU\_EV\_04, PMU\_EV\_10, PMU\_EV\_11, RPMU\_SYS\_1, RPMU\_SYS\_2, RPMU\_SYS\_3, RPMU\_SYS\_4) (443).
- Errata for 7.0 - PMU branch events (PMU\_EV\_03, PMU\_EV\_10) (538).
- Errata for 7.0 - RAS rules errata (RAS\_02, RAS\_02A, RAS\_07, RAS\_08) (320).
- Errata for 7.0 - RAS\_01 errata (500).
- Errata for 7.0 - Enhancements to RAS\_01 and recommended RAS features for CE counter in self-correcting memory (517).
- Errata for 7.0 - RAS\_02 errata (501).
- Errata for 7.0 - RAS\_06 clarification (470).
- Errata for 7.0 - SYS\_RAS\_3 poison requirements seems to cancel each other out (326).
- Errata for 7.0 - S\_L7RAS\_1 typo (369).
- Errata for 7.0 - S\_L5SM\_04 typo (426).
- Errata for 7.0 - S\_L7TME\_2 typo (427).
- Errata for 7.0 - S\_L6PE\_01 is a duplicate of B\_SEC\_01-05 (552).
- Errata for 7.0 - Update PCIe terms and abbreviations (464).
- Errata for 7.0 - Add IHI 0091 reference (469).
- Errata for 7.0 - Fix RAS features heading section (472).
- Errata for 7.0 - Fix informative statements that are using must (487).
- Errata for 7.0 - Change Arm recommends and Arm Strongly Recommends terminology (523).
- Errata for 7.0 - Fix invalid rule ID links in the checklist
- SBSA future requirements Checklist (451).
- PE future requirements (378).
- Armv9 self-hosted debug (377).
- SPE PMBIDR\_EL1.F (336).
- PBHA not used in PE (406).
- Armv9 requirements with respect to SBSA Levels (459).
- RAS access rule for error nodes on shared resources (257)
- PCIe integration for GPU accelerated compute (440).
- GPU accelerated compute changes (529).
- PCIe MCTP VDM (288).
- PCIe RC external abort on errors (305).
- PCIe I/O Coherency rule (447).
- PCIe Steering Tag handling (424).
- PCIe IDE/CMA recommendation (482).
- PCIe AMBA Integration Guide reference (436).
- CXL related terminology and definitions (468).

### Version 7.0 (31 Jan 2021)

- Major additions are Level-7, Appendix B and C.
- Armv8.6 requirements (248).
- DMA masters behind SMMU (258).
- Entropy/TRNG requirement (270).
- SMMU-ATS when CXL (273).
- PCIe error handling (279).
- MPAM Level-7 (281).
- PCIe RP EP interoperability (289).

- PCIe miscellaneous rules (287).
- PMU requirements (156).
- RAS requirements (306).
- Branch PMU events (318).
- ERRADDR.AI bit in RAS (319).
- Errata for 6.0 - Pointer authentication (330).

#### **Version 6.1 (15 Sep 2020)**

- The document is reorganized to be a supplement of the Arm BSA document.

#### **Version 6.0 (16 Sep 2019)**

- Armv8.5 requirements (115).
- Instruction to data to instruction coherency (175).
- Nanosecond units for system counter (159).
- Rules for SMMU to ease page sharing between PE and SMMUs (158).
- MSI(-X) must be mapped to LPI (173).
- Mention ServerReady and GIC level 3 clarification (166).
- Clarification on this meaning MSI(-X) (103).
- Typo in word address, changed to addresses (102).
- Mislabelling of UARTIMSC in table 15 (101).
- IO virtualization clarifications (112).
- Armv8 RAS extension requirements (133).
- Relaxation of SVE requirement (162).
- Errata clarification on level 5 interrupt controller section of SBSA5.0B (104).
- Level 5 - Section 4.4.2 Interrupt Controller(117).
- SMMU and HTTU support (134).
- SBSA typo on Level 3 PE requirements for SVE (152).
- PCI enumeration related requirements (110).
- PCIe RCiEP and iEP related requirements (108).
- MSI support in SMMU for events (194).

#### **Version 5.0 (30 May 2018)**

- RAS requirements (5).
- ROP and JOB requirements for SBSA (6).
- SBSA and SVE (7).
- Nested virtualization and SBSA (8).
- MPAM requirements for SBSA (9).
- Invisible caches (ban type 2 caches) AKA Forced WB (10).
- Add activity monitor requirements to SBSA (11).
- Crypto requirements to SBSA (12).
- Add support for 48-bit operating systems booting on Armv8.2+ systems with 52-bit PA (13).
- Ban non-standard interrupt controller (14).
- Base frequency standardization (15).
- Assign PPIs for new timers in v8.4 (16).
- Secure EL2 not required (17).
- TLBI range homogeneity (18).
- SVE heterogeneity (19).
- PCIe clarifications (21).
- UART clarifications (22).
- PCIe requirements for assignable devices (23).
- PTM for SBSA-based systems (24).
- PCIe deadlocks (25).
- ACS and Peer-to-peer (26).
- TPM guidance for SBSA (27).
- Deprecate levels 0, 1, and 2 (28).

- Errata - remove references to Prince Algorithm (29).
- Errata - UART trigger levels (30).
- Heterogeneity in server (31).
- Clean to point of persistence support (33).
- Watchdog scaling in SBSA (36).
- Clarify that PEs must be able to access all Non-secure address space (37).
- Appendix I needs to be clearer about RID spaces do not strictly need to supply all 16 bits of width (38).
- Address space input into an SMMU from a device is entirely contiguous, no holes (39).
- Clarify the uniqueness requirements of SPIs that are used to handle legacy PCIe interrupts (A/B/C/D) (40).
- SBSA Armv8.4 and SMMUv3.2 rules on break before make for page tables (41).
- Add reservation for trace buffer overflow PPI (42).
- FP16 support in SBSA (43).
- Correction on stating devices describe their ability to be virtualized through FW (44).
- Clarify IO BAR usage (45).
- Generalize SMMU requirements (46).
- Relaxations in future GIC PPI reservation for Level 5 (79).

**Version 3.1 (02 Feb 2016)**

- Change of Proprietary Notice. Addition of release history. Non-confidential.

**Version 3.0 (01 Feb 2016)**

- Initial Release. Non-confidential.

## Arm Non-Confidential Document Licence (“Licence”)

This Licence is a legal agreement between you and Arm Limited (“**Arm**”) for the use of Arm’s intellectual property (including, without limitation, any copyright) embodied in the document accompanying this Licence (“**Document**”). Arm licenses its intellectual property in the Document to you on condition that you agree to the terms of this Licence. By using or copying the Document you indicate that you agree to be bound by the terms of this Licence.

“**Subsidiary**” means any company the majority of whose voting shares is now or hereafter owner or controlled, directly or indirectly, by you. A company shall be a Subsidiary only for the period during which such control exists.

This Document is **NON-CONFIDENTIAL** and any use by you and your Subsidiaries (“Licensee”) is subject to the terms of this Licence between you and Arm.

Subject to the terms and conditions of this Licence, Arm hereby grants to Licensee under the intellectual property in the Document owned or controlled by Arm, a non-exclusive, non-transferable, non-sub-licensable, royalty-free, worldwide licence to:

- (i) use and copy the Document for the purpose of designing and having designed products that comply with the Document;
- (ii) manufacture and have manufactured products which have been created under the licence granted in (i) above; and
- (iii) sell, supply and distribute products which have been created under the licence granted in (i) above.

**Licensee hereby agrees that the licences granted above shall not extend to any portion or function of a product that is not itself compliant with part of the Document.**

Except as expressly licensed above, Licensee acquires no right, title or interest in any Arm technology or any intellectual property embodied therein.

THE DOCUMENT IS PROVIDED “AS IS”. ARM PROVIDES NO REPRESENTATIONS AND NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, INCLUDING, WITHOUT LIMITATION, THE IMPLIED WARRANTIES OF MERCHANTABILITY, SATISFACTORY QUALITY, NON-INFRINGEMENT OR FITNESS FOR A PARTICULAR PURPOSE WITH RESPECT TO THE DOCUMENT. Arm may make changes to the Document at any time and without notice. For the avoidance of doubt, Arm makes no representation with respect to, and has undertaken no analysis to identify or understand the scope and content of, third party patents, copyrights, trade secrets, or other rights.

NOTWITHSTANDING ANYTHING TO THE CONTRARY CONTAINED IN THIS LICENCE, TO THE FULLEST EXTENT PERMITTED BY LAW, IN NO EVENT WILL ARM BE LIABLE FOR ANY DAMAGES, IN CONTRACT, TORT OR OTHERWISE, IN CONNECTION WITH THE SUBJECT MATTER OF THIS LICENCE (INCLUDING WITHOUT LIMITATION) (I) LICENSEE’S USE OF THE DOCUMENT; AND (II) THE IMPLEMENTATION OF THE DOCUMENT IN ANY PRODUCT CREATED BY LICENSEE UNDER THIS LICENCE). THE EXISTENCE OF MORE THAN ONE CLAIM OR SUIT WILL NOT ENLARGE OR EXTEND THE LIMIT. LICENSEE RELEASES ARM FROM ALL OBLIGATIONS, LIABILITY, CLAIMS OR DEMANDS IN EXCESS OF THIS LIMITATION.

This Licence shall remain in force until terminated by Licensee or by Arm. Without prejudice to any of its other rights, if Licensee is in breach of any of the terms and conditions of this Licence then Arm may terminate this Licence immediately upon giving written notice to Licensee. Licensee may terminate this Licence at any time. Upon termination of this Licence by Licensee or by Arm, Licensee shall stop using the Document and destroy all copies of the Document in its possession. Upon termination of this Licence, all terms shall survive except for the licence grants.

Any breach of this Licence by a Subsidiary shall entitle Arm to terminate this Licence as if you were the party in breach. Any termination of this Licence shall be effective in respect of all Subsidiaries. Any rights granted to any Subsidiary hereunder shall automatically terminate upon such Subsidiary ceasing to be a Subsidiary.

The Document consists solely of commercial items. Licensee shall be responsible for ensuring that any use, duplication or disclosure of the Document complies fully with any relevant export laws and regulations to assure that the Document or any portion thereof is not exported, directly or indirectly, in violation of such export laws.

This Licence may be translated into other languages for convenience, and Licensee agrees that if there is any conflict between the English version of this Licence and any translation, the terms of the English version of this Licence shall prevail.

The Arm corporate logo and words marked with ® or ™ are registered trademarks or trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere. All rights reserved. Other brands and names mentioned in this document may be the trademarks of their respective owners. No licence, express, implied or otherwise, is granted to Licensee under this Licence, to use the Arm trade marks in connection with the Document or any products based thereon. Visit Arm's website at <http://www.arm.com/company/policies/trademarks> for more information about Arm's trademarks.

The validity, construction and performance of this Licence shall be governed by English Law.

Copyright © 2016-2022 Arm Limited (or its affiliates). All rights reserved.

Arm Limited. Company 02557590 registered in England.

110 Fulbourn Road, Cambridge, England CB1 9NJ.

LES-PRE-21585 version 4.0

## About this document

### Terms and abbreviations

Term	Meaning
ACS	Access Control Services. A set of features intended to ensure that uncontrolled peer-to-peer transaction cannot occur. See [1].
AER	Advanced Error Reporting. A PCIe feature that enables software to isolate and analyze errors with fine granularity. See [1].
Arm ARM	Arm Architecture Reference Manual. See [2].
ATS	Address Translation Services
Base Server System	A system that is compliant with the Server Base System Architecture. See [3]
BBR	Base Boot Requirements [4].
CMA	PCIe Component Measurement and Authentication. See [1].
Completer	An agent in a computing system that responds to and completes a memory transaction that was initiated by a Requester.
CXL	Compute Express Link. See [5].
DMA	Direct Memory Access
ECAM	Enhanced Configuration Access Mechanism
FRU	Field Replaceable Unit
GIC	Generic Interrupt Controller
HDM	Host-managed Device Memory. See [5].
I/O coherent	A device is I/O coherent with the PE caches if its transactions snoop the PE caches for cacheable regions of memory. The PE does not snoop the device cache.
IDE	PCIe Integrity and Data Encryption. See [1].
LPI	Locality-specific Peripheral Interrupt (GICv3 [6]).
MEFN	Memory Error Forward Notification. See [5].
MMIO	Memory Mapped Input Output
MSC	Memory System Component.
NUMA	Non-uniform memory access
P2P or Peer-to-peer	See PCIe specification [1] for more details.
PCIe Host Bridge (PHB)	See PCIe specification [1] for more details.
PE	Processing Element, as defined in the Arm ARM. Typically a single hardware thread of a PE.
PMU	Performance Monitor Unit
PPI	Private Peripheral Interrupt

Term	Meaning
PRI	Page Request Interface
RCEC	Root Complex Event Collector. See [1].
RCiEP	Root Complex integrated End Point. See [1].
RCRB	Root Complex Register Base. Memory-mapped register space defined by a PCIe or CXL root complex.
Requester	An agent in a computing system that is capable of initiating memory transactions.
Root Complex (RC)	See PCIe specification [1] for more details.
Root Port (RP)	See PCIe specification [1] for more details.
SBBR	Server Base Boot Requirements [4].
SBSA	Server Base System Architecture.
SGI	Software-Generated Interrupt
SPI	Shared Peripheral Interrupt
SR-IOV	Single Root I/O virtualization. This is a method for a PCIe device to be virtualized. See [1].
System firmware data	System description data structures, for example ACPI or Flattened Device Tree.
VDM	PCI Express Vendor-defined Messages. See [1].
Viral	A hardware-based containment mechanism used in CXL. See [5].
VM	Virtual Machine

## References

This section lists publications by Arm and by third parties.

See Arm Developer (<http://developer.arm.com>) for access to Arm documentation.

- [1] *PCI Express Base Specification Revision 5.0, version 1.0*. PCI-SIG.
- [2] *DDI 0487 Arm® Architecture Reference Manual ARMv8, for the ARMv8-A architecture profile*. Arm Ltd.
- [3] *DEN 0094A Arm® Base System Architecture version 1.0*. Arm Ltd.
- [4] *DEN 0044 Arm® Base Boot Requirements*. Arm Ltd.
- [5] *CXL 2.0 specification*. [computeexpresslink.org](http://computeexpresslink.org).
- [6] *IHI 0069 Arm® Architecture Specification, GIC architecture version 3.0 and version 4.0*. Arm Ltd.
- [7] *DEN 0029C Arm® Server Base System Architecture Version 6.0*. Arm Ltd.
- [8] *Enhanced Allocation (EA) for Memory and I/O Resources*. PCI-SIG.
- [9] *DEN 0068 CoreSight Base System Architecture*. Arm Ltd.
- [10] *DDI 0587 Arm® Reliability, Availability and Serviceability (RAS) specification Armv8, for the Armv8-A architecture profile*. Arm Ltd.
- [11] *NIST 800-90*. NIST. "<https://csrc.nist.gov/publications/detail/sp/800-90b/final>".
- [12] *SBSA ACS*. Arm Ltd. "<https://github.com/ARM-software/sbsa-acs>".

- [13] *IHI 0070 Arm® System Memory Management Unit Architecture Specification, SMMU architecture version 3.3*. Arm Ltd.
- [14] *DEN 0114 ARM® PCIe AMBA Integration Guide*. Arm Limited.
- [15] *Arm® Architecture Reference Manual Supplement, The Scalable Vector Extension (SVE), for Armv8-A (ARMDDI 0584)*. Arm Ltd.
- [16] *IHI 0091 Arm® CoreSight Performance Monitoring Unit Architecture*. Arm Limited.
- [17] *DDI 0608 Arm® Architecture Reference Manual Supplement ARMv9, for ARMv9-A architecture profile*. Arm Ltd.
- [18] *IHI 0054 ARM® System Trace Macrocell Programmers Model Architecture*. Arm Limited.
- [19] *DDI 0416 Arm® CoreSight Trace Memory Controller Technical Reference Manual*. Arm Limited.
- [20] *ARM® Embedded Trace Router Architecture Specification*. Arm Limited.
- [21] *ARM® CoreSight System-on-Chip SoC-600 Technical Reference Manual*. Arm Limited.
- [22] *IHI 0029 ARM® CoreSight Architecture Specification*. Arm Limited.
- [23] *DEN 0034 ARM® Debug and Trace Configuration and Usage Models*. Arm Limited.

## Rules-based writing

This specification consists of a set of individual rules. Each rule is clearly identified by the letter R.

Rules must not be read in isolation, and where more than one rule relating to a particular feature exists, individual rules are grouped into sections and subsections to provide the proper context. Where appropriate, these sections contain a short introduction to aid the reader. An implementation which is compliant with the architecture must conform to all of the rules in this specification.

Some architecture rules are accompanied by rationale statements which explain why the architecture was specified as it was. Rationale statements are identified by the letter X.

Some sections contain additional information and guidance that do not constitute rules. This information and guidance is provided purely as an aid to understanding the architecture. Information statements are clearly identified by the letter I.

Implementation notes are identified by the letter U.

Software usage descriptions are identified by the letter S.

Arm strongly recommends that implementers read *all* chapters and sections of this document to ensure that an implementation is compliant.

Rules, rationale statements, information statements, implementation notes and software usage statements are collectively referred to as *content items*.

## Identifiers

Each content item may have an associated identifier which is unique within the context of this specification.

When the document is prior to beta status:

- Content items are assigned numerical identifiers, in ascending order through the document (*0001, 0002, ...*).
- Identifiers are volatile: the identifier for a given content item may change between versions of the document.

After the document reaches beta status:

- Content items are assigned random alphabetical identifiers (*HJQS, PZWL, ...*).
- Identifiers are preserved: a given content item has the same identifier across versions of the document.

## Examples

Below are examples showing the appearance of each type of content item.

R	This is a rule statement.
R <sub>x001</sub>	This is a rule statement.
I	This is an information statement.
X	This is a rationale statement.
U	This is an implementation note.
S	This is a software usage description.

## Feedback

Arm welcomes feedback on its documentation.

If you have comments on the content of this manual, send an e-mail to [errata@arm.com](mailto:errata@arm.com). Give:

- The title (Server Base System Architecture).
- The document ID and version (DEN0029H 7.1).
- The page numbers to which your comments apply.
- A concise explanation of your comments.

Arm also welcomes general suggestions for additions and improvements.

# 1 Server Base System Architecture

## 1.1 Arm Base System Architecture

This document is a supplement to the Arm Base System Architecture document version 1.0 [3].

## 1.2 Introduction

Server Base System Architecture (SBSA) specifies a hardware system architecture, based on Arm 64-bit architecture, that server system software, for example operating systems, hypervisors, and firmware can rely on. SBSA extends the requirements specified in the Arm BSA [3].

The Server Base System Architecture embeds the notion of levels of functionality. Each level adds functionality over and above a previous level. Unless explicitly stated, all specification items that belong to level N apply to levels greater than N.

All views of a level are required to be implemented to be compliant to a level.

SBSA Compliance is a requirement for the SR and LS bands of the Arm SystemReady program.

An implementation is consistent with a level of the Server Base System Architecture if it implements all of the functionality of that level at performance levels that are appropriate for the target uses of that level. This means that all functionality of a level can be exploited by software without unexpectedly poor performance.

The SBBR recipe in the Arm Base Boot Requirements (Arm BBR) specification [4] describes firmware requirements for an Arm server system. Where this specification refers to system firmware data, it refers to firmware specified in the Arm BBR.

### 1.2.1 SBSA levels and minimum component versions

Table 2 summarizes the minimum architecture versions that map to the rules that are mentioned in an4 SBSA level.

This table is indicative only. The rules in each level describe the specific features that are required to be compliant to that level. For a checklist of each level's minimum rules, refer to Section 1.9.

**Table 2: SBSA level mapping summary**

Level	PE: A profile	SMMU	GIC
3	v8.0	v2 or v3	v3.0
4	v8.3	v3.0	v3.0
5	v8.4	v3.2	v3.0
6	v8.5 or v9.0	v3.2	v3.0
7	v8.6 or v9.1	v3.2	v3.0

#### Note

Each SBSA level requires a set of PE features. These features were introduced in different revisions of the architecture. Generally, features that are available in version X of the architecture extension can be optionally implemented in version X-1. Some features can be backported even further. Please refer to extension documents for more details. Table 2 indicates the version of the architecture extension at which the required features were introduced.

---



---

#### Note

The levels between Level 3 and Level 6 that are presented in this document are semantically equivalent to the corresponding levels in SBSA 6.0 [7].

---

## 1.3 Level 3

- R<sub>S</sub>\_L3\_01 The rules from Arm BSA, [3] as listed in Section 1.9.1 must be implemented.
- I The PE Security requirements in Arm BSA [3] are relaxed to become a recommendation. This is to maintain consistency with systems which are certified as ServerReady v1.0.

### 1.3.1 PE architecture

- R<sub>S</sub>\_L3PE\_01 PEs must support 4KB and 64KB translation granules at stage 1 and stage 2.
- R<sub>S</sub>\_L3PE\_02 PEs must implement 16-bit ASID support.
- R<sub>S</sub>\_L3PE\_03 PEs must implement AArch64 at all implemented Exception levels.
- I Level 3 systems can be implemented using the Armv8.0 revision of the architecture, or higher. FEAT\_LPA [2], large Physical Address (PA) and Intermediate Physical Address (IPA) support, expands the maximum physical address width from 48 to 52 bits. Using 52-bit PA requires 64KB translation granule.
- R<sub>S</sub>\_L3PE\_04 Server base systems that make use of FEAT\_LPA must provide a mode where the memory map is wholly contained inside  $2^{48}$  bytes (256TB). This is required to support operating systems that do not use the 64KB granule.
- I It is consistent with this specification to implement PEs with support for the AArch32 Execution state.

### 1.3.2 Memory map

- R<sub>S</sub>\_L3MM\_01 To enable Non-secure EL2 hypervisors to use a 64KB translation granule at stage 2 MMU translation, the base system must ensure that all memory and peripherals can be mapped using 64KB stage 2 pages and must not require the use of 4KB pages at stage 2.
- R<sub>S</sub>\_L3MM\_02 Peripherals that will be assigned to different virtual machines will be situated within different 64KB regions of memory.

### 1.3.3 Interrupt controller

- R<sub>S</sub>\_L3GI\_01 A base server system must implement an interrupt controller that is compliant with the GICv3 or higher architecture [6].
- R<sub>S</sub>\_L3GI\_02 All MSI and MSI-X targeting hypervisor and operating system software must be mapped to LPI.

### 1.3.4 PPI assignments

R<sub>S</sub>\_L3PP\_01 The Interrupt IDs must be the same as the recommended values mentioned in Arm BSA [3].

### 1.3.5 System MMU and device assignment

I Armv8.4 introduces TLB Invalidation instructions which apply to a range of input addresses, instead of just to a single address.

R<sub>S</sub>\_L3SM\_01 If PEs that are used by the base system support TLB range instructions, then all OS visible requesters that contain a TLB must support range invalidates. See FEAT\_TLBIRANGE in [2].

### 1.3.6 Watchdog

R<sub>S</sub>\_L3WD\_01 The server base system must implement a Non-secure Generic watchdog as described in the Watchdogs section in Arm BSA [3].

## 1.4 Level 4

The list of rules to be implemented for Level 4 is available in Section 1.9.2.

### 1.4.1 PE architecture

In addition to the level 3 requirements, the following must be true of the PEs in the base server system:

R<sub>S</sub>\_L4PE\_01 All PEs must implement the RAS extension that is introduced in Armv8.2. See FEAT\_RAS in [2].

R<sub>S</sub>\_L4PE\_02 If the system contains persistent memory that is exposed to the OS, all PEs must support the clean to point of persistence instruction (DC CVAP). The instruction must be able to perform a clean to the point of persistence for all memory that is exposed as persistent memory to the OS.

R<sub>S</sub>\_L4PE\_03 All PEs must implement FEAT\_VMID16.

R<sub>S</sub>\_L4PE\_04 All PEs must implement FEAT\_VHE.

### 1.4.2 System MMU and device assignment

R<sub>S</sub>\_L4SM\_01 Stage 1 System MMU functionality must be provided by a System MMU that is compliant with the Arm SMMUv3, or higher, architecture revision.

R<sub>S</sub>\_L4SM\_02 Stage 2 System MMU functionality must be provided by a System MMU that is compliant with the Arm SMMUv3, or higher, architecture revision.

R<sub>S</sub>\_L4SM\_03 The integration of the System MMUs is compliant with the rules in SMMUv3 integration appendix in Arm BSA [3].

### 1.4.3 Peripheral subsystems

R<sub>S</sub>\_L4PCI\_1 All peripherals that are intended for assignment to a virtual machine or a user space device driver must be based on PCI Express.

R<sub>S</sub>\_L4PCI\_2 There must be no OS observable use of PCIe Enhanced Allocation [8].

## 1.5 Level 5

The list of rules to be implemented for level 5 is available in Section 1.9.3.

### 1.5.1 PE architecture

In addition to the level 4 requirements, the following must be true of the PEs in the base server system:

R <sub>S_L5PE_01</sub>	All PEs must support changing of translation table mapping size (FEAT_BBM) using the Level 1 or Level 2 solution that is proposed in the Armv8.4 extension. Level 2 is recommended. See Section 1.5.3 for the equivalent requirements for the SMMU.
R <sub>S_L5PE_02</sub>	All PEs must implement address authentication using the standard algorithm that is defined by the Arm architecture [2], as indicated by ID_AA64ISAR1_EL1.APA. See FEAT_PAuth in [2].
R <sub>S_L5PE_03</sub>	PEs that are based on Armv8.4 must implement the requirements of the CS-BSA combination C [9].
I	CS-BSA combination C[9] requires that PEs based on Armv8.x implement the Armv8.4 Self-hosted trace extension. See FEAT_TRF in [2].
I	For PEs that are based on Armv9 architecture, see Section D.
R <sub>S_L5PE_04</sub>	All PEs must implement the Activity Monitors Extension (AMU).
R <sub>S_L5PE_05</sub>	Where export control allows, all PEs must implement cryptography support for SHA3 and SHA512. See FEAT_SHA3 and FEAT_SHA512 in [2].

---

#### Note

It is recommended that FEAT\_SM3 and FEAT\_SM4 is also supported in hardware aimed at the Chinese market.

---

R <sub>S_MPAM_PE</sub>	Implementation of the MPAM extension (FEAT_MPAM) is OPTIONAL, however if implemented, the following minimal requirements must be met by the implementation: <ul style="list-style-type: none"> <li>• provide a minimum of 16 physical partition IDs.</li> <li>• provide virtualization support with a minimal of 8 virtual partition IDs.</li> <li>• provide a minimal of 2 performance monitor groups.</li> <li>• provide cache portion partitioning in the last level cache.</li> </ul>
I	Last level cache refers to an on SoC cache, as opposed to an off chip DRAM cache.

---

#### Note

It is recommended that the number of partition IDs scales with the number of PEs.

---

R <sub>S_L5PE_06</sub>	All PEs must provide support for stage 2 control of memory types and cacheability, as introduced by the Armv8.4 extensions. See FEAT_S2FWB in [2].
R <sub>S_L5PE_07</sub>	All PEs must implement enhanced nested virtualization, that is provided by HCR_EL2.NV2 and the VNCR_EL2 register. See FEAT_NV2 in [2].

### 1.5.2 Interrupt controller

R<sub>S\_L5GTI\_01</sub> Only a GIC v3, or higher, interrupt controller will be visible to operating system software. Other forms of interrupt controller, for example interrupt combining or forwarding engines, are not permissible, if they require a platform-specific kernel driver.

### 1.5.3 System MMU and device assignment

R<sub>S\_L5SM\_01</sub> SMMU implementations must be compliant with the Arm SMMUv3.2 architecture revision or higher.

R<sub>S\_L5SM\_02</sub> SMMU implementations must provide level 1 or level 2 support for translation table resizing.

---

#### Note

It is recommended that the SMMU implements Level 2. If the SMMU implementation provides Level 2, then it is recommended that the PE also provides level 2.

---

X MPAM architecture requires that all requesters that can access an MPAM-controlled resource, must support passing MPAM ID information.

R<sub>S\_L5SM\_03</sub> An SMMUv3.2 implementation must support the MPAM extension if the requests it serves access MPAM-controlled resources.

### 1.5.4 Clock and Timer subsystem

#### 1.5.4.1 Operating system

R<sub>S\_L5TI\_01</sub> A system that is compatible with level 5 will implement a generic counter which counts in nanosecond units. This means that, to the operating system, the reported frequency will be 1GHz.

I Systems are permitted to use the counter scaling (FEAT\_CNTSC) that is introduced in Armv8.4 as a method to implement this.

---

#### Note

It is recommended that this type of system uses revision 1 of the Generic watchdog. This is because at 1Ghz the watchdog timeout refresh period is limited to just over 4s.

---

### 1.5.5 Watchdog

I It is recommended that SBSA Level 5 systems use revision 1 of the Generic watchdog. This is because at 1Ghz the watchdog timeout refresh period is limited to just over 4s.

### 1.5.6 PPI assignments

R<sub>S\_L5PP\_01</sub> In addition to the PPI assignment specified in Arm BSA [3], the following PPIs are reserved by the SBSA specification:

**Table 3: PPI assignments**

Interrupt ID	Interrupt	Description
1056-1071	Reserved	Reserved for future SBSA usage.
1088-1103	Reserved	Reserved for future SBSA usage.

## 1.6 Level 6

The list of rules to be implemented for Level 6 is available in Section 1.9.4.

### 1.6.1 PE architecture

In addition to the Level 5 requirements, the following must be true of the PEs in the base server system:

- R The PE security requirements B\_SEC\_01, B\_SEC\_02, B\_SEC\_03, B\_SEC\_04, and B\_SEC\_05 as described in [3] are mandatory.
- R<sub>S\_L6PE\_02</sub> PEs must provide support for Branch Target Identification. Support is indicated by register ID\_AA64PFR1\_EL1.BT==b0001. See FEAT\_BTI in [2].
- R<sub>S\_L6PE\_03</sub> PEs must protect against timing faults being used to guess translation table mappings by implementing the TCR\_EL1.E0PD0 and TCR\_EL1.E0PD1 controls, and the same in TCR\_EL2. See FEAT\_E0PD in [2]. Support is indicated by ID register ID\_AA64MMFR2\_EL1.E0PD==b0001.
- R<sub>S\_L6PE\_04</sub> All PEs must implement FEAT\_PMUv3p5 [2].
- R<sub>S\_L6PE\_05</sub> Hardware updates to the Access flag and Dirty state in translation tables, as indicated by ID\_AA64MMFR1\_EL1.HAFDBS = 0b0010, must be supported. See FEAT\_HAFDBS in [2].

---

#### Note

It is recommended that the server base system removes the need for data cache clean for instruction to data coherency, and instruction invalidation for instruction to data coherency, as indicated by CTR\_EL0.IDC == 0b1 and CTR\_EL0.DIC == 0b1 respectively.

---

- R<sub>S\_L6PE\_06</sub> PEs must provide support for enhanced virtualization traps as indicated by ID\_AA64MMFR2\_EL1.EVT==b0010. See FEAT\_EVT in [2].

### 1.6.2 System MMU and device assignment

- R<sub>S\_L6SM\_01</sub> The SMMU must implement coherent access to memory structures, queues, and page tables as indicated by SMMU\_IDR0.COHAAC = 0b1.
- R<sub>S\_L6SM\_02</sub> The SMMU must support hardware translation table update (HTTU) of the Access flag and the Dirty state of the page for AArch64 translation tables, as indicated by SMMU\_IDR0.HTTU = 0b10.
- R<sub>S\_L6SM\_03</sub> The SMMU must support Message Signaled Interrupts as indicated by SMMU\_IDR0.MSI = 0b1.

### 1.6.3 Watchdog

- R<sub>S</sub>\_L6WD\_01** The Generic watchdog revision must be v1, W\_IIDR==0001b.
- X** In Generic watchdog v0, at 1 GHz the watchdog timeout refresh period is limited to just over 4 seconds.

### 1.6.4 Armv8 RAS extension requirements

- R<sub>S</sub>\_RAS\_01** PEs and other system components that implement the Armv8 RAS extension [10] must:
- Use Private Peripheral Interrupts for ERI or FHI if the only interface available for a RAS node is System register based.
  - Use the generic counter time base if the RAS Timestamp Extension is implemented. This means that ERR<n>FR.TS must be either b00 or b01.
- R<sub>S</sub>\_RAS\_02** Support for error injection is OPTIONAL, however it is recommended that if error injection is supported, the standard programming model that is described in the Arm RAS System Architecture version 1.1 is followed. Support for error injection is indicated by ERR<n>FR.INJ==b01.
- I** It is recommended that the RAS Timestamp Extension is implemented.

## 1.7 Level 7

The list of rules to be implemented for Level 7 is available in Section 1.9.5.

### 1.7.1 PE architecture

In addition to the Level 6 requirements, the requirements in this section must be true of all the PEs in the base server system.

- R<sub>S</sub>\_L7PE\_01** PEs must implement fine-grained trap support as indicated by ID\_AA64MMFR0\_EL1.FGT== b0001. See FEAT\_FGT in [2].
- R<sub>S</sub>\_L7PE\_02** PEs must implement the enhanced counter virtualization functionality as indicated by ID\_AA64MMFR0\_EL1.ECV == b0010. See FEAT\_ECV in [2].
- R<sub>S</sub>\_L7PE\_03** PEs must implement the enhancements to the Activity Monitor Unit (AMU) as indicated by ID\_AA64PFR0\_EL1.AMU == b0010. See FEAT\_AMUv1p1 in [2].
- R<sub>S</sub>\_L7PE\_04** PEs must implement the Advanced SIMD Int8 matrix multiply extension as indicated by: ID\_AA64ISAR1\_EL1.I8MM == b0001. See FEAT\_I8MM in [2].
- R<sub>S</sub>\_L7PE\_05** PEs must implement the BFLOAT16 extension. See FEAT\_BF16 in [2].
- R<sub>S</sub>\_L7PE\_06** PEs must implement the EnhancedPAC2 and FPAC extensions as indicated by ID\_AA64ISAR1\_EL1.APA == b0101. See FEAT\_PAuth2 in [2].
- R<sub>S</sub>\_L7PE\_07** Support for SVE is OPTIONAL. Where implemented, PEs must implement the SVE Int8 matrix multiply extension as indicated b ID\_AA64ZFR0\_EL1.I8MM ==b0001 See FEAT\_I8MM in [2].
- R<sub>S</sub>\_L7PE\_08** It is recommended that PEs implement the data gathering hint feature as indicated by ID\_AA64ISAR1\_EL1.DGH == b001. See FEAT\_DGH in [2].
- R<sub>S</sub>\_L7PE\_09** It is recommended that PEs implement the WFE fine-tuning delay feature as indicated by ID\_AA64MMFR0\_EL1.TWED == b0001. See FEAT\_TWED in [2].

R<sub>S\_L7PE\_10</sub> It is recommended that the enhanced PAN feature is implemented, which is introduced in Arm v8.7-A but can be implemented from v8.1-A. Support for enhanced PAN is indicated by ID\_AA64MMFR1\_EL1.PAN == b0011.

### 1.7.1.1 RAS

This section describes the requirements for the Reliability, Availability and Serviceability (RAS) extension. See FEAT\_RAS in [2].

R<sub>S\_L7RAS\_1</sub> For containable errors, error exceptions on reads from Normal memory must be taken as synchronous Data Abort exceptions.

R<sub>S\_L7RAS\_2</sub> Errors that are signaled to a PE on speculative accesses must not generate Abort exceptions.

I Containable errors are those which have not been silently propagated by the PE and can be taken as a containable error. Other error exceptions are taken as either synchronous data aborts or asynchronous SError interrupts.

### 1.7.1.2 Transactional Memory Extension

Implementation of Transactional Memory Extension (TME) is OPTIONAL. However, if implemented, the rules in this section apply.

R<sub>S\_L7TME\_1</sub> All PEs must have the same value of ID\_AA64ISAR0\_EL1.TME.

R<sub>S\_L7TME\_2</sub> The latency of starting and committing a transaction must not be higher than the latency of the code sequence that is recommended for acquiring and releasing a spinlock.

R<sub>S\_L7TME\_3</sub> For adequate performance of applications written in Java and C/C++, hardware must support a read set size of at least 512 objects and a write set size of at least 300 objects, assuming average object size to be 128 bytes.

R<sub>S\_L7TME\_4</sub> It is recommended that the hardware cache coherency facilities of the processor are used to detect transactional conflicts. This is also known as eager conflict detection.

R<sub>S\_L7TME\_5</sub> It is recommended that implementations do not generate a transactional conflict when a read generated by a PRFM instruction or by hardware prefetching, accesses a location within the transactional write set of a transaction.

## 1.7.2 MPAM

R<sub>S\_L7MP\_01</sub> PEs must implement the MPAM extension. See FEAT\_MPAM in [2].

R<sub>S\_L7MP\_02</sub> PEs must implement a minimum of 16 physical partition IDs.

R<sub>S\_L7MP\_03</sub> The implementation must provide MPAM Cache Storage Usage monitors for the last-level cache.

R<sub>S\_L7MP\_04</sub> Last-level cache must provide a minimum of 16 Cache Storage Usage monitors.

R<sub>S\_L7MP\_05</sub> The implementation must provide MPAM Memory Bandwidth Usage monitors (MBWUs) for the interfaces that provide general purpose memory.

I General purpose memory refers to normal DRAM or similar, in contrast to other memory interfaces, for example flash memory or storage class memories.

R<sub>S\_L7MP\_06</sub> The MBWUs must implement the MPAM v1.1 64-bit MBWU extension. See FEAT\_MPAMv1p1 in [2].

R<sub>S\_L7MP\_07</sub> The MBWUs for an interface must be sized so that they can count the total number of bytes that are transferred in a ten-second window when operating at maximum capacity.

- X The 64-bit MBWU extension allows counters to be either 44 bits or 63 bits. The effect of the above constraint is that:
- A 44-bit counter is sufficient for a single interface of up-to 1.6TB/s.
  - A 63-bit counter is necessary otherwise.
- I It is recommended that the number of MBWU scale with the number of PEs in the system.
- R<sub>S\_L7MP\_08</sub> The memory map of each Memory-System Component (MSC) that is accessible to Normal world software must be in global address space and have no overlap with other MSCs or peripherals.

### 1.7.3 Entropy source

- R<sub>S\_L7ENT\_01</sub> To support key and nonce generation, a system must have a hardware entropy source. This source must be a true random number generator that is visible to PE software and meets the requirements that are specified in the NIST SP 800-90 series of specifications [11], or the corresponding national equivalent.
- I NIST 800-22 statistical test suite for the validation of random numbers is included in the SBSA Architecture Compliance Suite version 2.4 and higher [12].

### 1.7.4 SMMU and device assignment

- I If an SMMU is in the path of a CXL device [5], the requirements are described in the SMMUv3 Architecture specification [13] chapter called SMMU interactions with CXL.
- X The path from a DMA requester to memory must carry the same protection as provided to a PE. An SMMU in the path from DMA requester to memory provides a standard interface to achieve this. [S\\_L7SM\\_01](#) enables a server to be deployed without limiting the use-cases involving I/O virtualization.
- R<sub>S\_L7SM\_01</sub> All DMA capable requesters that are visible to the normal world PE software must be behind an SMMU. This rule applies to types of requesters such as PCIe Root Ports and DMA capable requesters, for example USB, network, disk and accelerators.
- Here is an indicative list of the requesters that are exempt from this rule;
- Any on-chip requester whose resources cannot be controlled by PE software, for example system controllers or power management controllers.
  - Secure-world requesters.
  - Autonomous, DMA capable requesters which are part of the Trusted computing base of the system.
- This rule does not apply to SMMUs, interrupt controllers, or debug access ports which are DMA capable, but are explicitly forbidden from being behind an SMMU by design.
- I It is recommended that an ETR [9] is placed behind a SMMU.
- R<sub>S\_L7SM\_02</sub> If there is no SMMU in the path of an ETR, then a CATU as defined in the CoreSight-BSA specification [9] must be implemented for address translation.
- I If a CATU is implemented and the device driver for CATU is not present in the OS, trace cannot be captured.
- R<sub>S\_L7SM\_03</sub> SMMU must implement the SMMUv3 Performance Monitors Extension.
- R<sub>S\_L7SM\_04</sub> All SMMU Performance Monitor Counter Groups must implement at least four counters.
- I The SMMUv3 Architecture Specification [13] defines which events the System MMU Performance Monitor Extension must implement.

## 1.7.5 Performance Monitoring Unit

R<sub>S\_L7PMU</sub> The rules in Section B of this specification must be implemented.

## 1.7.6 System RAS

R<sub>SYS\_RAS</sub> The rules in Section C of this specification must be implemented.

### 1.7.6.1 Scrubbing

R<sub>SYS\_RAS\_1</sub> Each NUMA node in the base server system that implements error detection must support patrol scrubbing, supporting both:

- Continuous background scrub of the memory that is connected to the memory controller.
- Targeted scrub that allows run-time configuration of at least:
  - Patrol speed.
  - The region of physical memory connected to the memory controller being scrubbed.

Enabling targeted scrub should not require disabling of background scrub, nor affect the background scrub rate.

I For implemented patrol scrubbing, the following aspects are IMPLEMENTATION DEFINED:

- Programmable features might only be accessible to a system control agent, for example a system control processor, or might be accessible to the PEs.
- The supported patrol speed levels.
- The number of scrubbing channels and whether all channels are programmable.
- The smallest address space granule supported by the scrubbing engine.

It is recommended that control features and IMPLEMENTATION DEFINED capabilities are accessible to an operating system through a common interface, for example ACPI RASF. ACPI RASF supports three patrol speed levels: slow, medium, and fast.

### 1.7.6.2 Poison

R<sub>SYS\_RAS\_2</sub> The system must support the storage and forwarding of poisoned values.

I It is recommended that poison is generated when an uncorrectable error is detected.

For example, when:

- Corruption is detected in data being evicted from a cache, the cache evicts poison.
- Corruption is detected in data being read from external memory, the memory controller returns poison.
- Corruption is detected in data being received by a component on a write, the error is deferred if possible.
- A poisoned location is modified by a write that does not mask the poison value entirely. In this situation, the location remains poisoned.
  - If the write completely masks the poison, the poison can be removed. In either case, the location is modified and must be treated as a modified location.
- Poison is passed between components in the system with different poison schemes and/or error protection granule. The poison is propagated and expanded if required.

Note: Generating poison can *defer* the error.

- U Poison propagation can be achieved in different ways. For example:
1. Setting a poison indicator in the data
  2. Carry the error detection bits (for example, ECC, parity, or CRC) with the data
- R<sub>SYS\_RAS\_3</sub> The system should support propagating error indication with a transaction (per the definition in [10]) as the transaction propagates through the system. If a component propagating the transaction is unable to propagate the error indication, then the component must record the error as uncorrectable. On reads or non-posted writes, the component must generate one or more in-band error response(s), for example an External Abort or Completer Error. If enabled, an error recovery interrupt would be generated on both posted and non-posted transactions, as described by [10]. All coherent on-chip caches that hold dirty data and coherent addressable memories shall support the ability to retain a previously detected error with the erroneous data.
- I It is recommended that a distinct poison indication is saved with the data to avoid the same error being detected again in the future (as opposed to just “leaving bad ECC” with the data in the storage). Examples of a “distinct poison indication” include a separate poison bit or an invalid ECC code.
- X The intention of this rule is to require end-to-end propagation of some error indication between all PEs and all coherent on-chip caches that hold dirty data and coherent addressable memories.
- I Details of poisoning schemes, for example, the error protection granule and poison granule sizes, are IMPLEMENTATION DEFINED.
- I There is no architecturally-defined method to remove poison from a location, for example by restoring it to a known uncorrupted value. There is no requirement for this type of a method to be possible.
- It is recommended that if removing poison is supported, systems provide details of the method or a means to execute it in publicly available documentation.

## 1.7.7 PCIe

In addition to the rules in Arm BSA [3] for PCIe root ports and endpoints, the following rules must be implemented.

### 1.7.7.1 PCIe integration

- R<sub>S\_PCIE\_01</sub> The system must support the translation of PE writes with all byte enable patterns to PCIe write requests. The translation must be done in compliance with PCIe byte enable rules.
- R<sub>S\_PCIE\_02</sub> The Root Port must support the following:
- 1B and 2B read from Prefetchable and Non-prefetchable address spaces.
  - 1B and 2B write to Prefetchable and Non-prefetchable address spaces.
- This must hold true for accesses to downstream functions as well as to the Root Port itself. This must hold true even when the read or write address is not DW (4 Byte) aligned.
- R<sub>S\_PCIE\_03</sub> The Root Complex must:
- Send 2B PE writes that are 2B aligned as 2B PCIe writes.
  - Send 4B PE writes that are 4B aligned as 4B PCIe writes.
  - Send 8B PE writes that are 8B aligned as 8B PCIe writes.
- R<sub>S\_PCIE\_04</sub> The System must ensure that:
- Aligned 2B writes from Endpoints reach the target as 2B writes.
  - Aligned 4B writes from Endpoints reach the target as 4B writes.
  - Aligned 8B writes from Endpoints reach the target as 8B writes.

### 1.7.7.2 i-EP

- X This rule prevents the Root port PHY LTSSM from waiting indefinitely for the reception of a TS1 ordered set, with the Disable Link bit set before considering its lanes as disabled.
- R<sub>S\_PCIE\_05</sub> Root Port PHY LTSSM must have the capability to consider its lanes as Disabled when all the following conditions are met:
- Root port transmit side has sent TS1 ordered set with Disable bit set
  - Root port transmit side has sent the EIOSQ
  - Root port receive side has received EIOSQ from the downstream device or switch
- See section 4.2.6.9 in [1].

### 1.7.7.3 Error handling

- R<sub>PCI\_ER\_01</sub> Root Port must support the Advanced Error Reporting feature as described in section 6.2 and 7.8.4 of [1].
- R<sub>PCI\_ER\_02</sub> The Root Port must implement the MSI capability.
- R<sub>PCI\_ER\_03</sub> The Root Port must report reception of PCIe error messages using MSI interrupts. See section 6.2.6 and 6.2.4.2 of [1].
- I The Root Port error report includes errors from downstream devices, or error detected within the Root port.
- R<sub>PCI\_ER\_04</sub> The Root Port must log and report the PCIe errors it detects using the AER mechanism. See figure 6-3, section 6.2.6 in [1].
- R<sub>PCI\_ER\_05</sub> The Root Port must implement Downstream Port Containment feature. See section 6.2.10 in Reference [1].
- R<sub>PCI\_ER\_06</sub> The Root Port must comply with the rules stated in the PCIe specification for transaction layer behavior during DPC. See section 2.9.3 in [1].

## 1.8 Future requirements

The following section lists a preview of new rules that will be required in a future SBSA Level, which will be published in a future version of this document.

The list of rules to be implemented for future SBSA Level is available in Section 1.9.6.

### 1.8.1 PE architecture

In addition to the Level 7 requirements, the following must be true of the PEs in the base server system:

- R<sub>S\_L8PE\_01</sub> PEs must implement the XS attribute (FEAT\_XS) functionality as indicated by ID\_AA64ISAR1\_EL1.XS = 0b0001.
- R<sub>S\_L8PE\_02</sub> PEs must implement the support for WFET and WFIT instructions (FEAT\_WFXT) functionality as indicated by ID\_AA64ISAR1\_EL2.WFXT = 0b0001.
- R<sub>S\_L8PE\_03</sub> PEs must implement the support for atomic single-copy 64-byte stores with and without return (FEAT\_LS64, FEAT\_LS64\_V) as indicated by ID\_AA64ISAR1.LS64 == 0b0010.
- R<sub>S\_L8PE\_04</sub> PEs must implement the Enhanced PAN support (FEAT\_PAN3) as indicated by ID\_AA64MMFR1\_EL1.PAN == 0b0011.
- R<sub>S\_L8PE\_05</sub> PEs must implement the Armv8.7 enhancements for PMU (FEAT\_PMUv3p7) as indicated by ID\_DFR0.PerfMon == 0b0111.

**R<sub>S\_L8PE\_06</sub>** Implementation of FEAT\_BRBE that is introduced in Armv9.2 is OPTIONAL for PEs that are based on Armv9 architecture. However, if FEAT\_BRBE is implemented, the branch record buffer must support at least 32 branch records.

**R<sub>S\_L8PE\_07</sub>** PEs must not implement functionality for Page-Based Hardware Attributes (PBHA) bits of the VMSAv8-64 block and page descriptors. This is applicable to Stage-1 and Stage-2.

## 1.8.2 Self-hosted debug

The following must be true of the PEs in the base server system that are based on Armv9 architecture:

**R<sub>S\_L8SHD\_1</sub>** PEs must implement ETE Level-1 as described in Self-hosted Debug chapter, see Section D.5.3.

**I** Implementation of STM Level-1 as described in Self-hosted Debug, Section D.6, is OPTIONAL.

## 1.8.3 System MMU and device assignment

**R<sub>S\_L8SM\_01</sub>** SMMU implementations must be compliant with the Arm SMMUv3.3 architecture revision or higher.

## 1.8.4 System RAS

**R<sub>SYS\_RAS\_4</sub>** Resources that are shared by two or more PEs, and are implementing Armv8 RAS Extensions, must minimally support the memory-mapped view of the error nodes.

## 1.8.5 PCIe

### 1.8.5.1 PCIe integration

**R<sub>S\_PCIE\_06</sub>** For systems that support MCTP over PCIe VDM, the following rules must be implemented:

- Root Port to Root Port routing of MCTP VDMs must be supported by an SoC if it has more than one Root Port.
  - This rule also applies to the case where a single Root Port is split into multiple Root Ports with smaller link width by bifurcation/quadfurcation.
  - This rule applies even if the Root Ports are in different segments.

**R<sub>S\_PCIE\_07</sub>** All observers in the system must observe Inbound writes (for example, writes from downstream devices) with RO=0 in the order in which they were received by the Root Port.

**R<sub>S\_PCIE\_08</sub>** All RO=0 writes from a PCIe initiator (that is, under a root-port) must push all previous writes from that initiator to a point of visibility for all observers regardless of the RO value of the prior writes. PCIe write-ordering needs to apply for all observers in the system for requests initiated by PCIe.

**I** For more information and guidance on integrating PCIe in AMBA AXI based systems, refer to [14].

**R<sub>S\_PCIE\_09</sub>** All I/O subsystems that support PCI Express components must support I/O Coherency. See “I/O coherency: memory types and attributes for PCI Express” in [3].

**I** Either explicitly or via inheritance from the BSA [3], the SBSA has always required I/O subsystems that support PCI Express components to also support I/O Coherency. This explicit rule coincides with relaxing the BSA requirements.

**R<sub>S\_PCIE\_10</sub>** If Steering Tags are supported by the system, the STE.DCP (Directed Cache Prefetch) control bit in the System MMU must be honored to ensure the feature can be enabled and disabled architecturally.

R<sub>S\_PCIE\_11</sub> The mechanisms for supporting Steering Tags are IMPLEMENTATION DEFINED, but must have the following properties:

- Steering Tags must be treated as cache allocation hints and must only apply to Normal Cacheable memory transactions in the shareability domain of the I/O Subsystem.
- Steering Tags must not alter coherency guarantees.
- Steering Tags are permitted to be ignored.
- Steering Tag value properties for the platform must be discoverable by software as advertised by firmware.

I A Steering Tag value of 0 indicates no Steering Tag preference and is treated as if no Steering Tag is provided.

I It is recommended that the system support PCIe Integrity and Data Encryption (IDE) and PCIe Component Measurement and Authentication (CMA) to provide security features related to PCIe.

### 1.8.5.2 Error handling

R<sub>PCI\_ER\_07</sub> The system must either return all 1s data to the requestor PE or trigger an External abort in the requestor PE for each of the following cases:

- PE reads targeting MMIO space mapped as Normal memory gets a CA or UR response from the completer or times out (that is, has a Completion Time out (CTO)).
- PE reads targeting MMIO space mapped as Device memory gets a CA or UR response from the completer or times out (that is, has a Completion Time out).
- PE reads targeting configuration space gets a CA response from the completer or times out (that is, has a Completion Time out).
- These requirements apply even if the UR/CA is a response synthesized by the Root Port.
  - Note: The Root Port can synthesize an UR or CA response for outstanding read requests in multiple situations including:
    1. When the Root Port is in Downstream Port Containment (DPC) state (see section 2.9.3 in [1]).
    2. When the Root Port is in DL\_Down status (see section 2.9.1 in [1]).
- The selection between returning all 1 data and triggering an External abort must be independently configurable for MMIO space mapped as Normal memory, MMIO space mapped as Device memory and Configuration space for each error type (UR, CA and CTO for MMIO space and CA and CTO for Configuration space).

Note: MMIO space can be prefetchable or non-prefetchable memory space.

If an External abort is triggered in response to a PE request because of a PCIe error, then:

- Information about the error must be logged in an Arm RAS architecture compliant error record.
- For CA response, UR response and Completion time out, the SERR code to be logged in the error record is 23.
- UE and ER bits must be set to 1.
- DE and PN bits must be set to 0.
- UET value must be 0b11 (uncorrected error, signalled or recoverable).

Note: In addition to logging in the Arm RAS Error record, the error must be logged in AER and/or RP-PIO error logs (if RP-PIO feature is present) as mandated in the PCIe specification [1]

R<sub>PCI\_ER\_08</sub> The system must either return all 1s data or forward the poisoned data to the requestor PE for each of the following cases:

- PE reads targeting MMIO space mapped as Normal memory gets poisoned data back from the completer.
- PE reads targeting MMIO space mapped as Device memory gets poisoned data back from the completer.
- PE reads from configuration space gets poisoned data back from the completer.

Note: MMIO space can be prefetchable or non-prefetchable memory space.

The selection between returning all 1s data and forwarding poisoned data must be configurable.

- The default behavior (in the absence of any programming by software/firmware) must be all 1s data return.
- The selection between returning all 1s data and forwarding poisoned data must be independently configurable for MMIO space mapped as Normal memory, MMIO space mapped as Device memory and Configuration space.

If the system forwards poisoned response data received from downstream PCIe hierarchy to the requestor PE for a read request, then the event must be logged in an Arm RAS error record.

- The SERR code to be used is 25.
- UE, PN and ER bits must be set to 1.
- DE bit must be set to 0.
- UET value must be 0b11 (uncorrected error, signalled or recoverable).

Note: In addition to logging in the Arm RAS Error record, the error must be logged in AER error logs as mandated in the PCIe specification [1].

**R<sub>PCI\_ER\_09</sub>** The value of “RP Extensions for DPC” bit in the DPC Capability register (see section 7.9.15.2 in [1]) must be configurable by firmware.

Note: The OS must still see this bit as Read Only.

**R<sub>PCI\_ER\_10</sub>** The Root Port must be capable of triggering DPC if the system returns all 1s data for the following error cases:

- PE reads targeting MMIO space mapped as Normal memory gets a CA or UR response from the completer.
  - If RP-PIO is not implemented or not exposed to the OS, then triggering DPC must be achieved by treating the CA or UR as an uncorrectable internal error (see section 6.2.9 in [1]).
- PE reads targeting MMIO space mapped as Device memory gets a CA or UR response from the completer.
  - If RP-PIO is not implemented or not exposed to the OS, then triggering DPC must be achieved by treating the CA or UR as an uncorrectable internal error.
- PE reads targeting configuration space gets a CA response from the completer.
  - If RP-PIO is not implemented or not exposed to the OS, then triggering DPC must be achieved by treating the CA as an uncorrectable internal error.

**I** For a posted request from downstream PCIe requestors, the following is recommended:

- If the Root Complex detects that the request’s data payload is corrupted, then:
  - Root Complex should poison the data and forward it to the target if the system supports poisoned data forwarding from the Root Complex.
  - Root Complex should log and report an error if the system does not support poisoned data forwarding from the Root Complex.
- If the Root Complex detects an error that cannot be isolated to the write data, then it should log and report the error.
- The logging and reporting can be through the PCIe error reporting mechanisms or through the Arm RAS architecture specified mechanisms.

### 1.8.6 GPU accelerated compute

**R** Systems with GPU accelerated compute must implement the rules in Section E.

## 1.9 SBSA checklist

This section lists the minimum hardware requirements that are required to install, boot, and run a server operating system on bare-metal or within a virtualization environment.

All rule IDs which begin with 'B' are defined in the Arm BSA document [3].

### 1.9.1 SBSA Level 3 checklist

PE	Rule ID
Operating system	
	B_PE_01
	B_PE_02
	B_PE_03
	B_PE_04
	B_PE_05
	B_PE_06
	B_PE_07
	B_PE_08
	B_PE_09
	B_PE_10
	B_PE_11
	B_PE_12
	B_PE_13
	B_PE_14
	<a href="#">S_L3PE_01</a>
	<a href="#">S_L3PE_02</a>
	<a href="#">S_L3PE_03</a>
	<a href="#">S_L3PE_04</a>
Hypervisor	
	B_PE_18
	B_PE_19
	B_PE_20
	B_PE_21
	B_PE_22
Platform security	
	B_PE_23
	B_PE_24
Memory map	
Rule ID	
Operating system	

Memory map	Rule ID
	B_MEM_01
	B_MEM_02
	B_MEM_03
	B_MEM_04
	B_MEM_05
	B_MEM_06
	B_MEM_07
	S_L3MM_01
	S_L3MM_02
Platform security	
	B_MEM_08
	B_MEM_09
Interrupts	Rule ID
Operating system	
	S_L3GI_01
	S_L3GI_02
	B_GIC_03
	B_GIC_04
	B_GIC_05
	S_L3PP_01
	B_PPI_01
Hypervisor	
	B_PPI_02
Platform security	
	B_PPI_03
SMMU	Rule ID
Operating system	
	B_SMMU_01
	B_SMMU_02
	B_SMMU_06
	B_SMMU_07

SMMU	Rule ID
	B_SMMU_08
	B_SMMU_12
	S_L3SM_01
Hypervisor	
	B_SMMU_16
	B_SMMU_17
	B_SMMU_18
	B_SMMU_19
	B_SMMU_21

Timer subsystem	Rule ID
Operating system	
	B_TIME_01
	B_TIME_02
	B_TIME_03
	B_TIME_04
	B_TIME_05
	B_TIME_06
	B_TIME_07
	B_TIME_08
	B_TIME_09
	B_TIME_10

Power and wakeup	Rule ID
	B_WAK_01
	B_WAK_02
	B_WAK_03
	B_WAK_04
	B_WAK_05
	B_WAK_06
	B_WAK_07
	B_WAK_08
	B_WAK_09

Power and wakeup	Rule ID
	B_WAK_10
	B_WAK_11

  

Peripherals	Rule ID
Operating system	B_PER_01
	B_PER_02
	B_PER_03
	B_PER_04
	B_PER_05
	B_PER_06
	B_PER_07
	B_PER_08
	B_PER_09
	B_PER_10
	B_PER_12
Platform security	B_PER_11

  

Watchdog	Rule ID
	B_WD_01
	B_WD_02
	B_WD_03
	B_WD_04
	B_WD_05

### 1.9.2 SBSA Level 4 checklist

In addition to the SBSA Level 3 rules in Section 1.9.1, the following additional rules are required.

Module	Rule ID
PE	<a href="#">S_L4PE_01</a>
PE	<a href="#">S_L4PE_02</a>
PE	<a href="#">S_L4PE_03</a>

Module	Rule ID
PE	S_L4PE_04
SMMU	<del>B_SMMU_08</del>
SMMU	S_L4SM_01
SMMU	S_L4SM_02
SMMU	S_L4SM_03
PCIe	S_L4PCI_1
PCIe	S_L4PCI_2

### 1.9.3 SBSA Level 5 checklist

In addition to the SBSA Level 4 rules in Section 1.9.2, the following additional rules are required.

Module	Rule ID
PE	S_L5PE_01
PE	S_L5PE_02
PE	S_L5PE_03
PE	S_L5PE_04
PE	S_L5PE_05
PE	S_L5PE_06
PE	S_L5PE_07
PE	S_MPAM_PE
GIC	S_L5GI_01
SMMU	<del>S_L4SM_01</del>
SMMU	S_L5SM_01
SMMU	S_L5SM_02
SMMU	S_L5SM_03
SMMU	B_SMMU_09
SMMU	B_SMMU_11
SMMU	B_SMMU_20
SMMU	B_SMMU_22
Timer	S_L5TI_01
Interrupt	S_L5PP_01

### 1.9.4 SBSA Level 6 checklist

In addition to the SBSA Level 5 rules in Section 1.9.3, the following additional rules are required.

Module	Rule ID
PE	B_PE_16
PE	B_PE_17
PE	S_L6PE_01
PE	S_L6PE_02
PE	S_L6PE_03
PE	S_L6PE_04
PE	S_L6PE_05
PE	S_L6PE_06
PE	B_SEC_01
PE	B_SEC_02
PE	B_SEC_03
PE	B_SEC_04
PE	B_SEC_05
SMMU	B_SMMU_03
SMMU	B_SMMU_04
SMMU	B_SMMU_05
SMMU	B_SMMU_13
SMMU	B_SMMU_14
SMMU	B_SMMU_23
SMMU	S_L6SM_01
SMMU	S_L6SM_02
SMMU	S_L6SM_03
Watchdog	S_L6WD_01
RAS	S_RAS_01
PCI	B_REP_1
PCI	B_IEP_1

### 1.9.5 SBSA Level 7 checklist

In addition to the SBSA Level 6 rules in Section 1.9.4, the following additional rules are required.

Module	Rule ID
PE	S_L7PE_01
PE	S_L7PE_02
PE	S_L7PE_03
PE	S_L7PE_04

Module	Rule ID
PE	S_L7PE_05
PE	S_L7PE_06
PE	S_L7PE_07
PE	S_L7PE_08
PE	S_L7PE_09
PE	S_L7PE_10
PE RAS	S_L7RAS_1
PE RAS	S_L7RAS_2
PE TME	S_L7TME_1
PE TME	S_L7TME_2
PE TME	S_L7TME_3
PE TME	S_L7TME_4
PE TME	S_L7TME_5
PE MPAM	S_L7MP_01
PE MPAM	S_L7MP_02
PE MPAM	S_L7MP_03
PE MPAM	S_L7MP_04
PE MPAM	S_L7MP_05
PE MPAM	S_L7MP_06
PE MPAM	S_L7MP_07
PE MPAM	S_L7MP_08
Entropy	S_L7ENT_1
SMMU	S_L7SM_01
SMMU	S_L7SM_02
SMMU	S_L7SM_03
SMMU	S_L7SM_04
PMU	S_L7PMU
RAS	SYS_RAS
RAS	SYS_RAS_1
RAS	SYS_RAS_2
RAS	SYS_RAS_3
PCIe	S_PCIE_01
PCIe	S_PCIE_02
PCIe	S_PCIE_03
PCIe	S_PCIE_04

Module	Rule ID
PCIe	S_PCIe_05
PCIe	PCI_ER_01
PCIe	PCI_ER_02
PCIe	PCI_ER_03
PCIe	PCI_ER_04
PCIe	PCI_ER_05
PCIe	PCI_ER_06

## 1.9.6 SBSA Future Level checklist

### 1.9.6.1 Base checklist for SBSA Future Level

In addition to the SBSA Level 7 rules in Section 1.9.5, the following additional rules are required.

Module	Rule ID
PE	S_L8PE_01
PE	S_L8PE_02
PE	S_L8PE_03
PE	S_L8PE_04
PE	S_L8PE_05
PE	S_L8PE_07
SMMU	S_L8SM_01
RAS	SYS_RAS_4
PCIe	S_PCIe_06
PCIe	S_PCIe_07
PCIe	S_PCIe_08
PCIe	S_PCIe_09
PCIe	S_PCIe_10
PCIe	S_PCIe_11
PCIe	PCI_ER_07
PCIe	PCI_ER_08
PCIe	PCI_ER_09
PCIe	PCI_ER_10
GPU	GPU_01
GPU	GPU_02
GPU	GPU_03

---

Module	Rule ID
GPU	<a href="#">GPU_04</a>

---

### **1.9.6.2 Additional checklist for SBSA Future Level for ARMv9-A profile**

In addition to the SBSA Future Level rules in Section [1.9.6.1](#), the following rules are required for servers with PEs that are based on ARMv9 architecture.

---

Module	Rule ID
PE	<a href="#">S_L8PE_06</a>
PE	<a href="#">S_L8SHD_1</a>

---

## A Level 3 - firmware

Level 3 - firmware is an optional additional set of requirements for a Level 3 system. Level 3 - firmware is designed to give a base set of functionality that standard platform firmware can rely on. A system that is compliant with level 3 and not compliant with level 3 - firmware is still a fully compliant level 3 system. The system has all the features required by the operating systems and hypervisors.

### A.1 Memory map

The system must provide some memory mapped in the Secure address space. The memory must not be aliased in the Non-secure address space. The amount of Secure memory provided is platform-specific as the intended use of the memory is for platform-specific firmware.

All Non-secure on-chip masters in a base server system that are expected to be used by the platform firmware must be capable of addressing all of the Non-secure address space. If the master goes through a SMMU then the master must be capable of addressing all of the Non-secure address space even when the SMMU is off.

### A.2 Clock and Timer Subsystem

A system compatible with level 3 - firmware must also include a Secure wakeup timer in the form of the memory mapped timer described in the Armv8 ARM [2]. This timer must be mapped into the Secure address space, and the timer expiry interrupt must be presented to the GIC as an SPI. This timer is called the Secure system wakeup timer.

The Secure wakeup timer does not require a virtual timer to be implemented. The virtual offset register can Read-As-Zero, where writes to the virtual offset register in CNTCTLBase frame are ignored. The timer is not required to have a CNTELOBase frame.

The following table summarizes which address space the register frames related to the Secure wakeup timer should be mapped on to.

Register Frame	
CNTControlBase	Secure
CNTReadBase	Not required
CNTCTLBase	Secure
CNTBaseN	Secure

CNTCTLBase might be shared among multiple timers, including various Secure and Non-secure timers. The SBSA specification does not require this.

#### Note

- GICv3 does not support Secure LPI. Therefore, the Secure system timer interrupt must not be delivered as LPI.
- It is recognized that in a large system, a shared resource like the system wakeup timer can create a

system bottleneck. This is because, access to it must be arbitrated through a system-wide lock. Level 3-firmware requires just a single timer so that standard firmware implementations have a guaranteed timer resource across platforms. We anticipate that large PE systems will implement a more scalable solution like one timer for each PE.

---

## A.3 Watchdog

The required behavior of watchdog signal 1 of the Non-secure watchdog is modified in level 3 - firmware and is required to be routed as an SPI to the GIC. It is expected that this SPI be configured as an EL3 interrupt, directly targeting a single PE.

A system compatible with level 3- firmware must implement a second watchdog, and is referred to as the Secure watchdog. It must have both its register frames mapped in the Secure memory address space and must not be aliased to the Non-secure address space.

Watchdog Signal 0 of the Secure watchdog must be routed as an SPI to the GIC and it is expected this will be configured as an EL3 interrupt, directly targeting a single PE.

---

### Note

- GICv3 does not support Secure LPI. The Secure watchdog interrupts must not be delivered as LPI.
  - Only directly-targeted SPIs are required to wake a PE. Programming the watchdog SPI to be directly targeted ensures delivery of the interrupt independent of PE power states. However, it is possible to use a 1 of N SPI to deliver the interrupt, if one of the target PEs is running. See Arm BSA [3] for information about SPI waking a PE.
- 

Watchdog Signal 1 of the Secure watchdog must be routed to the platform. In this context, platform means any entity that is more privileged than the code running at EL3. Examples of the *platform* component that services Watchdog Signal 1 are a system control processor, or dedicated reset control hardware.

The action that is taken on the raising of Watchdog Signal 1 of the Secure watchdog is platform-specific.

## B Performance Monitoring Unit

Performance monitoring features have two main use models:

- Profiling and software optimization
- Monitoring

Profiling is a tool that is used in software optimization. Performance Monitors are used by profiling. Monitoring is a tool that is used in system operations. Performance Monitors are a hardware feature used in monitoring.

Both approaches typically use sampling to read the PMU.

### B.1 Sampling

There are two main sampling models for Hardware Performance Monitors (HPM):

- Time-based sampling:

Software periodically records values from the HPM and records the location in the program. The changes are tracked over time through phases of software execution. The engineer looks for correlations between phases and recorded events.

- Event-based sampling:

The location in the program, or some other measurement, is recorded when a *sampled* event occurs. This builds up a statistical model of where events occur.

### B.2 PE PMU

The requirements in this section must be met by all the PEs in the base server system.

R<sub>PMU\_PE\_01</sub>

The PE implements the Performance Monitors Extension.

R<sub>PMU\_PE\_02</sub>

The PMU overflow signal from each PE must be wired to a unique PPI interrupt with no intervening logic.

R<sub>PMU\_PE\_03</sub>

Each PE must implement a minimum of six PMU event counters and the PMU cycle counter.

I

Performance monitors are required for PE local caches and the last level cache. They are OPTIONAL for intermediate cache levels that are not local to the PE.

I

The Performance Monitors Extension requires that at least one of the INST\_RETIRED and INST\_SPEC events is implemented. It is recommended that the INST\_RETIRED event is implemented.

R<sub>PMU\_EV\_11</sub>

PEs in the base server system must either:

- Not implement any multithreaded PMU extension. `PMEVTYPEPER<n>_EL0.MT` are `RES0`. `ID_AA64DFR0_EL1.MTPMU == 0b1111`.
- Implement the ARMv8.6-MTPMU extension. `ID_AA64DFR0_EL1.MTPMU == 0b0001`.

R<sub>PMU\_SPE</sub>

Implementation of the SPE is optional. However, if SPE is implemented, the following requirements must be met:

- PEs must support Hardware management of the Access Flag and dirty state for accesses made by the SPE[2], as indicated by `PMBIDR_EL1.F == 01b`.
- If the SPE implementation samples micro-operations, the implementation provides public documentation of the effect of such sampling on the weighting of instructions in the sample population.

- If the SPE implementation samples the Data Source indicator, the implementation provides public documentation of the mappings of the Data Source values to data sources.
- If the SPE implementation includes IMPLEMENTATION DEFINED packets or IMPLEMENTATION DEFINED events in the events packet, the implementation provides public documentation of these packets and events.

I It is recommended that the SPE is implemented.

---

#### Note

PMU\_SPE relates to the following from [2]:

- An architecture instruction might create one or more micro-ops at any point in the execution pipeline.
  - The definition of a micro-op is IMPLEMENTATION DEFINED.
  - An architecture instruction might create more than one micro-op for each instruction.
  - A micro-op might also be removed or merged with another micro-op in the execution stream. This means that an architecture instruction might create no micro-ops for an instruction.
  - Any arbitrary translation of architecture instructions to an equivalent sequence of micro-ops is permitted.
  - In some implementations, the relationship between architecture instructions and micro-ops might vary over time. For example, an instruction that generates two micro-ops is twice as likely to be sampled as an instruction that generates a single micro-op.
  - It is recommended that applicable implementation details are provided in freely-usable, machine-readable standard formats.
- 

## B.3 PE PMU events

R<sub>PMU\_EV\_01</sub> PEs in the base server system must implement the PMU events that are shown in the following table to measure IPC:

Number	Mnemonic	Description
0x0008	INST_RETIRED	Instruction architecturally executed
0x0011	CPU_CYCLES	Cycles

R<sub>PMU\_EV\_02</sub> PEs in the base server system must implement the PMU events that are shown in the following table to measure cache effectiveness, for each PE local cache and the last level cache:

Number	Mnemonic	Description
0x0004	L1D_CACHE	Level 1 data cache access
0x0013	MEM_ACCESS	Data memory access
0x0014	L1I_CACHE	Level 1 instruction cache access
0x0016	L2D_CACHE	Level 2 data cache access
0x0027	L2I_CACHE	Level 2 instruction cache access

Number	Mnemonic	Description
0x002B	L3D_CACHE	Level 3 data cache access
0x0036	LL_CACHE_RD	Last Level cache memory read
0x0037	LL_CACHE_MISS_RD	Last Level cache memory read miss
0x0039	L1D_CACHE_LMISS_RD	Level 1 data cache long-latency read miss
0x0040	L1D_CACHE_RD	Level 1 data cache access, read
0x0050	L2D_CACHE_RD	Level 2 data cache access, read
0x0066	MEM_ACCESS_RD	Data memory access, read
0x00A0	L3D_CACHE_RD	Level 3 data cache access, read
0x4006	L1I_CACHE_LMISS	Level 1 instruction cache long latency miss
0x4009	L2D_CACHE_LMISS_RD	Level 2 data cache long-latency read miss
0x400A	L2I_CACHE_LMISS	Level 2 instruction cache long latency miss
0x400B	L3D_CACHE_LMISS_RD	Level 3 data cache long-latency read miss

L<n>D and L<n>I cache events that are not for PE local caches are not required.

#### Note

The definitions of which caches are PE local and which cache is the last level cache are IMPLEMENTATION DEFINED. The last level cache usually refers to the last level of cache before memory.

PE events are attributable to the PE counting the events.

X

This is the rationale for the rule [PMU\\_EV\\_03](#).

- Effectiveness: For example, a cache is *effective* if accesses to the cache have low latency because they do not miss in the cache.

Counting events, for example a cache miss, that indicate effectiveness or ineffectiveness might be more cost efficient than measuring actual latency.

- TLB effectiveness

TLBs are designed to improve performance by caching the results of translation table walks. Therefore, a TLB is not effective if accesses miss in the TLB and cause a translation table walk.

The following events are designed for monitoring effectiveness of the TLBs:

- L1D\_TLB
  - DTLB\_WALK
- L1I\_TLB
  - ITLB\_WALK

R<sub>PMU\_EV\_03</sub>

The events that are listed in the following table are recommended to be implemented by PEs in the base server system for measuring TLB effectiveness, for each applicable level of PE TLB:

Number	Mnemonic	Description
0x0025	L1D_TLB	Level 1 data TLB access

Number	Mnemonic	Description
0x0026	L1I_TLB	Level 1 instruction TLB access
0x0034	DTLB_WALK	Data TLB access with at least one translation table walk
0x0035	ITLB_WALK	Instruction TLB access with at least one translation table walk

R<sub>PMU\_EV\_04</sub>

If the L<n>D and L<n>I cache and TLB events count both demand and non-demand accesses, then the PE must implement additional events that count only demand accesses.

It is recommended that following events are defined for this:

Number	Mnemonic	Description
0x8130	L1D_TLB_RW	Level 1 data TLB demand access
0x8131	L1I_TLB_RD	Level 1 instruction TLB demand access
0x813C	DTLB_WALK_RW	Data TLB demand access with at least one translation table walk
0x813D	ITLB_WALK_RD	Instruction TLB demand access with at least one translation table walk
0x8140	L1D_CACHE_RW	Level 1 data cache demand access
0x8141	L1I_CACHE_RD	Level 1 instruction cache demand access
0x8148	L2D_CACHE_RW	Level 2 data cache demand access
0x8149	L2I_CACHE_RD	Level 2 instruction cache demand access
0x8150	L3D_CACHE_RW	Level 3 data cache demand access

However, implementations can also meet this requirement using IMPLEMENTATION DEFINED events. For the purposes of this rule, demand access is as defined for these events in [2].

X

The following events are designed for accounting for cycles.

- CPU\_CYCLES
  - STALL
    - \* STALL\_FRONTEND
  - STALL
    - \* STALL\_BACKEND
      - STALL\_BACKEND\_MEM
  - Not stalling

R<sub>PMU\_EV\_05</sub>

The events listed in the following table must be implemented for cycle accounting:

Number	Mnemonic	Description
0x0011	CPU_CYCLES	Cycle
0x0023	STALL_FRONTEND	No operation sent for execution due to the frontend
0x0024	STALL_BACKEND	No operation sent for execution due to the backend
0x003C	STALL	No operation sent for execution

Number	Mnemonic	Description
0x4005	STALL_BACKEND_MEM	Memory stall cycles

### B.3.1 Fractional cycle accounting

X

This is the rationale for the rule [PMU\\_EV\\_06](#).

The following events are designed for accounting for fractional cycles:

- CPU\_CYCLES\_SLOTS
  - STALL\_SLOT
    - \* STALL\_SLOT\_FRONTEND
  - STALL\_SLOT
    - \* STALL\_SLOT\_BACKEND
  - *Not stalling*

The IMPLEMENTATION\_DEFINED constant SLOTS is discoverable from the system register PMMIR\_EL1.SLOTS. It is the maximum number that STALL\_SLOT can increment by in a single [CPU\\_CYCLE](#).

Fractional events are similar in concept to cycle accounting events, but focus on the resources of the CPU. For example, if a CPU can issue up to three instructions in a cycle, then the STALL\_SLOT events can count up-to three per cycle.

---

#### Note

Another way to view this is that the CPU is being profiled like it issues one instruction per cycle, but operates at three times the frequency.

---

These resources are referred to as *slots*.

This top-down approach allows for subdivision of these resources, but does not necessarily require strict accounting of resources.

### B.3.2 PE utilization

It is useful to break down the operations that are issued for execution into:

- Operations that are architecturally retired
- Operations that only perform wasted work, for example, because of speculation.

The following events are designed for analyzing the impact of speculation:

- OP\_SPEC
  - OP\_RETIRED

Here are some measurements to profile the PE utilization.

- Fraction of operations retired

$$f_{\text{retired}} = \frac{\text{OP\_RETIRED}}{\text{OP\_SPEC}}$$

- Operations wasted, for example, due to branch misprediction

$$OP\_WASTED = OP\_SPEC - OP\_RETIRED$$

- Fraction of operations wasted

$$f_{\text{wasted}} = 1 - \frac{OP\_RETIRED}{OP\_SPEC}$$

- Utilization of CPU

$$\rho_{\text{CPU}} = \left(1 - \frac{STALL\_SLOT}{CPU\_CYCLES \times SLOTS}\right) \times \left(\frac{OP\_RETIRED}{OP\_SPEC}\right)$$

### B.3.3 Top-down accounting

Top-down accounting combines the Section B.3.1 and Section B.3.2 methodologies. Together these methodologies allow for four key measurements to be extracted:

$$f_{\text{frontend-bound}} = \frac{STALL\_SLOT\_FRONTEND}{CPU\_CYCLES \times SLOTS}$$

$$f_{\text{backend-bound}} = \frac{STALL\_SLOT\_BACKEND}{CPU\_CYCLES \times SLOTS}$$

$$f_{\text{retired}} = \left(\frac{OP\_RETIRED}{OP\_SPEC}\right) \times \left(1 - \frac{STALL\_SLOT}{CPU\_CYCLES \times SLOTS}\right)$$

$$f_{\text{wasted}} = \left(1 - \frac{OP\_RETIRED}{OP\_SPEC}\right) \times \left(1 - \frac{STALL\_SLOT}{CPU\_CYCLES \times SLOTS}\right)$$

R<sub>PMU\_EV\_06</sub>

The events in the following table must be implemented for Section B.3.3:

Number	Mnemonic	Description
0x0011	CPU_CYCLES	Cycle
0x003A	OP_RETIRED	Micro-operation architecturally executed
0x003B	OP_SPEC	Micro-operation speculatively executed
0x003D	STALL_SLOT_BACKEND	No operation sent for execution on a slot due to the backend
0x003E	STALL_SLOT_FRONTEND	No operation sent for execution on a slot due to the frontend
0x003F	STALL_SLOT	No operation sent for execution on a slot

### B.3.4 Workload events (SVE)

X

This is the rationale for the rule [PMU\\_EV\\_07](#).

SVE [15] defines pairs of events for measuring vector and scalar operation workloads.

- FP\_FIXED\_OPS\_SPEC

- FP\_HP\_FIXED\_OPS\_SPEC
- FP\_SP\_FIXED\_OPS\_SPEC
- FP\_DP\_FIXED\_OPS\_SPEC
- INT\_FIXED\_OPS\_SPEC
- FP\_SCALE\_OPS\_SPEC
  - FP\_HP\_SCALE\_OPS\_SPEC
  - FP\_SP\_SCALE\_OPS\_SPEC
  - FP\_DP\_SCALE\_OPS\_SPEC
- INT\_SCALE\_OPS\_SPEC

These events count the number of operations performed:

- SCALE events count variable-length vector operations. Software must multiply these by the number of 128-bit containers in the vector.
- FIXED events count scalar and fixed-length vector operations, for example Advanced SIMD vector operations.
- For vector operations, the counter increments by an amount that is scaled according to the container (element) size. For example, a single-precision operation on a vector counts twice as many operations as a double-precision operation, because the vector contains twice as many elements.
- Compound operations, for example multiply-accumulate and dot products, count as multiple operations.

Although these events are defined by [15], the FIXED events can be implemented on a PE that does not include [15].

These events allow classification by:

- Integer and floating-point.
- For floating-point only, by data width.

R<sub>PMU\_EV\_07</sub>

The events listed in the following table must be implemented for measuring workload:

Number	Mnemonic	Description
0x80C1	FP_FIXED_OPS_SPEC	Non-scalable floating-point element operations speculatively executed
0x80C9	INT_FIXED_OPS_SPEC	Non-scalable integer element operations speculatively executed

If SVE is implemented, the events listed in the following table must be implemented for measuring workload:

Number	Mnemonic	Description
0x80C0	FP_SCALE_OPS_SPEC	Scalable floating-point element operations speculatively executed
0x80C8	INT_SCALE_OPS_SPEC	Scalable integer element operations speculatively executed

### B.3.5 Branch predictor effectiveness

X

This is the rationale for the rule [PMU\\_EV\\_08](#).

The following events are designed for monitoring branch predictor effectiveness.

- BR\_RETIRED
  - BR\_PRED\_RETIRED
    - \* BR\_IMMED\_PRED\_RETIRED
    - \* BR\_IND\_PRED\_RETIRED
      - BR\_INDNR\_PRED\_RETIRED
      - BR\_RETURN\_PRED\_RETIRED
  - BR\_PRED\_RETIRED
    - \* BR\_TAKEN\_PRED\_RETIRED
    - \* BR\_SKIP\_PRED\_RETIRED
  - BR\_MIS\_PRED\_RETIRED
    - \* BR\_IMMED\_MIS\_PRED\_RETIRED
    - \* BR\_IND\_MIS\_PRED\_RETIRED
      - BR\_INDNR\_MIS\_PRED\_RETIRED
      - BR\_RETURN\_MIS\_PRED\_RETIRED
  - BR\_MIS\_PRED\_RETIRED
    - \* BR\_TAKEN\_MIS\_PRED\_RETIRED
    - \* BR\_SKIP\_MIS\_PRED\_RETIRED

R<sub>PMU\_EV\_08</sub>

The events listed in the following table are recommended to be implemented for measuring Section B.3.5:

Number	Mnemonic	Description
0x000C	PC_WRITE_RETIRED	Instruction architecturally executed, condition code check pass, software change of the PC.
0x000D	BR_IMMED_RETIRED	Branch instruction architecturally executed, immediate
0x000E	BR_RETURN_RETIRED	Branch instruction architecturally executed, procedure return, taken.
0x0021	BR_RETIRED	Instruction architecturally executed, branch
0x0022	BR_MIS_PRED_RETIRED	Branch instruction architecturally executed, mispredicted
0x8110	BR_IMMED_PRED_RETIRED	Branch instruction architecturally executed, predicted immediate
0x8111	BR_IMMED_MIS_PRED_RETIRED	Branch instruction architecturally executed, mispredicted immediate
0x8112	BR_IND_PRED_RETIRED	Branch instruction architecturally executed, predicted indirect
0x8113	BR_IND_MIS_PRED_RETIRED	Branch instruction architecturally executed, mispredicted indirect
0x8114	BR_RETURN_PRED_RETIRED	Branch instruction architecturally executed, predicted procedure return
0x8115	BR_RETURN_MIS_PRED_RETIRED	Branch instruction architecturally executed, mispredicted procedure return
0x8116	BR_INDNR_PRED_RETIRED	Branch instruction architecturally executed, predicted indirect excluding procedure return.
0x8117	BR_INDNR_MIS_PRED_RETIRED	Branch instruction architecturally executed, mispredicted indirect excluding procedure return
0x8118	BR_TAKEN_PRED_RETIRED	Branch instruction architecturally executed, predicted branch, taken
0x8119	BR_TAKEN_MIS_PRED_RETIRED	Branch instruction architecturally executed, mispredicted branch, taken

Number	Mnemonic	Description
0x811A	BR_SKIP_PRED_RETIRE	Branch instruction architecturally executed, predicted branch, not taken
0x811B	BR_SKIP_MIS_PRED_RETIRE	Branch instruction architecturally executed, mispredicted branch, not taken
0x811C	BR_PRED_RETIRE	Branch instruction architecturally executed, predicted branch
0x811D	BR_IND_RETIRE	Branch instruction architecturally executed, indirect branch
0x811E	BR_INDNR_RETIRE	Branch instruction architecturally executed, indirect excluding procedure return

R<sub>PMU\_EV\_09</sub>

The BR\_RETIRE event must count unconditional taken branches.

### B.3.6 Latency

X

This is the rationale for the rule [PMU\\_EV\\_10](#).

The base PMU architecture defines *events* that can be counted, but is also useful to measure the *duration*, or latency, of an event. Arm ARM [2] recommends the following are included as IMPLEMENTATION DEFINED events:

Cumulative occupancy for resource queues, like data access queues, and entry or exit counts, so that average latencies can be determined, separating out counts for key resources that might exist.

Summing cumulative occupancy allows *average* latency to be calculated.

$$\text{Average latency} = \frac{\sum \{\text{OCCUPANCY}\}}{\#\{\text{ACCESSES}\}}$$

Armv8.6 does not define a common microarchitectural event for instruction fetches. However, an event for counting instruction fetches, which is analogous to MEM\_ACCESS and its associated latency event are recommended.

R<sub>PMU\_EV\_10</sub>

The events listed in the table below are recommended for PEs in the base server system to measure bandwidth, latency, and utilization:

Number	Mnemonic	Description
0x0019	BUS_ACCESS	Bus access
0x0034	DTLB_WALK	Data TLB access with at least one translation table walk
0x0035	ITLB_WALK	Instruction TLB access with at least one translation table walk
0x0060	BUS_ACCESS_RD	Bus access, read
0x0061	BUS_ACCESS_WR	Bus access, write
0x0066	MEM_ACCESS_RD	Data memory access, read
0x0067	MEM_ACCESS_WR	Data memory access, write
0x8120	INST_FETCH_PERCYC	Event in progress, INST_FETCH
0x8121	MEM_ACCESS_RD_PERCYC	Event in progress, INST_FETCH

Number	Mnemonic	Description
0x8124	INST_FETCH	Instruction memory access
0x8125	BUS_REQ_RD_PERCYC	Event in progress, BUS_REQ_RD
0x8126	BUS_REQ_WR_PERCYC	Event in progress, BUS_REQ_WR
0x8128	DTLB_WALK_PERCYC	Event in progress, DTLB_WALK
0x8129	ITLB_WALK_PERCYC	Event in progress, ITLB_WALK
0x818D	BUS_REQ_RD	Bus request, read
0x818E	BUS_REQ_WR	Bus request, write

### B.3.7 Memory workload

X

The rationale for memory workload events is as follows.

The SVE [15] feature defines a pair of vector and scalar memory workload events.

- LDST\_FIXED\_OPS\_SPEC
  - LD\_FIXED\_OPS\_SPEC
  - ST\_FIXED\_OPS\_SPEC
- LDST\_SCALE\_OPS\_SPEC
  - LD\_SCALE\_OPS\_SPEC
  - ST\_SCALE\_OPS\_SPEC
- LDST\_FIXED\_BYTES\_SPEC
  - LD\_FIXED\_BYTES\_SPEC
  - ST\_FIXED\_BYTES\_SPEC
- LDST\_SCALE\_BYTES\_SPEC
  - LD\_SCALE\_BYTES\_SPEC
  - ST\_SCALE\_BYTES\_SPEC

These events can count either the number of operations or the number of bytes transferred:

- SCALE events and FIXED events determine how the counters increment:
  - SCALE events count variable-length vector loads and stores. Software must scale these events by number of 128-bit containers in the vector.
  - FIXED events count scalar and fixed-length vector loads and stores.
- OPS and BYTES determines what the counter counts:
  - For vector loads and stores, OPS events increment by an amount scaled according to the container (element) size.
  - For scalar loads and stores, pairwise load/store OPS events count multiple accesses. All other Load/Store and atomic operations OPS events count a single access.
  - BYTES events count all load/stores, but the counter increments by the number of bytes that are transferred.

Although these events are defined by [15], the FIXED events can be implemented on a PE that does not include SVE.

These events can be used to derive bandwidth performance figures.

I It is recommended that the events listed in the table below are implemented for measuring Section [B.3.7](#):

Number	Mnemonic	Description
0x80CB	LDST_FIXED_OPS_SPEC	Non-scalable load/store element operations speculatively executed
0x80CD	LD_FIXED_OPS_SPEC	Non-scalable load element operations speculatively executed
0x80CF	ST_FIXED_OPS_SPEC	Non-scalable store element operations speculatively executed
0x80DB	LDST_FIXED_BYTES_SPEC	Non-scalable load/store bytes speculatively executed
0x80DD	LD_FIXED_BYTES_SPEC	Non-scalable load bytes speculatively executed
0x80DF	ST_FIXED_BYTES_SPEC	Non-scalable store bytes speculatively executed

I If SVE is implemented, it is recommended that the events listed in the following table are implemented for measuring Section B.3.7:

Number	Mnemonic	Description
0x80CA	LDST_SCALE_OPS_SPEC	Scalable load/store element operations speculatively executed
0x80CC	LD_SCALE_OPS_SPEC	Scalable load element operations speculatively executed
0x80CE	ST_SCALE_OPS_SPEC	Scalable store element operations speculatively executed
0x80DA	LDST_SCALE_BYTES_SPEC	Scalable load/store bytes speculatively executed
0x80DC	LD_SCALE_BYTES_SPEC	Scalable load bytes speculatively executed
0x80DE	ST_SCALE_BYTES_SPEC	Scalable store bytes speculatively executed

I It is recommended that the events in the following table are implemented for refining [workload](#) measurements:

Number	Mnemonic	Description
0x80C3	FP_HP_FIXED_OPS_SPEC	Non-scalable half-precision floating-point element operations speculatively executed
0x80C5	FP_SP_FIXED_OPS_SPEC	Non-scalable single-precision floating-point element operations speculatively executed
0x80C7	FP_DP_FIXED_OPS_SPEC	Non-scalable double-precision floating-point element operations speculatively executed

I If SVE is implemented, it is recommended that the events in the following table are implemented for refining [workload](#) measurements:

Number	Mnemonic	Description
0x80C2	FP_HP_SCALE_OPS_SPEC	Scalable half-precision floating-point element operations speculatively executed
0x80C4	FP_SP_SCALE_OPS_SPEC	Scalable single-precision floating-point element operations speculatively executed

Number	Mnemonic	Description
0x80C6	FP_DP_SCALE_OPS_SPEC	Scalable double-precision floating-point element operations speculatively executed

Note: If half-precision operations are not implemented, then the FP\_HP\*\_OPS\_SPEC events count no operations.

## B.4 System performance monitors

- I For the purposes of these requirements, an interface is considered to support read and write requests:
- A read is a transaction from a requester to a completer where the completer responds with data.
  - A write is a transaction from a requester to a completer where the requester sends data. The completer might acknowledge the request.
- Other types of traffic can be supported, but are categorized as reads or writes based on the traffic characteristics.
- I Examples of types of traffic that might share an interface include:
- Data transfers
  - Cache coherency messages
  - Distributed virtual memory messages
  - Message-signaled interrupts
- I For the purposes of these requirements, for a monitor located on one side of an interface, accesses are either:
- Outbound, that is originating on the same side as the monitor, or
  - Inbound, that is originating on the other side of the interface.
- I For the purposes of these requirements, an interface is considered to be either unidirectional or bidirectional:
- A unidirectional interface supports requests originating on only one side of the interface. That is, either inbound traffic or outbound traffic, but not both.
  - A bidirectional interface supports requests originating on either side of the interface. That is, both inbound and outbound traffic.
- I A bidirectional read/write interface therefore supports all of inbound reads, inbound writes, outbound reads, and outbound writes.
- The interface between two nodes in a NUMA system is an example of a bidirectional read/write interface.
- I Bandwidth is measured in bytes per second. This means that the monitor must be capable of measuring bytes transferred. However, this might be quantized in larger units.
- If an interface always transfers larger units, a monitor might count transfers in these units. Either the monitor itself or software using the monitor scales the counted value to bytes transferred.
- Furthermore, where an interface usually transfers larger units, a monitor might count transfers in these units, including rounding up smaller transfers to a multiple of the unit. For example, if most traffic to a conventional memory interface is in quantities of 64 bytes, the monitor might present a value that, when scaled, gives an upper-bound for the number of bytes transferred.
- However, if the interface often transfers different sized amounts, the monitor should present a reasonably accurate count.

- I Total bandwidth is the sum of read bandwidth and write bandwidth. However, monitors should allow monitoring of this type of aggregate events using a minimum set of monitored events. Monitors should not require software to monitor many small, detailed events and compute aggregates.
- I Average latency is calculated by accumulating the number of open transactions on each cycle and dividing this by the total number of transactions.
- I A single interface might support multiple devices.  
PCIe is an example of an interface that supports multiple devices on a single interface.
- I The measurements that this section describes should be available to at least the host operating system or hypervisor to allow profiling and monitoring of the base server system. Although this section describes these measurements in terms of monitors that perform the measurements, this does not imply any specific approach to generation of the measurements.
- I These requirements do not specify how the measurements are generated and collected.
- I It is recommended that measurements that are generated by event monitors are implemented in a way that is compatible with the *Arm CoreSight Performance Monitoring Unit Architecture* [16].

The base server system must implement bandwidth monitors for:

- R<sub>PMU\_BM\_1</sub> • Memory interface.
- R<sub>PMU\_BM\_2</sub> • PCIe interface.
- R<sub>PMU\_BM\_3</sub> • External accelerator interface.
- R<sub>PMU\_BM\_4</sub> • Chip-to-chip interface.

R<sub>PMU\_MEM\_1</sub> The base server system must implement average latency monitors for each memory interface.

- R<sub>PMU\_SYS\_1</sub> Each monitor for an interface must be capable of measuring all of the following measurements simultaneously:
- If the interface supports outbound read traffic, average outbound read latency.
  - If the interface supports outbound traffic, total outbound bandwidth.
  - If the interface supports inbound read traffic, average inbound read latency.
  - If the interface supports inbound traffic, total inbound bandwidth.

That is, up to 4 measurements (typically requiring up to 6 event counters) simultaneously.

- R<sub>PMU\_SYS\_2</sub> Each monitor for an interface must be capable of measuring all of the following measurements simultaneously:
- If the interface supports outbound read traffic, outbound read bandwidth and average outbound read latency.
  - If the interface supports outbound write traffic, outbound write bandwidth.
  - If the interface supports inbound read traffic, inbound read bandwidth and average inbound read latency.
  - If the interface supports inbound write traffic, inbound write bandwidth.

That is, up to 6 measurements (typically requiring up to 8 event counters) simultaneously.

- I Each monitor that monitors devices on an interface might limit the number of devices that can be monitored simultaneously.

- R<sub>PMU\_SYS\_5</sub> For NUMA systems, each monitor must be capable of collecting measurements for each of the following traffic groups simultaneously:
- Local node traffic and remote node traffic.
  - All traffic.

- R<sub>PMU\_SYS\_6</sub> For interfaces that carry multiple types of traffic, each monitor must be capable of filtering monitored traffic that is based on its traffic type.

I If FEAT\_MPAM is implemented, it is permitted but not required for the bandwidth monitors to be implemented as MPAM bandwidth monitors and support MPAM filtering. Software should be able to monitor each of the measurements for all traffic simultaneously with the requirements to monitor filtered traffic as part of bandwidth partitioning.

R<sub>PMU\_SYS\_7</sub> Each significant cache in the base server system must be capable of measuring cache effectiveness.

I For example, each significant cache in the system implements hardware performance monitors capable of counting the events listed in the following table:

Mnemonic	Description
CACHE	Cache access.
CACHE_MISS	Cache miss.

---

**Note**

Significant cache means any large system cache, for example the last level cache, that might have a significant impact on performance.

These event monitors are in addition to those provided by a PE.

---

### B.4.1 Recommendations

It is recommended that the base server system implement the monitors that are specified in this section.

I It is recommended that the base server system implement average latency monitors for each:

- PCIe interface
- External accelerator interface
- Chip-to-chip interface

I It is recommended that the base server system implement bandwidth monitors for each:

- External PCIe device
- High-bandwidth internal device

I It is recommended that the base server system implement average latency monitors for each:

- External PCIe device
- High-bandwidth internal device

---

**Note**

High-bandwidth internal device means any device integrated in the SoC device that is likely to use a significant amount of available bandwidth. Examples might include an integrated GPU, NPU, or other data processing engine. These devices might or might not be PCIe devices.

---

#### B.4.1.1 Security

R<sub>PMU\_SEC\_1</sub> When deployed in production systems, performance monitors must not expose Secure data to untrusted software.

- I The definitions of Secure data and untrusted software are IMPLEMENTATION DEFINED and relate to how the information encoded in the data relates to the threat model for the system.  
  
For example, in a typical system that supports both Secure and Non-secure memory, data that is stored in or related to Secure memory is considered Secure data. Other data is considered Non-secure data.
- I The above statement requires that the monitor separately measures Secure and Non-secure events, and either ignores Secure events, or does not expose Secure measurements to untrusted software. Alternatively, untrusted software might need to use a firmware or other proxy to access the performance measurements.

## C Server RAS

The Arm RAS Extension and Arm RAS System Architecture as described in the Arm RAS for Armv8-A supplement[10] describe frameworks for RAS features in Arm PE and Arm-based SoC implementations. Neither standard requires any level of RAS features. Both standards are defined to allow scaling from systems where reliability is not a key concern to systems where reliability is a key concern.

### C.1 Justifications and impact

This section is informative. The requirements are presented in Section C.2.

For Arm server systems, setting a higher standard of features is recommended, both to ensure greater consistency across Arm implementations and to guide designers of Arm server systems. Arm engages in conversation with cloud vendors on RAS, and the requirements presented in this section of the specification are the result of these discussions.

#### C.1.1 PE architecture

##### C.1.1.1 PE error exception handling

The RAS Extension allows error exceptions to be handled in a variety of ways, ranging from:

- Synchronous data abort exceptions. These are precise exceptions that are taken before the corrupt data is consumed by the PE.
- Uncontainable asynchronous SError interrupt exceptions. These might be imprecise exceptions, taken after the corrupt data is consumed and possibly propagated by the PE.

Designers trade off multiple factors when determining how to take error exceptions. If errors occur very late in the execution of an instruction, a power, performance, and area cost might be associated with retaining the state that is required to generate a precise exception.

Nevertheless, having a variety of possible error exceptions creates challenges for portable software. For reliability and availability, the CPU should take exceptions in manner that allows software to contain and possibly recover from an error.

#### C.1.2 System architecture components

This section includes components of a CPU that the RAS System Architecture section in Arm RAS for Armv8-A supplement[10] treats as system components, for example caches.

##### C.1.2.1 Error counters

Some components might generate large amounts of corrected errors if, for example, these components manage a lot of storage. Error counters provide additional information for fault analysis when polling.

Counters should support some form of thresholding. This means that the counter can be configured to generate an event for fault analysis software only after a configured number of corrected errors have been recorded.

Corrected error counters are described in the RAS System Architecture section in the Arm RAS for Armv8-A supplement[10]. These support overflow detection, but not thresholding.

### **C.1.2.2 Interrupt enable controls**

Architectural means to control the interrupts from error nodes are required. In particular, control is required to disable generation of interrupts on error correction, or to switch between polled and event-driven methods.

### **C.1.2.3 Standardized FRU and error location identification**

The cost of going to firmware to decode FRU and location information can be high. This includes using an SMC or interpreter, for example an ASL interpreter to do this on the error recovery path.

Some users do not want firmware to be able to mask errors from the kernel. Other system vendors have the requirement for firmware to mask the errors. Therefore, this requirement must be controllable.

An ASL or similar interpreter might not be available.

---

#### **Note**

An error location might be within a FRU, or might be within a Field Non-replaceable Unit (FNRU), for example within the SoC. A recorded error location is used for fault analysis.

---

### **C.1.2.4 Memory and cache ECC features**

Memory, system caches, and large caches can have a long window of vulnerability for the data they store. Window of vulnerability is the period between accesses to the data, during which the data is liable to corruption. For a cache, a lower access rate for a location implies a larger window of vulnerability for that data.

Using Error Detection and Correction Code (EDAC) features is therefore important for the reliability of such memories.

If the data that is stored in a cache is always clean, then an error can be corrected by invalidating the copy in the cache. Otherwise, an error has to be correctable in-place.

Scrubbing interrupts the window of vulnerability of data. This reduces the chance that single bit errors become multi-bit errors.

Poison is a technique that defers errors to the point of their consumption. Poison provides the following advantages to a system:

- Poison provides a technique to contain latent errors. For example, when dirty data is evicted from a cache, the active process or virtual machine evicting the line is not consuming that data and might not even own the data. If an error is detected when dirty data is evicted from a cache, poison allows the error to be deferred to the real consumer of the data.
- Poison helps prevent false errors from generating failure. For example, a false error is generated if a failure occurs when an error is detected on a memory location that is not accessed, for example a location in a cache line that is not the location being accessed. If the poison granule is smaller than a cache line, poison can be allocated into a cache without causing software that does not access the error to fail.

If the location is overwritten before it is accessed, no failure occurs. However, the error is recorded for fault analysis.

Poison therefore improves availability.

However, the details of EDAC features are a quality of implementation issue. This means that the details vary greatly depending on target customer markets and other design details. There is no one-size-fits-all solution to RAS features.

Therefore, SBSA has a mix of requirements and recommendations for base EDAC. Implementations are expected to go beyond these recommendations for certain markets.

### C.1.3 Software faults

Historically it has been common to return an in-band error response in response to a software fault. (An in-band error response is the RAS term for an external abort.)

Examples of software faults include:

- Access to memory or device register that is not present. This includes cases where Secure and Non-secure memory are physically aliased.
- Access to a device that is not permitted at the device. For example, a Non-secure access to a Secure register.
- Access to a device that is in an inaccessible state or other illegal access. For example, the device is powered down, or the value written is not supported.

Some hardware faults are similar to software faults. For example, when a device surprise removal occurs, software is not at fault for attempting to access the now not-present device.

The Arm Architecture has no such classification for External Abort exceptions that result from in-band error responses. This means that the historical approach leads to software faults being handled by the same software that processes RAS faults. In addition, there is no architecture to record syndrome for a software fault. This means that a conservative error handler might need to treat any software fault as a serious, potentially critical, error.

However, software faults are entirely predictable, usually containable, and avoidable. They should be treated as bugs.

Furthermore, if a device can be relied upon to return an in-band error response, that device must never be made available to any untrusted software, including a guest VM, if it could provide a denial-of-service attack vector.

The RAS System Architecture section in Arm RAS for Armv8-A supplement[10] issue C.b made the recommended that errors classifiable as software faults generate software fault interrupts and do not return in-band error responses. However, this statement is replaced in issue D.a with the following recommendations:

- Where another standard defines a rule or sets a convention for a device, that should be followed. For example:
  - For a PCIe device, certain illegal accesses are RAO/WI.
  - [2] requires that *reserved accesses* to a component, such as reads and writes of unallocated or unimplemented registers and writes to read-only registers, behave as RAZ/WI.
  - [2] requires that under certain conditions accesses to certain debug registers return an error response.
- Accesses to a memory location that is not present can return an in-band error response when all of the following are true:
  - The location is not present due to a configuration of the physical address map that is either static or controlled by trusted software. For example, a configuration choice made by the designer, set during initial system configuration, or reconfigured by trusted software.
  - Within the aligned page that contains the not-present location, all other locations are also *not present* and have the same behavior. The size of this page is the largest supported translation granule size of all PEs in the system.

That is, there is never any legitimate reason for software to access the page containing the location, and trusted software should set up the translation tables to prevent accesses from occurring.

For all other cases, the access should do one of the following:

- Return zeros for a read and ignore writes. This is the recommended behavior for reads and writes of unallocated or unimplemented registers, reads of write-only registers, and writes of read-only registers.

- Return all-ones for a read and ignore writes.
- Return an IMPLEMENTATION DEFINED value for a read and ignore writes.

However, a device might implement a RAS System Architecture error node and error records for recording software faults, for improved debuggability of the fault.

When a device implements a RAS System Architecture error node and error records for recording software faults, software faults can be recorded as an error. Faults are reported with an in-band error response or a fault handling interrupt (software fault interrupt).

This should be configurable through ERR<n>CTLR, allowing software to disable the feature. For example, if an error exception might cause an unrecoverable software state.

When the feature is disabled, accesses should behave as recommended in the “all other cases” description above.

The following ERR<n>STATUS.SERR values can be used to record software faults:

SERR	Description
13	Illegal address (software fault). For example, access to unpopulated memory.
14	Illegal access (software fault). For example, byte write to word register.
15	Illegal state (software fault). For example, device not ready.
25	Error recorded by PCIe error logs. Indicates that the node has recorded an error in a PCIe error log. This might be the PCIe device status register, AER, DVSEC, or other mechanisms defined by PCIe.

If a device does not support a means to record the software fault, it should not return an in-band error response.

## C.2 Server RAS architecture requirements

The following must be true of all PEs and other system components that implement the Arm RAS System Architecture [10] in the base server system.

R<sub>RAS\_01</sub>

If the component is a memory controller or significant cache, it must implement one or more error counters following the standard programming model described in the Arm RAS System Architecture version 1.1 for at least all Corrected errors.

Within the standard programming model:

- The size of the error counter is an IMPLEMENTATION DEFINED choice of 8 or 16 bits (including an overflow bit).
- It is IMPLEMENTATION DEFINED whether a single counter, or a counter pair comprising a base counter and a repeat counter is implemented. The choice of counter format should be based on the particular design parameters of the component.

I

Note that corrected errors in this context does not necessarily mean single bit error. From the memory controller perspective, multi-bit errors (depending on the Quality of Implementation (QoI)) could still be considered “corrected”. In addition, memory controller corrected error counters (CEC) may only count errors that are detected and corrected at the controller.

For example, on systems that use DDR5 or any self-correcting memory technology, software needs to comprehend that the memory controller CEC may not provide the comprehensive count of all corrected DRAM errors.

---

**Note**

Future revisions of this specification may define additional RAS rules and recommendations for systems with DDR5 or other self-correcting memory technology.

---

- $R_{RAS\_02}$  If the component is a memory controller or significant cache, it must follow the RAS System Architecture v1.1 as described in Arm RAS for Armv8-A supplement[10]. The component must implement the CFI, DUI, and UI controls described in the RAS System Architecture v1.1. That is, for the RAS node that records errors for the component,  $ERRFR.\{CFI, DUI, UI\}$  must each be either 0b10 or 0b11.
- $R_{RAS\_03}$  Each error record group implements a single fault handling interrupt for all the records that are contained in the group.
- $R_{RAS\_04}$  If any error record in an error record group is capable of generating an error recovery interrupt, the group implements a single error recovery interrupt for all the records contained in the group.
- $R_{RAS\_05}$  If any error record in an error record group is capable of generating a critical error interrupt, the group implements a single critical error interrupt for all the records contained in the group.
- $R_{RAS\_06}$  Each fault handling, error recovery, or critical error interrupt that is generated by a RAS node that is accessible to a PE in the base server system must be connected to the base system GIC interrupt controller. Unless otherwise required to be implemented as a PPI, the interrupt is an SPI.
- 

**Note**

This rule does not prevent the interrupt also being connected to, for example, a system control processor.

---

- $R_{RAS\_07}$  If the component records a physical address (as required by the RAS System Architecture) for a fault at a location, and the address located in system memory is accessible to PEs in the base system, then the  $ERR\langle n \rangle ADDR.AI$  bit, defined in Arm RAS for Armv8-A supplement[10], must be 0b0 if the address is the same as System Physical Address for the location, and 0b1 otherwise.
- I The System Physical Address (SPA) for a location is the physical address for the location that is used by PEs in the base server system.
- I If the component records errors relating to more than one FRU, it should be possible for an operating system kernel fault handler to identify the FRU based on the information in the error record without requiring active firmware assistance.
- I It should be possible for an operating system kernel fault handler to identify a fault location within a component based on the information in the error record without requiring active firmware assistance, for fault profiling and Predictive failure analysis. This includes locations in non-serviceable components.
- $R_{RAS\_08}$  If an error record of the component is accessible through an *error record group*, then the  $ERRGSR$  (Error Group Status Register) as described in [10] reports the status of the error record.

### C.2.1 Software faults

Examples of software faults include:

- Access to memory or device register that is not present. This includes cases where Secure and Non-secure memory are physically aliased.
- Access to a device that is not permitted at the device. For example, a Non-secure access to a Secure register.

- Access to a device that is in an inaccessible state or other illegal access. For example, the device is powered down, or the value written is not supported.

R<sub>RAS\_10</sub>

Where the PCIe standard [1], Arm architecture[2], or other standard defines a rule or sets a convention for a software fault at a device, that rule or convention must be followed.

I

For example:

- Arm ARM [2] requires that reserved accesses to a component, for example reads and writes of unallocated registers and writes to read-only registers, behave as RAZ/WI.
- [2] requires that under certain conditions accesses to certain debug registers return an error response.
- For a PCIe device, certain illegal accesses are RAO/WI.
- PCIe error rules are described in Section 1.7.7.3.

R<sub>RAS\_11</sub>

On a software fault error if none of the following apply, a read returns an IMPLEMENTATION DEFINED value, and a write is ignored:

- The access is to a not present location.
- The response to the access is defined by RAS\_10.

I

It is recommended that the IMPLEMENTATION DEFINED value that is returned on a read of an unallocated, unimplemented, or write-only register is zero.

R<sub>RAS\_12</sub>

For the purposes of the rule RAS\_11, a location is defined as not present, only if all of the following apply:

- The location is not present due to a configuration of the physical address map that is either static or is controlled by trusted software.
  - A static configuration is a configuration that is made by the system designer, system integrator, or set during initial system configuration.
  - Controlled by Trusted software means that the location might be present or not present, but this is configured by Trusted software.
  - The split between trusted and untrusted is IMPLEMENTATION DEFINED. However, Untrusted would typically include unprivileged software and, in systems that supports virtualization, guest operating systems.
  - Untrusted might or might not include Non-secure hypervisors.
- Within the aligned page that contains the not-present location, all other locations are also not present and have the same behavior. The size of this page is the largest supported translation granule size of all PEs in the system.

I

A device may also include a RAS error node to record a software fault error, to improve debugging of software faults. This node might also include controls to enable the return of an in-band error response, and a software fault interrupt, both of which should be disabled by default.

## C.2.2 Recommended RAS features

This section is informative. This section describes RAS features which are recommended for all PEs and other system components in the base server system.

I

It is recommended that each memory controller, system cache, and other large cache implements error detection.

---

### Note

A large cache is one in which data might reasonably be expected to have a high window of vulnerability. It is recommended that Level 1 PE cache is classified as a large cache.

---

- I It is recommended that caches that hold dirty data implement error correction with at least SECDED. It is also recommended that memory controllers protect external DRAM with higher level of ECC protection than that of the DRAM memory technology attached to it. For example: If the memory technology supports single error correction (SEC) on the DRAM memory, the memory controller is recommended to support at least one stronger ECC protection scheme on the data (for example, SECDED or Symbol ECC ).
- I It is recommended that each system cache and other large cache that can hold dirty data, and implements error detection, supports patrol scrubbing.
- I It is recommended that each system cache and other large cache that holds only clean data and implements error detection does so to at least a SED standard. Detected errors must result in invalidation of the data.

## D Self-hosted Debug for Armv9-A

### D.1 Goals

This chapter specifies hardware requirements for Armv9-A architecture revision and higher [17], where debug functionality running on an Operating System can rely on the hardware resources. It addresses PE features and key aspects of system architecture.

The primary goal is to ensure sufficient system architecture to enable a suitably-built driver and debug software framework to run on all hardware compliant with this section. It is anticipated that a machine-readable description of the hardware configuration is needed to ensure that the driver and debug software are appropriately configured for the specific system.

### D.2 Levels of functionality

The requirements are introduced through levels of functionality, where each level provides a specified set of capabilities that software can rely on.

An implementation is consistent with a level of the Architecture if it implements all of the functionality of a given level, at performance levels appropriate for that particular level. This means that all of the functionality of a level can be exploited by software without unexpectedly poor performance.

While the levels are numbered, the numbering scheme does not always mean that software written for a particular level will work on hardware designed for any lower-numbered or higher-numbered level. For example, some higher-numbered levels restrict the options provided at lower-numbered levels. As such, software written only for the higher level might not work with hardware compliant with a lower-numbered level.

### D.3 Self-hosted debug capabilities

A suitably-built debug software framework running on hardware compliant with one or more levels of this specification should be able to provide debug functionality consistent with the levels implemented. These capabilities include, but are not limited to:

- Tracing of PE program flow, providing detailed history of program execution. PE trace has multiple uses, including:
  - Post-mortem analysis.
  - Reverse debugging.
  - Performance analysis.
- Performance monitoring, providing detailed information about the performance and timing characteristics of programs running on a PE.

### D.4 External debug capabilities

This chapter does not specify any capabilities for external debug. It is anticipated that external debug scenarios might use the same hardware functions as self-hosted debug scenarios however, and both external debug software and self-hosted debug software must accommodate the possibility that hardware functions might not be available because they are being used by another agent. To provide external debug functions, a system might include functionality or components in addition to those specified which might require initialization or programming to provide the self-hosted debug capabilities.

These requirements do not specify any mechanism for software agents to arbitrate over the use of hardware functions.

## D.5 PE Trace

### D.5.1 Introduction

I PE trace provides a detailed history of program flow, and is useful for both debugging and performance analysis. The objective of PE trace is to provide a debug software framework such as gdb or Linux perf, with a history of executed instructions, branches, or function calls.

I This specification details the following sets of requirements:

- The trace information in the generated trace and the observation capabilities of the trace functionality in each PE.
- The methods of capturing the generated trace.

### D.5.2 Embedded Trace Extension (ETE)

I The levels of Embedded Trace Extension (ETE) provide the ability to generate a program trace and to set trace filtering trigger conditions.

### D.5.3 ETE Level 1

I ETE Level 1 provides basic program flow tracing capability for all PEs in the system, with tracing provided from all PEs concurrently.

R In a system with multiple PEs, each PE that is to be used in the same operating system or hypervisor must be compliant with at least the same ETE level.

R Each PE in the system must be provided with a trace unit compliant with the Embedded Trace Extension (FEAT\_ETE).[17]

R The trace unit must support the following ETE features:

- Cycle counting with a cycle counter that is at least 12-bits in size.
- At least one address comparator pair.
- At least one Context ID comparator.
- At least one Virtual context identifier comparator, if the PE implements EL2.
- At least one single-shot comparator control.
- At least one event in the trace.
- At least two counters.
- The sequencer state machine.
- At least four resource selection pairs.

I If a system provides control over the power to the trace unit, it is recommended that this control provides the means to conserve power when the trace unit is not used, including when transitioning between states where trace is used and not used.

R All trace units must share the same physical timestamp source.

R When TRFCR\_EL1 and TRFCR\_EL2 are used to select the time source, the selection of CoreSight time and Physical time must select the same time source. This means that there is no difference between CoreSight time and Physical time.

R The trace units for all PEs must be able to be enabled concurrently, and generate trace concurrently.

- R Each PE in the system must implement the Trace Buffer Extension (FEAT\_TRBE).[17]
- R All the TRBE trace buffers must implement Flag Updates, or all TRBE trace buffers must not implement Flag Updates. TRBIDR\_EL1.F must be the same for all TRBE trace buffers.
- R All the TRBE trace buffers must implement the same minimum alignment constraints. TRBIDR\_EL1.Align must be the same for all TRBE trace buffers.
- R An implementation must reserve a PPI for the TRBE interrupt to PE.
- I It is recommended that the TRBE and system interconnect infrastructure for capturing trace has sufficient bandwidth for common scenarios, while avoiding trace unit overflow scenarios. Common scenarios for tracing include:
- Continuous tracing of one PE, without cycle counts or branch broadcasting, with minimal impact on PE or system performance.
  - Continuous tracing of all PEs concurrently, without cycle counts or branch broadcasting.
  - Continuous tracing of one PE, with cycle counting and branch broadcasting enabled.

## D.6 System Trace Macrocell

### D.6.1 Introduction

- I Arm defines a System Trace Macrocell (STM) Architecture [18] which enables tracing of instrumented software and system activity. The STM is used to generate trace for use either by external debuggers or self-hosted debuggers.
- I Providing an STM implementation in a system provides a memory-mapped programmers model for software instrumentation libraries to target consistently across SoCs.

### D.6.2 Level 1

The Trace Generation requirements in Section D.6.2.1 define a central STM unit for use by instrumentation software running on all PEs in a system.

The Trace Capture requirements in Section D.6.2.2 define self-hosted capture of the trace from an STM. This capture method ensures the PE trace and STM trace are captured independently, ensuring trace buffers for PE tracing can be independently managed.

The Trace Generation requirements defined in Section D.6.2.1 are identical to STG Level 1 defined in [9]. The Trace Capture requirements defined in Section D.6.2.2 are identical to STC Level 2 defined in [9].

#### D.6.2.1 Trace Generation (STG)

- R The system must provide an STM unit compliant with the STMv1.1 architecture as defined in [18].
- R The STM unit must implement the following options:
- STPv2 trace protocol. See [18].
  - Absolute timestamping.
  - At least one stimulus port master.
  - Each master must implement at least 16384 Extended stimulus ports.
  - A 64-bit fundamental data size.
  - Invariant timing and guaranteed transaction types.
  - The following Configuration Register options:
    - STMTSFREQR is read-write.
    - STMTSSTIMR is implemented.

- STMSYNCR is implemented.
  - STMSPER is implemented.
  - STMSPTER is implemented.
  - Trigger control implements both multi-shot and single-shot.
  - STMSPOVERRIDER is implemented.
  - STMSPMOVERRIDER is implemented.
  - STMSPSCR is implemented.
  - STMSPMSCR is implemented.
  - STMTCSR.SYNCEN is always 0b1.
- Data compression on stimulus ports is either programmable or not implemented.

R If any PE implements EL3, the STM unit must implement at least one stimulus port master which is only accessible by Secure software.

R For systems with more than one PE, the system must use stimulus port masters according to one of the following:

- All PEs must be able to use the same stimulus port master.
- Each PE must be able to use a separate stimulus port master.

#### **D.6.2.2 Trace Capture (STC)**

I It is recommended that the system provides an independent trace buffer for the STM, based on shared system memory.

R A trace buffer must be provided in the system to capture trace from the STM.

R The STM must have a separate logical trace buffer based on shared system memory, based on one of the following:

- Shared system memory, based on a configuration of the CoreSight Embedded Trace Router (ETR), as defined in [19].
- Shared system memory, based on a configuration of the Embedded Trace Router (ETR) Architecture, as defined in [20].

R The trace buffer must be able to be located anywhere in normal memory accessible by the PE.

R The trace buffer must support the Circular Buffer mode of operation.

R The trace buffer for the STM must not allow trace from any PE trace unit to be captured.

R Where the trace buffer is based on shared system memory and supports a size greater than the smallest translation granule of the PEs, a page scattering mechanism must be provided to support the trace buffer being partitioned into separate physical pages each of which is no larger than the smallest translation granule of the system.

R Where a page scattering mechanism is required, this must be based on one of the following:

- A System Memory Management Unit (SMMU) based on at least SMMUv1.
- The translation service provided by a CoreSight Address Translation Unit [21].

R If the PEs in the system implement EL2, the page scattering mechanism must support two stages of translation using one of the following methods:

- Both stages are implemented using the CoreSight Address Translation Unit [21]. The CoreSight Address Translation Unit can provide both stages in a single step.
- Both stages are implemented using an SMMU.
- Stage 1 is implemented using the CoreSight Address Translation Unit, and stage 2 is implemented using an SMMU.

R Where any part of the page scattering mechanism is implemented using an SMMU, the SMMU infrastructure must support independent streams for each trace buffer.

- R Where the trace buffer is based on shared system memory, the trace buffer must also support storage of trace into a contiguous physically-addressed buffer without a page scattering mechanism.
- R Where the trace buffer is based on shared system memory and the PE supports both Secure and Non-secure memory, the trace buffer must be able to be located in Non-secure memory, and it is IMPLEMENTATION DEFINED whether the trace buffer is able to be located in Secure memory.
- R Where a trace buffer can capture trace from multiple trace sources, the trace buffer must support packing all of the trace streams in the trace buffer using one of the following:
- The CoreSight formatting protocol, see [22].
- R Any trace buffer control component, page scattering mechanism, and any components that need to be programmed to ensure trace can reach the shared system memory must be accessible without the need for any system specific software other than that required to provide power to the component.
- R Any trace buffer control component, page scattering mechanism, and any components that need to be programmed to ensure trace can reach the trace buffer must be accessible directly in the Physical Address space of every PE.
- I If a system provides control over the power to the trace buffer control components, page scattering mechanism, or any components that need to be programmed to ensure trace can reach the trace buffer memory, it is recommended that this control provides the means to conserve power when trace is not used, including when transitioning between states where trace is used and not used.
- R For all components that are involved in ensuring trace reaches the trace buffer, if the component implements the CoreSight authentication interface, system firmware must ensure the authentication interface is set to a condition which ensures trace will reach the trace buffer, before any Hypervisor or Operating System are started. If the PE implements EL3, only Non-secure debug is required to be enabled. If the PE does not implement EL3, debug must be enabled for the implemented security state. For more details on the authentication interface see [23].
- R Only the following programmable components must be involved in ensuring the trace reaches the trace buffer:
- CoreSight ATB Funnel.
  - CoreSight ATB Replicator.
  - CoreSight Trace Memory Controller, configured as an Embedded Trace FIFO (ETF).
- R For an STM with a dedicated output trigger signal, when this trigger signal is asserted, any trace buffer which captures trace from this STM must be able to detect the event has occurred, and must be able to use this event to stop trace capture. One of the following mechanisms must be used:
- Use of the trigger trace source ID value 0x7D as defined in [22], and use of the ATB Trigger Enable function in the STM.
  - Use of a dedicated mechanism when a trigger occurs in the STM.
- R If the system uses the trace source ID value 0x7D, the trace buffer must support one of the following:
- Ignoring the trigger trace source ID other than for the purposes of using the event to control trace capture. The CoreSight ETR does not obey this rule.
  - Packing of the trace streams such that the trigger source ID and its payload are embedded in the data stored in the trace buffer. The CoreSight ETR supports this operation.
- R Each trace buffer must be able to generate an interrupt to indicate when any of the following occur:
- In Circular Buffer mode, when the top of the trace buffer is reached.
  - The trace buffer has reached a pre-programmed watermark fill level.
- R The trace buffer interrupt must be at least one of an SPI or LPI.
- R The trace buffer control component, page scattering mechanism, and any components that need to be programmed to ensure the trace can reach the trace buffer from the STM must not be reset on a Warm reset of the PEs.

- R The path from the STM to its trace buffer must not rely on any programmable components that also control the path for a different trace unit to its trace buffer.
- R The registers to control each trace buffer must be located in separate 64KB pages.
- R Any components which control the path of trace from the STM to the trace buffer and the page scattering mechanism must be located in separate 64KB pages from components which control the path of trace from a different trace unit.
- R Where the CoreSight Address Translation Unit is used for stage 2 translations, each CoreSight Address Translation Unit must be located in a separate 64KB page from all other trace control components.
- I Where the CoreSight Address Translation Unit is used, it is recommended that the CoreSight Address Translation Unit is located in a separate 64KB page from all other trace control components.

## E GPU accelerated compute

This chapter lists a set of rules that must be met for correct and performant GPU accelerated compute. In addition, a set of recommendations are also given which when implemented would improve the performance of GPU accelerated compute.

- R<sub>GPU\_01</sub> PCIe integration requirements as defined in Section 1.7.7.1 must be implemented.
- R<sub>GPU\_02</sub> The system must support peer-to-peer read and write transactions between all Endpoints in the system, including those belonging to different Root Port downstream hierarchies. As described in the PCI Express Base Specification[1], there are two mechanisms by which a Root Complex may support peer-to-peer transactions:
1. The Root Complex “takes ownership” of the Transactions as a Completer and initiates separate Transactions as a Requester to the peer.
  2. The Root Complex routes TLPs directly between peers. Note that in some configurations, Completions may take a different route than Requests and may not pass through the Root Complex.
- R<sub>GPU\_03</sub> The peer-to-peer rules as defined in the Arm BSA [3] document are mandatory:
- PCI\_PP\_02
  - PCI\_PP\_03
  - PCI\_PP\_04
  - PCI\_PP\_05
- Additionally, the information statement PCI\_PP\_06 is a mandatory rule for all switches in the system.
- R<sub>GPU\_04</sub> The Root Ports and System MMUs in the system must support PCIe ATS.
- Root Port support is indicated by the following PCI Express capabilities:
- ATS Extended Capability
  - Page Request Extended Capability
- System MMU support is indicated by the following SMMUV3 capabilities:
- SMMU\_IDR0.ATS
  - SMMU\_IDR0.PRI
- I It is recommended that the system supports atomic transactions, as advertised by the following Root Port capabilities:
- Device Capabilities 2 Register (Offset 24h): 32-bit AtomicOp Completer Supported
  - Device Capabilities 2 Register (Offset 24h): 64-bit AtomicOp Completer Supported
  - Device Capabilities 2 Register (Offset 24h): 128-bit CAS Completer Supported
- I It is recommended that the system support peer-to-peer atomic transactions between all Endpoints in the system, extending GPU\_02.
- This support is achieved in two parts:
1. Ensure each Root Port advertises the following capability: Device Capabilities 2 Register (Offset 24h): AtomicOp Routing Supported.
  2. Ensure that there is no asymmetry of this capability and that AtomicOp Routing is supported between all Root Port pairs in the system. This strengthens the PCI Express specification, which allows for asymmetry.
- I It is recommended that all NUMA nodes in the system have a symmetric PCIe topology with the same number of Root Port instances and symmetric lane width configurations for attaching accelerator and storage devices.

I Unless system address space is exhausted, each GPU device should be able to allocate the following amount of system address space:

- At least 16 MB for use as 32-bit non-prefetchable BAR memory region.
- At least 256 GB for use as 64-bit prefetchable BAR memory region.
- At least 32 MB for use as a second 64-bit prefetchable BAR memory region.

Some GPUs may require a legacy IO port space for legacy features like VGA. See Arm BSA [3] rule PCI\_IO\_01 for legacy IO port space requirement. The typical size is 128B.

I It is recommended that the amount of non-prefetchable memory space and prefetchable memory space allocated to the PCIe hierarchy below each Root Port is independently configurable on a per Root Port basis. The system should allow such configuration to be done post discovery of the downstream hierarchy's memory space needs.

---

**Note**

The intent of this recommendation is that a Root Port with a complex hierarchy (for example, including switches and multiple devices) can have more memory space allocated to its downstream hierarchy when compared to a Root Port with a single device directly attached to it.

---

I If there is no target for an address range in the system address map and that address range is above 4GB and below the maximum physical address limit, then mapping that address range to PCIe for use as 64 bit prefetchable BAR memory should be permitted.

I For performance reasons, it is recommended that:

1. Root Ports allow Inbound writes with RO=1 to overtake previously received Inbound writes provided the Root Port complies with the recommendation (2) below.
2. A write request with RO=1 is not permitted to overtake another write request if both writes are writing to bytes that fall within the same 64-bytes aligned granule in the system address space.

---

**Note**

64 bytes is the minimum granule size, performance of GPU compute usage models can benefit from larger granule sizes of up to 256 bytes.

---

I For performance reasons, it is recommended that a read is not permitted to overtake another read request if both reads are accessing bytes that fall within the same 64-bytes aligned granule in the system address space.

---

**Note**

64 bytes is the minimum granule size, performance of GPU compute usage models can benefit from larger granule sizes of up to 256 bytes.

---

I For performance reasons, it is recommended that the system is capable of supporting full saturation (for example, full utilization of the available link bandwidth) of the Root Port Inbound and Outbound links across a variety of traffic patterns. The following cases are provided for guidance:

1. Device to host memory (Inbound for Root Ports) read only, write only, and concurrent read and write traffic for all supported packet sizes greater than or equal to 128 bytes.

2. Device to device (peer-to-peer) read only, write only and concurrent read and write traffic for all supported packet sizes greater than or equal to 128 bytes.

---

**Note**

While verifying performance, a latency of 2.5us should be assumed between a non-posted request entering requestor side Root Port pins to the last byte of the last completion for that request exiting requestor side Root Port pins.

---