



RISC-V Server SoC Specification

Server SoC Task Group

Version v0.0.3, 2023-09-01: Draft

Table of Contents

Preamble.....	1
Copyright and license information.....	2
Contributors.....	3
1. Introduction.....	4
1.1. Glossary.....	6
2. Server SoC Hardware Requirements.....	9
2.1. RISC-V Harts.....	9
2.2. Clocks and Timers.....	10
2.3. External Interrupt Controllers.....	11
2.4. Input-Output Memory Management Unit (IOMMU).....	12
2.5. PCIe Subsystem Integration.....	15
2.5.1. Enhanced Configuration Access Method (ECAM).....	16
2.5.2. PCIe Memory Space.....	18
2.5.3. Access Control Services (ACS).....	20
2.5.4. Address Routed Transactions.....	21
2.5.5. ID Routed Transactions.....	22
2.5.6. Cacheability and Coherence.....	22
2.5.7. Message signaled interrupts.....	23
2.5.8. Precision Time Measurement (PTM).....	24
2.5.9. Error and Event Reporting.....	24
2.5.10. Vendor Specific Registers.....	25
2.5.11. SoC-Integrated PCIe Devices.....	25
2.6. Reliability, Availability, and Serviceability (RAS).....	26
2.7. Quality of Service.....	29
2.8. Manageability.....	30
2.9. Debug.....	32
2.10. Trace.....	32
2.11. Performance Monitoring.....	33
2.12. Security Requirements.....	34
Bibliography.....	35

Preamble



This document is in the [Development state](#)

Assume everything can change. This draft specification will change before being accepted as standard, so implementations made to this draft specification will likely not conform to the future standard.

Copyright and license information

This specification is licensed under the Creative Commons Attribution 4.0 International License (CC-BY 4.0). The full license text is available at creativecommons.org/licenses/by/4.0/.

Copyright 2022 by RISC-V International.

Contributors

This RISC-V specification has been contributed to directly or indirectly by (in alphabetical order):

Aaron Durbin, Andrei Warkentin, Andrew Jones, Beeman Strong, Cameron McNairy, Greg Favor, Heinrich Schuchardt, Isaac Chute, Jon Masters, Ken Dockser, Krste Asanovic, Manu Gulati, Mark Hayter, Michael Klinglesmith, Paul Walmsley, Shubu Mukherjee, Sibaranjan Pattnayak, Ved Shanbhogue

Chapter 1. Introduction

The RISC-V server SoC specification defines a standardized set of capabilities that portable system software such as operating systems and hypervisors, can rely on being present in a RISC-V server SoC.

A server is a computing system designed to manage and distribute resources, services, and data to other computers or devices on a network. It is often referred to as a 'server' because it serves or provides information and resources upon request. Such computing systems are designed to operate continually and have higher requirements for capabilities such as RAS, security, performance, and quality of service. Examples of servers include web servers, file servers, database servers, mail servers, game servers, and more. This specification focuses on defining requirements for general-purpose server computing systems that may be used for one or more of these purposes.

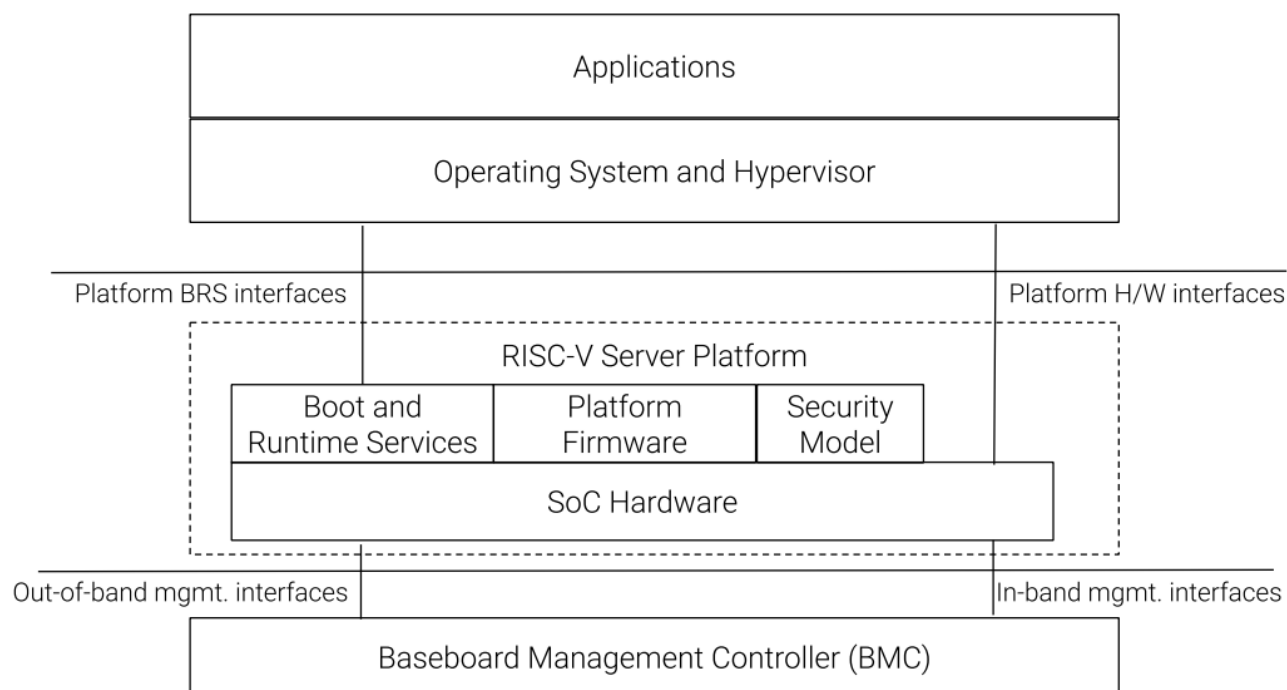


Figure 1. Components of a RISC-V Server Platform

The RISC-V server platform is defined as the collection of SoC hardware, platform firmware, boot/runtime services, and security services. The platform provides hardware interfaces (e.g., harts, timers, interrupt controllers, PCIe root ports, etc.) to portable system software. It also offers a set of standardized boot and runtime services based on the UEFI and ACPI standards. To support provisioning and platform management, it interfaces with a baseboard management controller (BMC) through both in-band and out-of-band (OOB) management interfaces. The in-band management interfaces support the use of standard manageability specifications like MCTP, PLDM, IPMI, and Redfish for provisioning and management of the operating system executing on the platform. The OOB interface supports the use of standard manageability specifications like MCTP, PLDM, Redfish, and IPMI for functions such as power management, telemetry, debug, and provisioning. The security model includes guidelines and requirements for aspects such as debug authorization, secure/measured boot, firmware updates, firmware resilience, and confidential computing, among others.

The platform firmware, typically operating at privilege level M, is considered part of the platform and is usually expected to be customized and tailored to meet the requirements of the SoC hardware (e.g., initialization of address decoders, memory controllers, RAS, etc.).

This specification standardizes the requirements for the hardware interfaces and capabilities (e.g., harts, timers, interrupt controllers, PCIe root complexes, RAS, QoS, in-band management, out-of-band (OOB) management, etc.) provided by the SoC to software executing on the application processor harts at privilege levels below M. It enables OS and hypervisor vendors to support such SoCs with a single binary OS image distribution model. The requirements posed by this specification represent a standard set of infrastructural capabilities, encompassing areas where divergence is typically unnecessary and where novelty is absent across implementations.

To be compliant with this specification, the SoC **MUST** support all mandatory requirements and **MUST** support the listed versions of the specifications. This standard set of capabilities **MAY** be extended by a specific implementation with additional standard or custom capabilities, including compatible later versions of listed standard specifications. Portable system software **MUST** support the specified mandatory capabilities to be compliant with this specification.

The requirements in this specification use the following format:

ID#	Requirement
CAT_NNN	<p>The CAT is a category prefix that logically groups the requirements and is followed by 3 digits - NNN - assigning a numeric ID to the requirement.</p> <p>The requirements use the key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" that are to be interpreted as described in RFC 2119 [1] when, and only when, they appear in all capitals, as shown here. When these words are not capitalized, they have their normal English meanings.</p>
<i>A requirement or a group of requirements may be followed by non-normative text providing context or justification for the requirement. The non-normative text may also be used to reference sources that are the origin of the requirement.</i>	

This specification groups the requirements in the following broad categories:

- Harts
- Clocks and Timers
- External Interrupt Controllers
- IOMMU
- PCIe subsystem
- Reliability, Availability, and Serviceability
- Quality of Service
- Debug
- Trace

- Performance monitoring
- Security

1.1. Glossary

Most terminology has the standard RISC-V meaning. This table captures other terms used in the document. Terms in the document prefixed by 'PCIe' have the meaning defined in the PCI Express (PCIe) Base Specification [2] (even if they are not in this table).

Table 1. Terms and definitions

Term	Definition
ACPI	Advanced Configuration and Power Interface [3].
ACS	Follows PCI Express. Access Control Services. A set of capabilities used to provide controls over routing of PCIe transactions.
AER	Advanced Error Reporting. Follows PCI Express. A PCIe defined error reporting paradigm.
AIA	RISC-V Advanced Interrupt Architecture.
ATS	Follows PCI Express. Address Translation Services.
BAR or Base Address Register	Follows PCI Express. A register that is used by hardware to show the amount of system memory needed by a PCIe function and used by system software to set the base address of the allocated space.
BMC	Baseboard Management Controller. A motherboard resident management controller that provides functions for platform management.
CXL	Compute Express Link bus standard.
DMA	Direct Memory Access.
DMTF	Distributed Management Task Force. Industry association for promoting systems management and interoperability.
ECAM	Follows PCI Express. Enhanced Configuration Access Method. A mechanism to allow addressing of Configuration Registers for PCIe functions. In addition to the PCI Express Base Specification, see the detailed requirements in this document.
EP, EP=1	Follows PCI Express. Also called Data Poisoning. EP is an error flag that accompanies data in some PCIe transactions to indicate the data is known to contain an error. Defined in PCI Express Base Specification 6.0 section 2.7.2. Unless otherwise blocked, the poison associated with the data must continue to propagate in the SoC internal interconnect.
GPA	Guest Physical Address: An address in the virtualized physical memory space of a virtual machine.
Guest	Software in a virtual machine.

Term	Definition
Hierarchy ID or Segment ID	Follows PCI Express. An identifier of a PCIe Hierarchy within which the Requester IDs are unique.
Host Bridge	Part of a SoC that connects host CPUs and memory to PCIe root ports, RCiEP, and non-Pcie devices integrated in the SoC. The host bridge is placed between the device(s) and the platform interconnect to process DMA transactions. IO Devices may perform DMA transactions using IO Virtual Addresses (VA, GVA or GPA). The host bridge invokes the associated IOMMU to translate the IOVA to Supervisor Physical Addresses (SPA).
HPM	Hardware Performance Monitor.
Hypervisor	Software entity that controls virtualization.
ID	Identifier.
IMSI	Incoming Message-signaled Interrupt Controller.
IO Bridge	See host bridge.
IOVA	I/O Virtual Address: Virtual address for DMA by devices.
MCTP	Follows DMTF Standard. Management Component Transport Protocol used for communication between components of a platform management system.
MSI	Message Signaled Interrupts.
NUMA	Non-uniform memory access.
OS	Operating System.
PASID	Follows PCI Express. Process Address Space Identifier: It identifies the address space of a process. The PASID value is provided in the PASID TLP prefix of the request.
PBMT	Page-Based Memory Types.
PRI	Page Request Interface. Follows PCI Express. A PCIe protocol that enables devices to request OS memory manager services to make pages resident.
RCiEP	Root Complex Integrated Endpoint. Follows PCI Express. An internal peripheral that enumerates and behaves as specified in the PCIe standard.
RCEC	Follows PCI Express. Root Complex Event Collector. A block for collecting errors and PME messages in a standard way from various internal peripherals.
RID or Requester ID	Follows PCI Express. An identifier that uniquely identifies the requester within a PCIe Hierarchy. Needs to be extended with a Hierarchy ID to ensure it is unique across the platform.
Root Complex, RC	Follows PCI Express. Part of the SoC that includes the Host Bridge, Root Port, and RCiEP.

Term	Definition
Root Port, RP	Follows PCI Express. A PCIe port in a Root Complex used to map a Hierarchy Domain using a PCI-PCI bridge.
P2P or peer-to-peer	Follows PCI Express. Transfer of data directly from one device to another. If the devices are under different PCIe Root Ports or are internal to the SoC this may involve data movement across the SoC internal interconnect.
PLDM	Follows DMTF standard. Platform Level Data Model.
PMA	Physical Memory Attributes.
PMP	Physical Memory Protection.
Prefetchable Non-prefetchable	Follows PCI Express. Defines the property of the memory space used by a device. For details see the PCIe Base Specification. Broadly, non-prefetchable space covers any locations where reads have side effects or where writes cannot be merged.
SMBIOS	System Management BIOS.
SoC	System on a chip, also referred as system-on-a-chip and system-on-chip.
SPA	Supervisor Physical Address: Physical address used to to access memory and memory-mapped resources.
SPDM	Follows DMTF Standard. Security Protocols and Data Models. A standard for authentication, attestation and key exchange to assist in providing infrastructure security enablement.
SR-IOV	Follows PCI Express. Single-Root I/O Virtualization.
TLP	Follows PCI Express. Transaction Layer Packet. Defined by Chapter 2 of the PCI Express Base Specification.
QoS	Quality of Service. Quality of Service (QoS) is defined as the minimal end-to-end performance that is guaranteed in advance by a service level agreement (SLA) to a workload.
UEFI	Unified Extensible Firmware Interface. [4]
UR, CA	Follows PCI Express. Error returns to an access made to a PCIe hierarchy.
VM	Virtual Machine.

Chapter 2. Server SoC Hardware Requirements

2.1. RISC-V Harts

A RISC-V SoC includes a RISC-V application processor and may include one or more service processors. These service processors may provide services such as security and power management to software executing on the application processors, and they may themselves implement the RISC-V ISA. The requirements in this section apply solely to harts in the application processors of the SoC.

ID#	Requirement
RVA_010	The RISC-V application processor harts in the SoC MUST support the RVA23 ISA profile [5].
<i>The next major release of the profiles is expected to be RVA24, which is still under construction. This specification should be updated to comply with the RVA24 profiles as the profile definition becomes more finalized.</i>	
RVA_020	<p>The RISC-V application processor harts in the SoC MUST support the following extensions:</p> <ul style="list-style-type: none">• Sv48• Sv48x4• Svadu• Sdtrig• Sdext• H• Sscofpmf• Zkr• Ssecorrupt• Ssccfg• Ssctr• Sscrind
<i>Ssccfg, Sscrind, and Ssctr are under construction.</i>	
<i>Many of these mandated extensions are optional in the RVA23 ISA profile. This requirement is placed here as a placeholder. These mandates may be moved downstream into a platform specification (as a chapter) or as a new ISA profile specification.</i>	
RVA_030	The ISA extensions and associated CSR field widths implemented by any of the RISC-V application processor harts in the SoC MUST be identical.

ID#	Requirement
<p><i>The RVA23 profile supports a set of optional extensions. The set of optional extensions implemented by the harts must be identical. Where the extension supports optionality in the form of field widths (e.g., ASIDLEN, VLEN, allowed vstart values, physical address width, debug triggers, cache-block size, etc.), the implementation of these must also be identical. Having an identical ISA on all harts allows system software to migrate tasks among the harts without constraints.</i></p>	
RVA_040	The RISC-V application processor harts in the SoC MAY support different power and performance characteristics but MUST be otherwise indistinguishable from each other from a software execution viewpoint.
<p><i>All harts in the SoC being indistinguishable from a software execution viewpoint allows system software to migrate tasks among the harts without constraints.</i></p>	
RVA_050	<p>The RISC-V application processor hart MUST support:</p> <ul style="list-style-type: none"> • Single stepping using the step bit in <code>dcsr</code> • Debug scratch register 0 (<code>dscratch0</code>)
RVA_060	<p>The RISC-V application processor hart MUST support:</p> <ul style="list-style-type: none"> • At least 6 instruction address match triggers. • At least 6 load/store address match triggers. • At least one icount trigger to support single stepping. • At least one interrupt trigger. • At least one exception trigger. • Trigger filtering using <code>hcontext</code>. • Trigger filtering using all VMID encodings supported by the hart. • Trigger filtering using <code>scontext</code>. • Trigger filtering using all ASID encodings supported by the hart.
RVA_070	The RISC-V application processor MUST support at least 6 hardware performance counters defined by the Zihpm extension in addition to the three counters defined by Zicntr extension.

2.2. Clocks and Timers

ID#	Requirement
CTI_010	The <code>time</code> CSR MUST increment at a constant frequency and the count MUST be in units of 1 ns. The frequency at which the CSR provides an updated time value MUST be at least 100 MHz.
CTI_020	The <code>time</code> counter MUST appear to be always on and MUST appear to not lose its count across power states, including when the hart is powered off.

2.3. External Interrupt Controllers

This section specifies the requirements on the interrupt controllers used to deliver external interrupts to the RISC-V application processor harts.

ID#	Requirement
IIC_010	The RISC-V Advanced Interrupt Architecture [6] MUST be supported.
IIC_020	External interrupts MUST be signaled to a hart as message-signaled interrupts (MSI).
<p><i>Since MSI are a memory write, it allows for a simplified implementation of producer-consumer ordering rules. For example, interrupts following a write from a device must only be observed after the writes have been observed, and the read completion generated by a device following an interrupt must only be observed after the interrupt has been observed.</i></p> <p><i>Message Signaled interrupts (MSI) are the preferred interrupt signaling mechanism in PCIe. Supporting MSI for harts allows low-latency interrupt signaling.</i></p>	
IIC_030	The Incoming Message-signaled Interrupt Controller (IMSIC) MUST implement an interrupt file for S-mode.
IIC_040	The IMSIC MUST support at least 5 VS-mode interrupt files.
<p><i>Supporting 5 VS-mode interrupt files for a hart allows context switching between up to 5 virtual CPUs (vCPU) on a hart without needing to swap the contents of the interrupt file out to memory. This is particularly beneficial when devices are directly assigned to virtual machines (VMs), as swapping out the context of an IMSIC interrupt file may result in longer latencies due to the need to redirect device interrupts to a memory-resident interrupt file.</i></p>	
IIC_050	The S-mode interrupt file MUST support at least 255 interrupt identities.
IIC_060	The VS-mode interrupt files MUST support at least 63 interrupt identities.
IIC_070	The memory regions for IMSIC interrupt files MUST have the following PMAs: <ul style="list-style-type: none">• Not cacheable, coherent, I/O region.• Support 4-byte aligned reads and writes.
IIC_080	<p>If the SoC implements devices that use wire-signaled interrupts then the SoC MUST implement an APLIC as specified by the RISC-V AIA specification and MUST use the APLIC to convert the wire-signaled interrupts into MSIs.</p> <p>If implemented, the APLIC MUST support:</p> <ul style="list-style-type: none">• Supervisor interrupt domain.• <code>GEILEN</code> equal to that implemented by the harts.• MSI delivery mode.• Extempore MSI generation using the <code>genmsi</code> register.

ID#	Requirement
<i>SoC devices using wire-signaled interrupts must implement the rules related to ordering of interrupts vs. older read/writes from devices as specified by the device and/or bus interface specifications that such devices conform to. See also SID_010.</i>	

2.4. Input-Output Memory Management Unit (IOMMU)

ID#	Requirement
IOM_010	All IOMMUs in the SoC MUST support the RISC-V IOMMU specification [7].
<i>The number of IOMMUs implemented in the SoC is UNSPECIFIED.</i>	
IOM_020	<p>All DMA capable peripherals (RCiEP and non-PCIe devices) and all PCIe root ports that are made available to software on the RISC-V application processor harts MUST be governed by an IOMMU.</p> <p>This requirement does not apply to platform devices such as the APLIC or the IOMMU itself. This requirement does not apply to memory accesses originated by a debug module using a System Bus Access block.</p>
<i>DMA capable peripherals being governed by an IOMMU allows OS/hypervisors to restrict DMA originating from such devices to a subset of memory to enhance security and software fault tolerance. The address translation capability provided by the IOMMU enables usages such as passthrough of such devices to virtual machines, shared virtual addressing, etc.</i>	
IOM_030	The IOMMU governing a PCIe root port MUST support at least 16-bit wide device IDs.
IOM_040	An IOMMU that does not govern a PCIe root port MUST support a device ID width REQUIRED to support all requester IDs originated by the devices governed by that IOMMU.
IOM_050	The IOMMU MUST implement all the page based virtual memory system modes and extensions that are supported by the IOMMU and are also implemented by the RISC-V application processor harts in the SoC.
<i>The page based virtual memory system modes that may be optionally supported by the IOMMU are defined in the IOMMU capabilities register.</i>	
IOM_060	<p>The IOMMU SHOULD support the following virtual memory extensions:</p> <ul style="list-style-type: none"> • Svadu (enumerated by 1 setting of <code>capabilities.AMO_HWAD</code>)
<i>Hardware A/D bit updates capability enables efficient support for use models such as VM migration, shared virtual addressing, and user space work submission.</i>	
IOM_070	The IOMMU SHOULD support pass-through mode and MRIF mode MSI address translation.
IOM_080	When MRIF mode MSI address translation is supported, the IOMMU MUST support atomic updates to the MRIF (enumerated by 1 setting of <code>capabilities.AMO_MRIF</code>).

ID#	Requirement
IOM_090	IOMMU governing PCIe root ports SHOULD support PCIe address translation services (ATS).
<p><i>High performance devices such as DPU/SmartNICs, GPUs, and FPGAs, utilized in server platforms rely on ATS and Page Request services to achieve high throughput and low-latency I/O. Supporting ATS is also required for efficiently accommodating usage models such as Shared Virtual Addressing and direct work submission from user mode.</i></p>	
IOM_100	IOMMU governing PCIe root ports SHOULD support the T2GPA mode of operation with ATS if ATS is supported.
IOM_110	IOMMU governing RCiEP MUST support PCIe address translation services (ATS) if any of the RCiEPs governed by the IOMMU support the ATS capability.
<p><i>The T2GPA control enables a hypervisor to prevent DMA from a device, even if the device misuses the ATS capability and attempts to access memory that is not explicitly authorized by the page tables governing that device's memory accesses. The threat model could also include a man-in-the-middle on the PCIe link inserting ATS-translated requests to access memory that was not previously authorized. As an alternative to setting T2GPA to 1, the hypervisor might establish a trust relationship with the device if authentication protocols such as SPDM are supported by the device. For PCIe, for example, the PCIe Component Measurement and Authentication (CMA) capability provides a mechanism to verify the device's configuration and firmware/executable (Measurement) and hardware identities (Authentication). This mechanism establishes such a trust relationship, and the PCIe link may be integrity-protected using PCIe integrity and data encryption (IDE) to defend against a man-in-the-middle adversary.</i></p>	
IOM_120	IOMMU governing RCiEP MAY support the T2GPA mode of operation with ATS if ATS is supported.
<p><i>The threats associated with misuse of ATS or malicious insertion of ATS translated requests by a man-in-the-middle may not be present with RCiEP being integrated in the SoC.</i></p>	
IOM_130	IOMMU MUST support MSI and MAY support wire-signaled interrupts for external interrupts originated by the IOMMU itself.
IOM_140	IOMMU MUST support little-endian memory access to its in-memory data structures.
IOM_150	IOMMU MAY support big-endian mode memory access to its in-memory data structures.
<p><i>The IOMMU memory-mapped registers always have a little-endian byte order.</i></p>	
IOM_160	IOMMU MAY support the PCIe PASID capability.
IOM_170	IOMMU that supports PASID capability MUST support 20-bit PASID width and MAY support 8-bit and 17-bit PASID widths.
<p><i>PCIe specification strongly recommends that hardware implement the maximum width of 20 bits to ensure interoperability with system software. See also the implementation note on PASID width homogeneity in the PCIe specification 6.0 section 6.20.2.2.</i></p>	
IOM_180	IOMMU SHOULD support a hardware performance monitor (HPM).

ID#	Requirement
<i>The HPM is a valuable tool for system integrators for performance monitoring and optimizations. An IOMMU is highly recommended to provide an HPM.</i>	
IOM_190	An IOMMU that supports an HPM MUST support the cycles counter and MUST incorporate at least 4 event counters.
IOM_200	The cycles counter and the event counters MUST be at least 40 bits wide.
IOM_210	The IOMMU SHOULD support the software debug capabilities enumerated by capabilities.DBG .
IOM_220	The physical address width supported by the IOMMU MUST be greater than or equal to the physical address width supported by the RISC-V application processor harts in the SoC.
<i>Having the physical address width greater than or equal to the width supported the harts in the SoC enables use of all addressable memory for I/O and facilitates the sharing of page tables between the hart MMU and the IOMMU.</i>	
IOM_230	The reset default of the iommu_mode MUST be Off .
<i>The IOMMU disallowing DMA unconditionally following reset due to the mode being Off allows the SoC firmware and software to enable DMA when suitable security protections as required have been established. The IOMMU mode being Off at reset does not pose a significant issue to SoC firmware that needs to employ DMA (e.g., for firmware loading) as that firmware may program the mode in the appropriate IOMMU prior to programming the peripheral governed by that IOMMU to perform a DMA.</i>	
IOM_240	The IOMMU SHOULD be implemented as an RCiEP with base class 08H and subclass 06H [8].
<i>The base class 08H and sub-class 06H are designated by PCIe for use by an IOMMU. Implementing the IOMMU as a PCIe device allows an operating system to determine a driver for the IOMMU and to assign resources such as interrupt vectors to the IOMMU in a PCIe compatible manner.</i>	
IOM_250	The host bridge MUST enforce the physical memory attribute checks and physical memory protection checks on memory accesses originated by the IOMMU and signal detected access violations to the IOMMU.
<i>These checks are analogous to the PMA and PMP checks performed by the RISC-V hart. The host bridge (also known as IO bridge) invokes the IOMMU for address translations. To perform the operations requested by the host bridge the IOMMU may need to access in-memory data structures such as the device directory table and page tables.</i>	
IOM_260	An IOMMU MUST support 24-bit device IDs if the IOMMU governs multiple PCIe root ports that may be part of different PCIe hierarchies.
<i>An IOMMU governing PCIe root ports uses requester ID (RID) - the tuple of bus/device/function numbers (or just bus/function numbers, if the PCIe ARI option is used) - to locate a device context to use for address translation and protection. The 16-bit RID uniquely identifies a requester within a hierarchy. This RID needs to be augmented with the Hierarchy ID (also known as segment ID) - an 8-bit number - to uniquely identify a requester across PCIe hierarchies.</i>	

ID#	Requirement
IOM_270	The host bridge MUST provide the PCIe RID as the bits 15:0 of the device_id to the IOMMU for requests from PCIe EPs and RCiEP.
IOM_280	When the IOMMU supports 24-bit device IDs, the host bridge MUST specify the segment number associated with the PCIe hierarchy from which requests were received as the bits 23:16 of the device_id to the IOMMU.
IOM_290	The determination of device_id input to an IOMMU for requests originating from non-PCIe devices is UNSPECIFIED . If PCIe and non-PCIe endpoints/RCiEP are governed by the same IOMMU, the SoC MUST ensure that there is no overlap between any device_id associated with non-PCIe devices with any device_id formed using the PCIe RID (and if applicable the segment ID).

2.5. PCIe Subsystem Integration

A PCIe subsystem consists of a root complex with a collection of root ports, root complex event collectors (RCECs), root complex register blocks (RCRBs), and root complex integrated end points (RCiEPs). The root complex implements a host bridge to connect the PCIe root ports, RCECs, RCRBs, and RCiEP, to the CPU and system memory in the SoC through an interconnect.

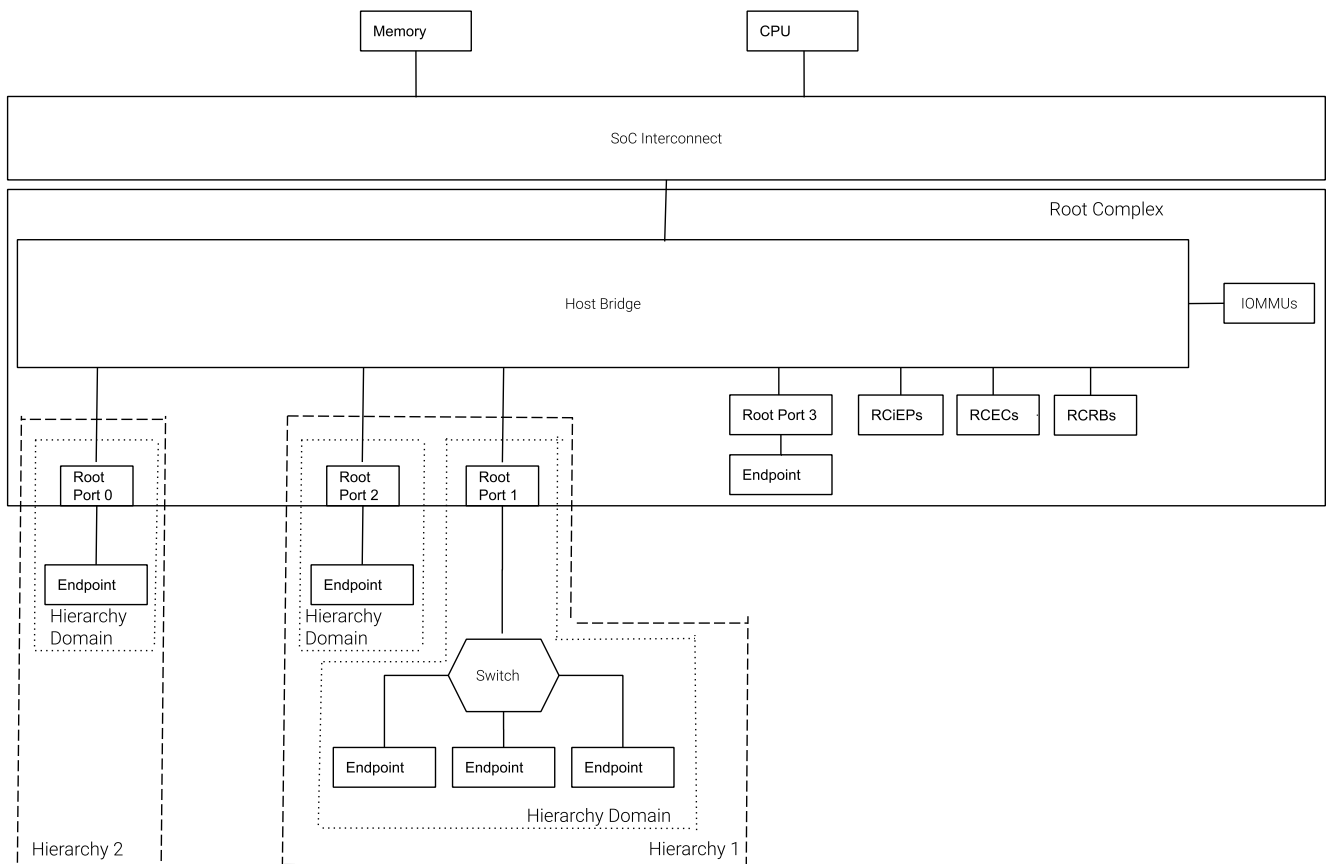


Figure 2. PCIe root complex

One or more root ports in a root complex may be part of a hierarchy where a hierarchy is a PCI

Express I/O interconnect topology, wherein the Configuration Space addresses, referred to as the tuple of Bus/Device/Function Numbers (or just Bus/Function Numbers, for PCIe ARI cases), are unique. These addresses are used for Configuration Request routing, Completion routing, some Message routing, and for other purposes. In some contexts a Hierarchy is also called a Segment, and in Flit Mode, the Segment number is sometimes also included in the ID of a Function. Each root port in a hierarchy originates a hierarchy domain i.e. a part of a Hierarchy originating from a single Root Port. The root ports are PCI-PCI bridges that bridge a primary PCIe bus to a range of secondary and subordinate buses.

In some SoCs, PCIe devices may be integrated in the same package/die as the root complex. Examples of such devices are network controllers, USB host controllers, NVMe controllers, AHCI controllers, etc. Such SoC integrated devices may be presented to software using one of the following options:

1. Presented to software as a PCIe endpoint (EP; See section 1.3.2.2 of the PCIe 6.0 specification) connected to a PCIe root port (See example of such an endpoint connected to root port 3 in [\[RISC-V-Server-RC\]](#)). Such PCIe endpoints must comply with the PCIe specified rules for endpoints.
2. Presented to software as a root complex integrated endpoint (RCiEP; See section 1.3.2.3 of the PCIe 6.0 specification). Such PCIe endpoints must comply with the PCIe specified rules for RCiEP.

Implementing integrated devices that perform as RCiEP or EP allows the use of standardized PCIe frameworks for memory and interrupt resource allocation, virtualization (SR-IOV), ATS/PRI for shared virtual addressing, trusted IO using SPDM/TDISP, RAS frameworks like data poisoning and AER, power management, etc.

The host bridge is placed between the device(s) and the system interconnect to process DMA transactions. Devices perform DMA transactions using IO Virtual Addresses (VA, GVA or GPA). The host bridge invokes the associated IOMMU to translate the IOVA to Supervisor Physical Addresses (SPA).

2.5.1. Enhanced Configuration Access Method (ECAM)

Each PCIe endpoint and the PCIe root port itself implement a set of memory mapped configuration registers that are accessed using the PCIe enhanced configuration access method (ECAM). The memory mapped ECAM address range for a hierarchy is up to 256 MiB in size and the base address of the range is naturally aligned to the size. Each PCIe function is associated with a 4 KiB page in this range such that the address bits (20+b):20 where b=0 to 7 identify the bus number of that function (see also recommendations in the PCIe specification 6.0 section 7.2.2), the address bits 19:15 identify the device number, and the address bits 14:12 identify the function number. The host bridge in conjunction with the SoC firmware maps the ECAM address range to the hierarchy domain originating at each PCIe root port.

ID#	Requirement
ECM_010	<p>The ECAM address ranges MUST have the following physical memory attributes (PMAs):</p> <ul style="list-style-type: none"> • Not cacheable, non-idempotent, coherent, strongly-ordered (I/O ordering) I/O region • One, two, and four byte naturally aligned read and write MUST be supported.
ECM_020	<p>Writes to the ECAM address range from a RISC-V hart MUST be non-posted and the write MUST complete at the hart only after a completion is received from the function hosting the accessed configuration register.</p>
<p><i>Besides performing a write, software executing on a hart must not require any additional actions to achieve this property.</i></p> <p><i>This requirement satisfies the processor and host bridge implementation requirement mentioned in the “Ordering Considerations for the Enhanced Configuration Access Mechanism” implementation note of the PCIe 6.0 specification.</i></p>	
ECM_030	<p>The ECAM address range for a hierarchy MUST be contiguous and the base address of the range MUST be naturally aligned to the size of the ECAM address range associated with the hierarchy.</p>
ECM_040	<p>A SoC MAY support multiple hierarchies. When multiple hierarchies are supported, the ECAM address range of the hierarchies MUST NOT overlap, but they are NOT REQUIRED to be contiguous.</p>
ECM_050	<p>The configuration space of the PCIe root ports MUST be associated with the primary bus number of the hierarchy associated with the root port.</p>
<p><i>PCIe root ports are PCI-PCI bridges that bridge the primary bus to the secondary/subordinate buses. The root port itself enumerates as a PCI-PCI bridge device on the primary bus. The collection of primary, secondary, and subordinate buses are part of a single hierarchy domain that originates at that PCIe root port.</i></p>	
ECM_060	<p>The configuration space of functions on the primary bus MUST be accessible irrespective of the state of the corresponding PCIe link.</p>
<p><i>Discovery and activation of the PCIe link requires accessing the configuration space registers of the PCIe root port itself and the PCIe root port is a PCI-PCI bridge device on the primary bus.</i></p>	
ECM_070	<p>The PCIe root port MUST enumerate as a PCI-PCI bridge to software and comply with the rules specified for PCIe root ports in PCIe specification 6.0.</p>
ECM_080	<p>The PCIe root port MUST support the PCIe Configuration RRS software (CRS) visibility enable control.</p>
<p><i>The number of times a configuration request is retried on an RRS response is UNSPECIFIED.</i></p>	
ECM_090	<p>Read and/or write to the ECAM range of the hierarchy domain originating at a root port MUST generate PCIe configuration transactions as type 0 or type 1 configuration transactions following the rules specified for ECAM in PCIe specification 6.0.</p>

ID#	Requirement
ECM_100	<p>Read access to ECAM address range from a RISC-V hart MUST be responded with all 1s data if any of the following conditions are TRUE:</p> <ul style="list-style-type: none"> • Access is to non-existent functions on the primary bus of a hierarchy domain. • Accessed bus is not part of any of the hierarchy domains. • An Unsupported Request or Completer Abort response was received. • A completion timeout occurs. • Access targets a function downstream of a root port whose link is not in DL_Active state. • A PCIe RRS response was received and CRS software visibility is not enabled. • PCIe CRS software visibility is enabled, but the access does not target the vendor ID register, and a RRS response was received on each retry of the configuration read.
<p><i>The data response to the Vendor ID register on receipt of an RRS response MUST follow the PCIe defined rules.</i></p> <p><i>See also the recommendations in PCIe specification 6.0 section 2.3.2.</i></p>	
ECM_110	<p>Write access from a RISC-V hart to configuration registers of non-existent functions on the primary bus MUST be dropped (silently ignored or discarded) and the write completed. Such accesses MUST NOT lead to any other behavior (e.g., hangs, deadlocks, etc.).</p>
ECM_120	<p>Poisoned data received from completers (EP=1) MUST be forwarded to the requesting RISC-V hart as poisoned data unless such forwarding is disallowed (e.g., SoC does not support data poisoning or forwarding of poisoned data is disabled though implementation defined means). If forwarding of poisoned data is disallowed then the poisoned data MUST be replaced with all 1s data.</p>

2.5.2. PCIe Memory Space

ID#	Requirement
MMS_010	<p>The SoC MUST support designating, for each hierarchy domain, one or more ranges of system physical addresses that may be used for mapping memory space of endpoints in that hierarchy domain using the 64-bit wide base address registers (BARs) of the endpoints.</p>
MMS_020	<p>SoC MUST support designating, for each hierarchy domain, at least one system physical address range for mapping memory space of endpoints in that hierarchy domain using 32-bit wide BARs of the endpoint.</p>
<p><i>The ranges suitable for mapping using 32-bit BARs are also sometimes termed as the low MMIO ranges and those suitable for use with 64-bit BARs termed as high MMIO ranges.</i></p>	

ID#	Requirement
MMS_030	<p>The system physical address ranges designated for mapping endpoint memory spaces MUST have the following physical memory attributes (PMAs):</p> <ul style="list-style-type: none"> • Not cacheable, non-idempotent, coherent, strongly-ordered (I/O ordering) I/O region. • MUST support all aligned and unaligned access sizes that can be generated by data requests from any of the RISC-V application processor harts in the SoC or by peer endpoints, including those of type RCiEP. • MAY support atomics, instruction fetch, and page walks.
<p><i>Software may use the Svpbmt extension to override the PMA to NC if such an override is compatible with the restricted programming model of the device.</i></p> <p><i>See also the implementation note on optimizations based on restricted programming mode in section 2.3.1 of PCIe specification 6.0.</i></p>	
MMS_040	<p>A load from a RISC-V application processor hart to memory ranges designated for the mapped memory space of endpoints or RCiEP MUST complete with an all 1s response and MUST NOT lead to any other behavior (e.g., hangs, deadlocks, etc.) if any of the following are TRUE:</p> <ul style="list-style-type: none"> • Address is not within any of the following: <ul style="list-style-type: none"> ◦ Memory base/limit or prefetchable memory base/limit of any root port. ◦ BAR (including when EA capability is used) mapped range of any RCiEP. ◦ BAR (including when EA capability is used) mapped range of any root port. • The PCIe link of the root port to which the access is routed is not active. <ul style="list-style-type: none"> ◦ Including due to the root port entering downstream port containment state. • A UR or a CA response is received from the completer. • A completion timeout occurs.
MMS_050	<p>A store from a RISC-V application processor hart to memory ranges designated for the mapped memory space of endpoints or RCiEP MUST be dropped (silently ignored or discarded) and MUST NOT lead to any other behavior (e.g., hangs, deadlocks, etc.) if any of the following are TRUE:</p> <ul style="list-style-type: none"> • Address is not within any of the following: <ul style="list-style-type: none"> ◦ Memory base/limit or prefetchable memory base/limit of any root port. ◦ BAR (including when EA capability is used) mapped range of any RCiEP. ◦ BAR (including when EA capability is used) mapped range of any root port. • The PCIe link of the root port to which the access is routed is not active. <ul style="list-style-type: none"> ◦ Including due to the root port entering downstream port containment state.

ID#	Requirement
MMS_060	Poisoned data received from completers (EP=1) MUST be forwarded to the requester PCIe device (a RCiEP or an endpoint) as poisoned data unless such forwarding is disallowed (e.g., poisoned TLP egress blocking).
MMS_070	Poisoned data received from completers (EP=1) MUST be forwarded to a requester RISC-V hart as poisoned data unless such forwarding is disallowed through implementation defined means. When such forwarding is disallowed, then the poisoned data MUST be replaced with all 1s data.
MMS_080	SoC MUST NOT use EA capability to indicate memory for allocation to endpoints downstream of a PCIe root port.

2.5.3. Access Control Services (ACS)

The PCIe ACS provides controls on routing of PCIe TLPs. ACS controls may be used to determine whether the TLP should be routed normally, blocked, or redirected. These controls may be applicable to the root complex, switches, multi-function devices, and SR-IOV capable devices.

ID#	Requirement
ACS_010	<p>PCIe root ports and SoC integrated downstream switch ports MUST support the following PCIe access control services (ACS) controls</p> <ul style="list-style-type: none"> • ACS source validation. • ACS translation blocking. • ACS I/O request blocking.
ACS_020	If a PCIe root port or a SoC-integrated downstream switch port implements memory BAR space, then it SHOULD support the PCIe ACS DSP memory target access control.
<i>The ACS DSP memory target access control can be used to prevent unauthorized accesses to protected memory spaces such as the PCIe root port's BAR mapped registers.</i>	
ACS_030	<p>Root ports and SoC-integrated downstream switch ports that support direct routing between root ports or direct routing from ingress to egress port of a root port MUST support the following PCIe ACS controls:</p> <ul style="list-style-type: none"> • ACS P2P request redirect. • ACS P2P completion redirect. • ACS upstream forwarding. • ACS direct translated P2P.
ACS_040	Root ports and SoC-integrated downstream switch ports that support direct routing between root ports or direct routing from ingress to egress port of a root port SHOULD also support ACS P2P egress control.

ID#	Requirement
	<i>More commonly, P2P routing is accomplished by forwarding the TLP to the host bridge for routing. For further information, refer to the application note accompanying Fig 2-14 and Section 1.3.1 of the PCIe specification 6.0.</i>
ACS_050	The ACS features, including the detection, logging, and reporting of ACS violations MUST adhere to the rules outlined in the PCIe specifications 6.0.

2.5.4. Address Routed Transactions

The rules in this section apply to treatment in the root complex of TLPs that are routed by address. An address carried in such transactions may be the address of a host memory location or the address of a location in the memory space of a peer endpoint or RCiEP.

ID#	Requirement
ADR_010	The host bridge MUST request IOMMU translations for addresses (Translated, Untranslated, or a PCIe ATS address translation request) used in the request by endpoints and RCiEPs.
	<i>The IOMMU must be invoked even for Translated requests to allow determination of whether the requester is configured by software to use Translated requests.</i>
	<i>When the IOMMU operates in the T2GPA mode, it provides a GPA as the translated address in response to a PCIe ATS address translation requests. In this mode of operation, the IOMMU must be invoked by the host bridge for Translated requests to translate the GPA to an SPA.</i>
	<i>When ACS direct translated P2P controls are enabled, the Translated requests may not be routed through the host bridge. In such cases, if direct P2P routing of these requests is not desired, due to security and/or functional reasons (e.g., when operating in T2GPA mode), software should utilize the ACS controls to direct these requests to the root complex.</i>
ADR_020	The host bridge MUST enforce physical memory attribute checks and physical memory protection checks on the translated address provided by the IOMMU and MUST treat violating requests as Unsupported Requests.
	<i>These checks are analogous to the PMA and PMP checks performed by the RISC-V hart.</i>
ADR_030	For Translated and Untranslated requests, the host bridge MUST use the translated addresses provided by the IOMMU to determine whether the transaction is targeting host memory or peer device memory.
ADR_040	The host bridge MAY support access to peer device memory. If peer device memory access is not enabled (by design or by configuration), then such accesses MUST be responded to with a UR/CA response. The host bridge MUST NOT cause any other error (e.g., hang, deadlock, etc.) when rejecting access to peer device memory.

ID#	Requirement
<i>A virtual machine may violate the peer-to-peer access policies and/or configurations established by the hypervisor and/or SoC firmware, which disallow peer device memory accesses. Such a VM may attempt to program devices passed through to the virtual machine to perform peer memory accesses. Such attempts to violate the peer-to-peer policies MUST NOT lead to system instabilities (e.g., hangs, deadlocks, etc.) or errors.</i>	
ADR_050	When a posted or non-posted-with-data request is allowed to access peer device memory, then any poisoned data (EP=1) MUST be forwarded as poisoned data, unless such forwarding is disallowed (e.g., poisoned TLP egress blocking or lack of support for data poisoning is the SoC).
ADR_060	Host memory writes caused by posted or non-posted-with-data requests with poisoned data (EP=1) MUST mark such data in the host memory as poisoned.
ADR_070	Host memory reads that have uncorrectable data errors detected within the SoC MUST cause a response with poisoned data (EP=1) if transmission of poisoned TLPs is not blocked (see also section 2.7.2.1 of PCIe specification 6.0).

2.5.5. ID Routed Transactions

The rules in this section apply to treatment in the root complex of TLPs that are routed by ID. Such requests may be ID Configuration requests, ID routed messages or completions.

ID#	Requirement
IDR_010	Configuration requests from endpoints and RCiEP MUST be treated as Unsupported Requests.
IDR_020	P2P routing of PCIe VDM between root ports within or across hierarchies SHOULD be supported.
<i>MCTP transport protocols using PCIe VDM are used by the BMC to manage PCIe/CXL devices. These messages are used to support manageability protocols such as PLDM, NVMe-MI, Redfish, etc.</i>	
IDR_030	P2P routing of PCIe VDM to/from RCiEP MAY be supported.

2.5.6. Cacheability and Coherence

ID#	Requirement
CCA_010	The host bridge MUST enforce PCIe memory ordering rules and SHOULD support the relaxed ordering (RO) and ID-based ordering (IDO).
<i>An implementation may occasionally or never permit the relaxations allowed by RO and/or IDO attributes. Such implementations will result in a more conservative interpretation of the ordering rules, but they will not result in a violation of the ordering rules.</i>	
CCA_020	Writes to host or device memory using the RO attribute set to 0 MUST be observed by other harts and bus mastering devices in the order in which the write was received by the PCIe root port or the host bridge such that all previous writes are globally observed before the RO=0 write is globally observed.

ID#	Requirement
CCS_030	The host bridge MUST enforce the idempotency, coherence, cacheability, and access type physical memory attributes of the accessed memory and perform any reordering or combining of PCIe transactions only if the combination of physical memory attributes and TLP specified memory ordering attributes allow it.
CCS_040	The host bridge SHOULD implement hardware enforced cache coherency, irrespective of the “No Snoop” attribute in the TLP, unless it has been configured through UNSPECIFIED means to not enforce coherency for TLPs with “No Snoop” attribute set to 1.
<p><i>A PCIe requester is permitted to set the “No Snoop” in transactions it initiates that do not require hardware enforced cache coherency. Host bridges that do not support isochronous VCs or can meet deadlines with hardware enforced coherency may always enforce coherency. Enforcing cache coherency is always conservative and will not lead to data corruption.</i></p> <p><i>Modern systems with integrated memory controllers and snoop directories may not require the use of “No Snoop” to meet the latency targets as memory regions accessed for isochronous operations would usually be device exclusive. PCIe requires a function to guarantee that addresses accessed using “No Snoop” set to 1 are not cached in any of the caches and software that instructs a device to perform “No Snoop” transactions must only do so when it can provide this guarantee.</i></p> <p><i>Some caches in a SoC may perform clean evictions to memory. In such SoCs, if the addresses used by the non-snooped transactions may be cached (e.g., due to speculative accesses from a hart), then such clean evictions may cause data corruption, even if the caches were explicitly cleaned by software using the cache management operations. To ensure data integrity, software should use memory that has such non-cacheable PMA or use the Svpbmt extension to override the PMA to NC/IO, thereby implementing the guarantee required by the PCIe specification when using the “No Snoop” attribute set to 1. If the Svpbmt extension was used to override the PMA, then use of cache management operations defined by Zicbom extension may be necessary to flush data that might already be cached.</i></p> <p><i>See also section 7.5.3.4 of PCIe specification 6.0.</i></p>	
CCS_050	The host bridge MUST NOT violate the coherence physical memory attribute if the “No Snoop” attribute in the TLP is 0.
CCS_060	The interpretation of the TLP processing hints (TPH) by the SoC is UNSPECIFIED .
<p><i>A future extension of the RISC-V IOMMU specification may define a standard interpretation of the TPH including the use of ATS memory attributes (AMA) for performing cache management.</i></p>	

2.5.7. Message signaled interrupts

A message signaled interrupt (MSI or MSI-X) is the preferred interrupt signaling mechanism in PCIe.

ID#	Requirement
MSI_010	Message Signaled Interrupts MUST be supported.

ID#	Requirement
MSI_020	SoC MUST NOT require any further action from the operating system besides configuring the MSI address register in devices with the address of an IMSIC interrupt register file (or a virtual interrupt file) and the MSI data register in devices with an external interrupt identity to enable the use of MSI or MSI-X.
MSI_030	SoC MUST NOT support INTx virtual wire based interrupt signaling.
<i>PCIe supports INTx emulation to support legacy PCI interrupt mechanisms. Modern SoC and devices are not expected be limited by the lack of this emulation mode.</i>	

2.5.8. Precision Time Measurement (PTM)

ID#	Requirement
PTM_010	PCIe root ports MAY support PCIe PTM capability.
<i>Several applications such as instrumentation, media servers, telecom servers, etc. require high precision monitoring and tracking of time. The PCIe PTM protocol supports synchronization of multiple devices/functions to a common shared PTM master time provided by the PTM root.</i>	
PTM_020	When PCIe PTM capability is supported, the SoC MUST make the PTM master time available to the operating system.
<i>The mechanism to make the master time available to the operating system is implementation specific.</i>	
<i>Making PTM master time available to software enables software to translate timing information between local time and PTM master time and thereby enable coordination of events across multiple PCIe devices.</i>	
PTM_030	When PCIe PTM capability is supported, the PTM master time MUST be 64-bit wide and MUST use the same or higher resolution clock than the clock used to provide time .

2.5.9. Error and Event Reporting

ID#	Requirement
AER_010	PCIe root ports MUST support advanced error reporting (AER) capability for reporting errors from connected devices or the errors detected by the root port itself.
<i>AER capability defines more robust error reporting as compared to the baseline error reporting capability.</i>	
AER_020	PCIe root ports MUST support the downstream port containment (DPC) capability.
AER_030	PCIe root ports MUST support the RP PIO controls.
<i>The root port programmed I/O (PIO) controls enable fine-grained control over handling of non-posted requests that encounter errors and allows handling of such errors as either uncorrectable or advisory based on policies established by the operating system.</i>	

ID#	Requirement
AER_040	A RCiEP in the SoC SHOULD support the AER capability if it detects any of the errors defined by PCIe specification 6.0 (See section 6.2.7).
AER_050	A RCiEP in the SoC MUST support the AER capability if it supports the ACS capability.
AER_060	SoC MUST implement one or more PCIe RCEC in the root complex if any of the RCiEP implement the AER capability or implement PME signaling.
AER_070	The PCIe RCEC implemented in a SoC MUST implement the RCEC endpoint association extended capability.

2.5.10. Vendor Specific Registers

ID#	Requirement
VSR_010	<p>Vendors specific registers in the root ports, host bridge, RCiEP, and RCRB MUST be implemented using one or more of the following capabilities:</p> <ul style="list-style-type: none"> • Vendor specific capability. • Vendor specific extended capability. • Designated Vendor Specific extended capability.
VSR_020	SoC MUST NOT require hypervisor and/or operating system interaction with configuration space registers that are not defined by an industry standards. Non-standard vendor specific registers, if implemented in the configuration space, must only be used by the SoC firmware.

Some industry standards such as CXL may define standard DVSEC structures in the configuration space.

The preferred way to implement device/SoC vendor specific registers that need to be used by drivers in the run-time environment is to implement them in the memory space of the device. Certain operating systems and hypervisors may disallow and/or require mediating access to the configuration space of devices. See also the implementation note in the PCIe specification 6.0 section 7.2.2.2.

2.5.11. SoC-Integrated PCIe Devices

ID#	Requirement
SID_010	SoC-integrated PCIe devices MUST implement all software visible rules defined by the PCIe specification 6.0 for an EP or RCiEP as applicable.

Implementing integrated devices as RCiEP or EP allows the use of standardized frameworks for memory and interrupt resource allocation, virtualization (SR-IOV), ATS/PRI, shared virtual addressing, trusted IO using SPDM/TDISP, participate in RAS frameworks like data poisoning and AER, power management, etc.

ID#	Requirement
SID_020	SoC-integrated PCIe devices MUST NOT use legacy PCI capabilities. They MUST NOT require the use of I/O space, I/O transactions, or the INTx virtual wire interrupt signaling mechanism.
SID_030	SoC integrated PCIe devices that cache address translations MUST implement the PCIe ATS capability if the address translation cache needs management by the operating system or hypervisors.
SID_040	SoC-integrated PCIe devices that support PCIe SR-IOV capability SHOULD support MSI-X capability.
<i>MSI-X capability enables virtual machines to assign interrupt resources to virtual functions without needing access to the configuration space of the function. Access to the configuration space of the virtual function is usually mediated by the hypervisor.</i>	
SID_050	SoC-integrated PCIe devices MAY support the PASID capability. When PASID capability is supported, the devices SHOULD support a 20-bit wide PASID.
<i>Endpoints are recommended to support a 20-bit wide PASID to ensure interoperability with system software. See also the implementation note on PASID width homogeneity in the PCIe specification 6.0 section 6.20.2.2.</i>	
SID_060	SoC-integrated PCIe devices (a multi-function device or an SR-IOV capable device) that support P2P traffic among functions (including among SR-IOV virtual functions) of the device MUST support the following PCIe ACS controls: <ul style="list-style-type: none"> • ACS P2P request redirect. • ACS P2P completion redirect. • ACS direct translated P2P.
SID_070	SoC-integrated PCIe devices MAY support programmable BAR registers. When programmable BAR registers are supported, the Memory Space Indicator (bit 0) of such BAR MUST be 1 and they SHOULD support being mapped anywhere in the 64 bit memory space.
SID_080	RCiEP MAY support the PCIe enhanced allocation (EA) capability for fixed allocation of memory resources.
SID_090	SoC-integrated PCIe devices MUST support the PCIe defined baseline error reporting capability and MAY support PCIe Advanced Error Reporting capability. If PCIe ACS controls are supported then the PCIe Advanced Error Reporting capability MUST be supported.
<i>See PCIe specification 6.0 section 7.5.1.1.14.</i>	
SID_100	A RCiEP that supports PCIe Advanced Error Reporting MUST be associated with a Root Complex Event Collector.

2.6. Reliability, Availability, and Serviceability (RAS)

ID#	Requirement
RAS_010	The level of RAS implemented by the SoC is UNSPECIFIED .
<p><i>The level of RAS implemented by an SoC depends on the reliability goals established for the SoC, which are commonly measured using metrics such as failure-in-time (FIT) and defects-per-million (DPM). Achieving these goals requires a combination of fault prevention, error detection, and error correction techniques.</i></p> <p><i>This specification strongly recommends the implementation of error detection and correction codes for storage elements like significant caches and memories. Furthermore, it suggests utilizing mechanisms such as single-symbol (SSC) ECC in DRAM controllers to address failure scenarios, such as when all bits in a single DRAM device experience a failure.</i></p> <p><i>Additionally, this specification encourages the adoption of mechanisms like periodic scrubbing, also known as patrol scrubbing. These mechanisms proactively identify and rectify errors before they accumulate to a critical point, surpassing the capability of the implemented error correction codes. For instance, this could involve addressing situations where single bit errors escalate into double bit errors, surpassing the correction code's capacity.</i></p>	
RAS_020	SoC SHOULD support the generation, storage, and forwarding of poisoned data. The granularity at which data is poisoned is UNSPECIFIED .
<p><i>When an uncorrected data error is detected by a component, it might allow potentially corrupted data to reach the data requester, but with an associated poison indicator. These errors are referred to as uncorrected deferred errors (UDE), as they enable the detecting component to continue functioning and postpone addressing the error until a later time, assuming the poisoned data gets consumed. If a component (such as a hart, an IOMMU, a device, etc.) consumes the poisoned data, it triggers an uncorrected urgent error (UUE), leading to the invocation of a recovery handler for immediate remedial actions, as further deferral of the error is not feasible.</i></p> <p><i>The technique of data poisoning facilitates delaying the handling of uncorrected errors until the moment the corrupted data is actually consumed. Data poisoning offers a more precise identification of the software and/or hardware component affected by the data corruption. This specificity allows for targeted recovery actions that impact only the affected components.</i></p> <p><i>To ensure the integrity of the poisoned data indicator when stored, error detection and correction codes should be applied. This practice prevents subsequent errors from leading to the silent consumption of the corrupted data.</i></p> <p><i>Data poisoning also empowers the implementation of error containment features supported by industry standards like PCIe and CXL.</i></p> <p><i>For more detailed discussions on the treatment of faults and errors, refer to the RISC-V RERI specification.</i></p>	

ID#	Requirement
RAS_030	If poisoned data needs to be transmitted from a first component to a second component that lacks the ability to manage poison, the first component MUST trigger an uncorrected urgent error report instead of silently transmitting the corrupted data.
<p><i>Some components serve as intermediaries through which data passes. For instance, a PCIe/CXL port acts as an intermediary that receives data from memory but doesn't consume it; rather, it forwards the data to an endpoint. In such cases, the intermediary component might encounter poisoned data. While this component can propagate the error and avoid logging an error, a different scenario arises when the destination component (such as a PCIe endpoint) cannot handle poison. In such situations, the originating component must trigger an urgent error signal instead of transmitting the poisoned data without the associated poison indicator. Failing to do so would breach the containment of the corrupted data during propagation.</i></p>	
RAS_040	The SoC SHOULD support the RISC-V RAS error record register interface (RERI) [9] for error logging and signaling.
<p><i>Note RERI is still under construction.</i></p>	
RAS_050	When RERI is supported, the RAS error records MUST include the capability to individually enable error signals for each severity (UUE, UDE, or CE) of error that could be logged in that specific error record.
<p><i>Configurable enables provide software with the flexibility of using an event-based or polling-based error logging for both corrected errors and deferred errors. Typically, software operates in an event-based mode for urgent errors, as these errors necessitate immediate remedial action when they arise.</i></p>	
RAS_060	If RERI is supported, RAS error records MUST preserve the state of logged error information (including status, address, information, supplemental information, and timestamp) across a RAS-initiated reset. The state of RAS error records MAY persist across other types of implementation-defined resets. After a reset, including a RAS-initiated reset, the state of the control register is considered UNSPECIFIED .
<p><i>Some errors may lead a hardware component to enter a failure mode in which it becomes incapable of servicing additional requests- colloquially termed 'jammed' or 'wedged'. In these situations, the SoC may require a reset to restore it to an operational state (a RAS-initiated reset). Preserving the RAS error records through such resets enables the SoC firmware and system software to retrieve these error records during boot following such a reset, facilitating logging and analysis.</i></p>	
RAS_070	If RERI is supported, the RAS error records MAY support error record injection, which is intended to facilitate RAS handler verification.

ID#	Requirement
<p><i>Verifying the correct implementation of RAS handlers presents a formidable challenge, given the impracticality of deterministically inducing all potential errors within the SoC to validate the RAS handler's adherence to desired recovery protocols. An unverified RAS handler can lead to undesired behavior during error occurrences, potentially reducing SoC availability or affecting its serviceability.</i></p> <p><i>To address this, error record injection offers a convenient method for conducting such verification. It allows the introduction of a range of error signatures, which can then be signaled and observed. While hardware error injection techniques also offer a means of verification (e.g., methods to intentionally corrupt a data location protected by an error detection code), providing open access to these capabilities for software use might not align with security and stability concerns.</i></p>	
RAS_080	<p>If RERI is supported, then the hardware components in the SoC that support error correction MUST incorporate a corrected error counter within their respective error records. Additionally, these components MUST support the signaling of counter overflows.</p> <p><i>Counting corrected errors offers a more precise assessment of system reliability. Enabling signaling upon counter overflow empowers software to define a suitable threshold for logging and analysis of these corrected errors.</i></p> <p><i>Certain hardware units might maintain a history of corrected errors and increment the corrected error counter only if the error differs from a previously reported one. Additionally, some hardware units could incorporate low-pass filters like leaky buckets, which regulate the rate at which corrected errors are reported and counted. This requirement pertains to corrected errors tracked by the error record once the hardware component determines reporting and counting based on its specific filtering rules.</i></p>

2.7. Quality of Service

Quality of Service (QoS) refers to the minimum end-to-end performance that a service level agreement (SLA) guarantees to an application in advance. QoS capabilities within the SoC offer mechanisms that system software can leverage to manage interference to an application, effectively diminishing performance variability caused by other applications' utilization of shared resources such as cache capacity, memory bandwidth, interconnect bandwidth, power consumption, and more.

ID#	Requirement
QOS_010	The SoC SHOULD incorporate QoS mechanisms to mitigate unwarranted performance interference that arises when multiple workloads access shared resources like caches and system memory.
QOS_020	The SoC SHOULD integrate support for the RISC-V capacity and bandwidth controller register interface (CBQRI) [10] in significant shared caches and the memory controllers.

ID#	Requirement
QOS_030	If CBQRI is supported, RISC-V harts within the application processors of the SoC MUST include support for the sqoscfg CSR. Furthermore, this CSR MUST support a minimum of 16 RCIDs and at least 32 MCIDs.
<i>The count of RCID and MCID that can be used in the SoC should scale with the number of RISC-V harts in the SoC.</i>	
QOS_040	If CBQRI is supported, the IOMMU in the SoC SHOULD incorporate support for the CBQRI-defined extension, enabling the association of RCID and MCID with requests initiated by devices and the IOMMU.
QOS_050	If CBQRI is supported, significant caches such as the last-level cache in the SoC SHOULD support cache capacity allocation.
QOS_060	If CBQRI is supported, significant caches such as the last-level cache in the SoC SHOULD incorporate support for monitoring cache capacity usage.
QOS_070	If CBQRI is supported, the memory controllers within the SoC SHOULD include support for bandwidth allocation.
QOS_080	If CBQRI is supported, the memory controllers in the SoC SHOULD include support for monitoring bandwidth usage.
<i>The method employed by the SoC for bandwidth throttling and control is specific to its implementation. It is advisable for the implementation to utilize a scheme that results in a deviation of no more than +/- 10 % from the target set by system software through the CBQRI interface.</i>	
QOS_090	If CBQRI is supported, the count of RCID and MCID supported by capacity controllers, bandwidth controllers, and all RISC-V application processor harts in the SoC MUST be consistent.
<i>Portable system software could opt to limit itself to accommodating the minimum count of RCID and MCID across the controllers. This approach avoids the complexity of dealing with unequal numbers of RCID and MCID across controllers, which would otherwise necessitate intricate allocations and constraints on workload placement.</i>	
QOS_100	If CBQRI is supported, the monitoring counters in the capacity and bandwidth controllers MUST be sufficiently wide to not overflow when sampled at a rate of 1 Hz.
<i>As an illustration, consider an HBM3 memory interface that can facilitate data transfers at a rate of up to 1 TB/s. This scenario would necessitate a 34-bit counter to prevent overflow when sampled at a frequency of 1 Hz.</i>	

2.8. Manageability

This section outlines the guidelines for RISC-V server SoCs to incorporate a standardized set of protocols and standards for server management. The SoC interfaces with a baseboard management controller (BMC) through in-band and out-of-band (OOB) management agents. The in-band management agents execute on the RISC-V application processor harts and the out-of-band management agents execute on a management controller in the SoC.

The out-of-band management interface facilitates the monitoring of sensors (e.g., temperature,

power, etc.), parameter control (e.g., power limits, etc.), and logging (e.g., RAS error records, etc.) by the BMC without participation of software on the application processor harts. The in-band management interface facilitates system configuration (e.g., boot order, memory domains, secure boot, network, etc.), and event log collection through management agents in the OS and/or firmware that executes on the application processor harts.

This specification strongly recommends the use of the DMTF Redfish [11], DMTF Platform Level Data Model (PLDM) [12], and DMTF Management Component Transport Protocol (MCTP) [13]) protocols for in-band and out-of-band server management.

This specification strongly recommends the use of DMTF specified Security Protocol and Data Model (SPDM) [14] for device attestation and using SPDM encrypted messages [15] for secure in-band and out-of-band communication with the BMC. SPDM authentication protocols support establishing a trust relationship between the manageability agents in the SoC and the BMC. Use of SPDM secured messages enables preserving the confidentiality and integrity of data exchanged between the BMC and the manageability agents in the SoC.

The specification recommends supporting Intelligent Platform Management Interface (IPMI) [16] due to the widespread use of this protocol for server management functions such as credentials provisioning and remote power control.

This specification recommends the RISC-V server SoC to support open standards for server management through supporting integration with technologies such as the datacenter-ready secure control module (DC-SCM) [17] specified by the Open Compute Project for server management, security, and control features.

Adhering to the industry standard management protocols such as those specified by DMTF and OCP allows server platforms built with RISC-V server SoCs to seamlessly integrate into the server management frameworks and tools employed by data centers and enterprises.

ID#	Requirement
MNG_010	The SoC SHOULD incorporate support for an x1 PCIe lane, preferably Gen 5, but at least Gen 3, to establish a connection with the BMC.
<p><i>This interface is commonly linked to a BMC as a PCIe endpoint, serving various purposes. These include facilitating host-to-BMC communication for tasks like video output (e.g., remote KVM support), MCTP transport over PCIe VDM, and hosting a USB controller. The BMC might also support remote presence capabilities, like remote media redirection and support for keyboard and mouse functions through virtual USB.</i></p> <p><i>The in-band network interfaces serve as communication channels for system software to interact with the BMC. This interaction employs protocols like the Redfish host interface.</i></p> <p><i>Furthermore, the PCIe interface to the BMC empowers the BMC, using SoC-routed PCIe VDMs, to utilize these VDMs for transmitting MCTP messages. These messages manage platform devices, including network controllers, NVMe controllers, FPGAs, GPUs, and more.</i></p>	
MNG_020	The SoC SHOULD support the use of I2C based IPMI SSIF for in-band management agents in the SoC to communicate with the BMC.

ID#	Requirement
MNG_030	The SoC SHOULD incorporate support for utilizing a UART connection to the BMC, enabling the provision of a host debug console.
MNG_040	The SoC SHOULD support remote debug using a JTAG interface with the BMC.
<i>The JTAG debug transport modules within the SoC can be utilized to access the RISC-V debug modules in the SoC through the Debug Module Interface.</i>	
MNG_050	The SoC SHOULD support the use of I3C for in-band and out-of-band management agents in the SoC to communicate with the BMC.
<i>The out-of-band interface enables monitoring of sensors (e.g., temperature, power, etc.), control of parameters (e.g, power limits, etc.), and logging (e.g., RAS error records, etc.) by the BMC.</i>	
<i>I3C supports high bandwidth communication with the BMC and is preferred over I2C to minimize server boot times.</i>	

2.9. Debug

ID#	Requirement
DBG_010	The SoC MUST support at least one RISC-V debug module as specified by the RISC-V debug specification [18].
DBG_020	<p>The debug modules MUST support the following capabilities:</p> <ul style="list-style-type: none"> • Program buffer to execute instructions in debug mode. • Support at least one halt group and at least one resume group, besides group 0. • Support debugging harts immediately out of reset. • Always perform program buffer and abstract memory access with a full and exact set of permissions (i.e., hardware <code>relaxedpriv</code> to 0). • Freezing hart local counters using <code>stopcount</code> control.
<i>The ability to halt harts as a group and the ability to halt at reset offers convenience in debugging system software.</i>	
<i>The program buffer provides greater flexibility and ease of use in comparison to abstract command access.</i>	
<i>Enforcing a complete and precise set of permissions on abstract and program buffer memory accesses mitigates security concerns.</i>	

2.10. Trace

ID#	Requirement
TRC_010	The SoC MUST support either the RISC-V E-trace [19] or the RISC-V N-trace [20] standard extension for instruction trace.
TRC_020	The SoC MUST support the RISC-V trace control interface specification.
<i>The RISC-V N-trace and the trace control interface specification are currently under development.</i>	
TRC_030	The trace control interface MUST support filtering by privilege levels.
TRC_040	The SoC SHOULD incorporate support for a System Memory trace sink.
<i>The system memory sink facilitates self-hosted trace decoders for software debugging and profiling purposes. Additionally, the system memory sink conveniently supports the off-chip transmission of trace data through a PCIe or USB ports.</i>	

2.11. Performance Monitoring

ID#	Requirement
SPM_010	<p>Significant caches within the SoC SHOULD incorporate an HPM capable of counting:</p> <ul style="list-style-type: none"> • Cache lookup • Cache miss <div> <p>If the SoC supports NUMA configurations, then the HPM SHOULD support filtering the counting based on whether the request originated in a local node or a remote node.</p> </div>
<i>It is recommended that a cache with a capacity larger than 32 KiB be considered a significant cache.</i>	
SPM_020	<p>The memory controllers within the SoC SHOULD incorporate an HPM capable of counting:</p> <ul style="list-style-type: none"> • Read bandwidth • Write bandwidth <div> <p>If the SoC supports NUMA configurations, then the HPM SHOULD support filtering the counting based on whether the request originated in a local node or a remote node.</p> </div>
SPM_030	<p>The PCIe ports within the SoC SHOULD incorporate an HPM capable of counting:</p> <ul style="list-style-type: none"> • Read bandwidth (from system memory) • Write bandwidth (to system memory)

ID#	Requirement
SPM_040	<p>The SoC SHOULD incorporate an HPM capable of counting the average latency of a read request from a memory requester (e.g., a hart, a PCIe host bridge, etc.) in the SoC.</p> <p>If the SoC supports NUMA configurations, then the HPM SHOULD support filtering the counting based on whether the request is to local memory or to remote memory.</p>
<i>Bandwidth and latency are the most commonly used performance metrics to guide workload placement and tuning.</i>	
SPM_050	The PCIe Gen6 ports within the SoC SHOULD incorporate support for the Flit performance measurement extended capability defined by PCIe specification 6.0.

2.12. Security Requirements

ID#	Requirement
SEC_010	The PCIe root ports within the SoC SHOULD support PCIe Integrity and Data Encryption (IDE) capability.
<i>The IDE extension adds optional capabilities to perform hardware encryption and integrity checks on packets transferred across PCIe links. This addition provides confidentiality, integrity, and replay protection against hardware-level attacks.</i>	
SEC_020	The SoC SHOULD support encryption of off-chip DRAM using a transient memory encryption key that has at least 192-bit key lengths.
<i>Off-chip memory encryption provides protection to critical assets in memory such as credentials, data encryption keys, and other secrets.</i>	
SEC_030	The cryptographic modules used to implement PCIe and off-chip DRAM encryption SHOULD comply with security requirements specified by standards such as FIPS 140-3.
SEC_040	The SoC should have the capability interfacing with a Trusted Platform Module (TPM) that adheres to the TPM 2.0 Library specification [21].
<i>A TPM enhances security by providing secure storage for sensitive information such as credentials and passwords, cryptographic operations and protection against tampering or unauthorized access.</i>	

Bibliography

- [1] “Key words for use in RFCs to Indicate Requirement Levels.” [Online]. Available: datatracker.ietf.org/doc/html/rfc2119.
- [2] “PCI Express® Base Specification Revision 6.0.” [Online]. Available: pcisig.com/pci-express-6.0-specification.
- [3] “Advanced Configuration and Power Interface (ACPI) Specification.” [Online]. Available: uefi.org/specifications.
- [4] “Unified Extensible Firmware Interface.” [Online]. Available: uefi.org/specifications.
- [5] “RVA23 Profiles.” [Online]. Available: github.com/riscv/riscv-profiles/blob/main/rva23-profile.adoc.
- [6] “RISC-V Advanced Interrupt Architecture.” [Online]. Available: github.com/riscv/riscv-aia.
- [7] “RISC-V IOMMU Architecture Specification.” [Online]. Available: github.com/riscv-non-isa/riscv-iommu.
- [8] “PCI Code and ID Assignment Specification Revision 1.1.” [Online]. Available: pcisig.com/sites/default/files/files/PCI_Code-ID_r_1_11_v24_Jan_2019.pdf.
- [9] “RISC-V RAS error record register interface.” [Online]. Available: github.com/riscv/riscv-ras-eri.
- [10] “RISC-V Capacity and Bandwidth QoS Register Interface.” [Online]. Available: github.com/riscv-non-isa/riscv-cbqri.
- [11] “Redfish specification 1.18.0.” [Online]. Available: www.dmtf.org/dsp/DSP0266.
- [12] “PLDM base specification 1.1.0.” [Online]. Available: www.dmtf.org/dsp/DSP0240.
- [13] “MCTP base specification 1.3.1.” [Online]. Available: www.dmtf.org/dsp/DSP0236.
- [14] “Security protocol and data model (SPDM) specification 1.2.1.” [Online]. Available: www.dmtf.org/dsp/DSP0274.
- [15] “Secured messages using SPDM specification 1.1.0.” [Online]. Available: www.dmtf.org/dsp/DSP0277.
- [16] “Intelligent Platform Management Interface (IPMI) 2.0.” [Online]. Available: www.intel.com/content/dam/www/public/us/en/documents/product-briefs/ipmi-second-gen-interface-spec-v2-rev1-1.pdf.
- [17] “Datacenter Secure Control Module Specification.” [Online]. Available: www.opencompute.org/documents/ocp-dc-scm-spec-rev-1-0-pdf.
- [18] “RISC-V Debug Specification.” [Online]. Available: github.com/riscv/riscv-debug-spec.
- [19] “RISC-V Processor Trace Specification.” [Online]. Available: github.com/riscv-non-isa/riscv-

[trace-spec](#).

[20] “RISC-V N-Trace Specification.” [Online]. Available: github.com/riscv-non-isa/tg-nexus-trace.

[21] “TPM 2.0 Library.” [Online]. Available: trustedcomputinggroup.org/resource/tpm-library-specification/.