

The RISC-V Instruction Set Manual

Volume II: Privileged Architecture

Privileged Architecture Version 1.11-draft

Document Version 1.11-draft

Warning! This draft specification may change before being accepted as standard by the RISC-V Foundation. While the editors intend future changes to this specification to be forward compatible, it remains possible that implementations made to this draft specification will not conform to the future standard.

Editors: Andrew Waterman¹, Krste Asanović^{1,2}

¹SiFive Inc.,

²CS Division, EECS Department, University of California, Berkeley

andrew@sifive.com, krste@berkeley.edu

October 17, 2018

Contributors to all versions of the spec in alphabetical order (please contact editors to suggest corrections): Krste Asanović, Peter Ashenden, Rimas Avižienis, Jacob Bachmeyer, Allen J. Baum, Paolo Bonzini, Ruslan Bukin, Christopher Celio, Chuanhua Chang, David Chisnall, Anthony Coulter, Palmer Dabbelt, Monte Dalrymple, Dennis Ferguson, Mike Frysinger, John Hauser, David Horner, Olof Johansson, Yunsup Lee, Andrew Lutomirski, Prashanth Mundkur, Jonathan Neuschäfer, Rishiyur Nikhil, Stefan O'Rear, Albert Ou, John Ousterhout, David Patterson, Dmitri Pavlov, Kade Phillips, Colin Schmidt, Wesley Terpstra, Matt Thomas, Tommy Thorn, Ray VanDeWalker, Megan Wachs, Steve Wallace, Andrew Waterman, Clifford Wolf, and Reinoud Zandijk.

This document is released under a Creative Commons Attribution 4.0 International License.

This document is a derivative of the RISC-V privileged specification version 1.9.1 released under following license: © 2010–2017 Andrew Waterman, Yunsup Lee, Rimas Avižienis, David Patterson, Krste Asanović. Creative Commons Attribution 4.0 International License.

Please cite as: “The RISC-V Instruction Set Manual, Volume II: Privileged Architecture, Version 1.10”, Editors Andrew Waterman and Krste Asanović, RISC-V Foundation, May 2017.

Preface

This is a **draft of** version 1.11 of the RISC-V privileged architecture proposal. Changes from version 1.10 include:

- Improvements to the description and commentary.
- Added a draft proposal for a hypervisor extension.
- Specified which interrupt sources are reserved for standard use.
- Added specification that xRET instructions may, but are not required to, clear LR reservations if A extension present.
- The virtual-memory system no longer permits supervisor mode to execute instructions from user pages, regardless of the SUM setting.
- Made the `mstatus.MPP` field **WARL**, rather than **WLRL**.
- Made the unused `xip` fields **WPRI**, rather than **WIRI**.
- Made the unused `misa` fields **WLRL**, rather than **WIRI**.
- Made the unused `pmpaddr` and `pmpcfg` fields **WARL**, rather than **WIRI**.
- Required all harts in a system to employ the same PTE-update scheme as each other.
- Rectified an editing error that misdescribed the mechanism by which `mstatus.xIE` is written upon an exception.
- Added the constraint that if a PMA supports LR, SC, or AMOs, but traps misaligned LR, SC, or AMOs, then it also must trap misaligned loads and stores.
- Specified the behavior of the `misa` and `xepc` registers in systems with variable IALIGN.
- Specified semantics for PMP regions coarser than four bytes.
- Specified contents of CSRs across XLEN modification.

Preface to Version 1.10

This is version 1.10 of the RISC-V privileged architecture proposal. Changes from version 1.9.1 include:

- The previous version of this document was released under a Creative Commons Attribution 4.0 International License by the original authors, and this and future versions of this document will be released under the same license.
- The explicit convention on shadow CSR addresses has been removed to reclaim CSR space. Shadow CSRs can still be added as needed.
- The `mvendorid` register now contains the JEDEC code of the core provider as opposed to a code supplied by the Foundation. This avoids redundancy and offloads work from the Foundation.
- The interrupt-enable stack discipline has been simplified.
- An optional mechanism to change the base ISA used by supervisor and user modes has been added to the `mstatus` CSR, and the field previously called Base in `misa` has been renamed to `MXL` for consistency.
- Clarified expected use of XS to summarize additional extension state status fields in `mstatus`.
- Optional vectored interrupt support has been added to the `mtvec` and `stvec` CSRs.
- The SEIP and UEIP bits in the `mip` CSR have been redefined to support software injection of external interrupts.
- The `mbadaddr` register has been subsumed by a more general `mtval` register that can now capture bad instruction bits on an illegal instruction fault to speed instruction emulation.
- The machine-mode base-and-bounds translation and protection schemes have been removed from the specification as part of moving the virtual memory configuration to `sptbr` (now `satp`). Some of the motivation for the base and bound schemes are now covered by the PMP registers, but space remains available in `mstatus` to add these back at a later date if deemed useful.
- In systems with only M-mode, or with both M-mode and U-mode but without U-mode trap support, the `medeleg` and `mideleg` registers now do not exist, whereas previously they returned zero.
- Virtual-memory page faults now have `mcause` values distinct from physical-memory access exceptions. Page-fault exceptions can now be delegated to S-mode without delegating exceptions generated by PMA and PMP checks.
- An optional physical-memory protection (PMP) scheme has been proposed.
- The supervisor virtual memory configuration has been moved from the `mstatus` register to the `sptbr` register. Accordingly, the `sptbr` register has been renamed to `satp` (Supervisor Address Translation and Protection) to reflect its broadened role.
- The SFENCE.VM instruction has been removed in favor of the improved SFENCE.VMA instruction.
- The `mstatus` bit MXR has been exposed to S-mode via `sstatus`.
- The polarity of the PUM bit in `sstatus` has been inverted to shorten code sequences involving MXR. The bit has been renamed to SUM.
- Hardware management of page-table entry Accessed and Dirty bits has been made optional; simpler implementations may trap to software to set them.

- The counter-enable scheme has changed, so that S-mode can control availability of counters to U-mode.
- H-mode has been removed, as we are focusing on recursive virtualization support in S-mode. The encoding space has been reserved and may be repurposed at a later date.
- A mechanism to improve virtualization performance by trapping S-mode virtual-memory management operations has been added.
- The Supervisor Binary Interface (SBI) chapter has been removed, so that it can be maintained as a separate specification.

Preface to Version 1.9.1

This is version 1.9.1 of the RISC-V privileged architecture proposal. Changes from version 1.9 include:

- Numerous additions and improvements to the commentary sections.
- Change configuration string proposal to be use a search process that supports various formats including Device Tree String and flattened Device Tree.
- Made `misa` optionally writable to support modifying base and supported ISA extensions. CSR address of `misa` changed.
- Added description of debug mode and debug CSRs.
- Added a hardware performance monitoring scheme. Simplified the handling of existing hardware counters, removing privileged versions of the counters and the corresponding delta registers.
- Fixed description of SPIE in presence of user-level interrupts.

Contents

Preface	i
1 Introduction	1
1.1 RISC-V Privileged Software Stack Terminology	1
1.2 Privilege Levels	2
1.3 Debug Mode	4
2 Control and Status Registers (CSRs)	5
2.1 CSR Address Mapping Conventions	5
2.2 CSR Listing	7
2.3 CSR Field Specifications	11
2.4 CSR Width Modulation	12
3 Machine-Level ISA, version 1.11-draft	13
3.1 Machine-Level CSRs	13
3.1.1 Machine ISA Register <code>misa</code>	13
3.1.2 Machine Vendor ID Register <code>mvendorid</code>	17
3.1.3 Machine Architecture ID Register <code>marchid</code>	17
3.1.4 Machine Implementation ID Register <code>mimpid</code>	18
3.1.5 Hart ID Register <code>mhartid</code>	18
3.1.6 Machine Status Register (<code>mstatus</code>)	18
3.1.7 Privilege and Global Interrupt-Enable Stack in <code>mstatus</code> register	19

3.1.8	Base ISA Control in <code>mstatus</code> Register	20
3.1.9	Memory Privilege in <code>mstatus</code> Register	21
3.1.10	Virtualization Support in <code>mstatus</code> Register	21
3.1.11	Extension Context Status in <code>mstatus</code> Register	22
3.1.12	Machine Trap-Vector Base-Address Register (<code>mtvec</code>)	25
3.1.13	Machine Trap Delegation Registers (<code>medeleg</code> and <code>mideleg</code>)	26
3.1.14	Machine Interrupt Registers (<code>mip</code> and <code>mie</code>)	28
3.1.15	Machine Timer Registers (<code>mtime</code> and <code>mtimecmp</code>)	30
3.1.16	Hardware Performance Monitor	31
3.1.17	Counter-Enable Registers (<code>[m s]counteren</code>)	32
3.1.18	Machine Scratch Register (<code>mscratch</code>)	33
3.1.19	Machine Exception Program Counter (<code>mepc</code>)	34
3.1.20	Machine Cause Register (<code>mcause</code>)	34
3.1.21	Machine Trap Value (<code>mtval</code>) Register	36
3.2	Machine-Mode Privileged Instructions	37
3.2.1	Environment Call and Breakpoint	37
3.2.2	Trap-Return Instructions	38
3.2.3	Wait for Interrupt	38
3.3	Reset	39
3.4	Non-Maskable Interrupts	40
3.5	Physical Memory Attributes	40
3.5.1	Main Memory versus I/O versus Empty Regions	42
3.5.2	Supported Access Type PMAs	42
3.5.3	Atomicity PMAs	42
3.5.4	Memory-Ordering PMAs	43
3.5.5	Coherence and Cacheability PMAs	44
3.5.6	Idempotency PMAs	45
3.6	Physical Memory Protection	45

3.6.1	Physical Memory Protection CSRs	46
3.6.2	Physical Memory Protection and Paging	50
4	Supervisor-Level ISA, Version 1.11-draft	51
4.1	Supervisor CSRs	51
4.1.1	Supervisor Status Register (sstatus)	51
4.1.2	Base ISA Control in sstatus Register	52
4.1.3	Memory Privilege in sstatus Register	53
4.1.4	Supervisor Trap Vector Base Address Register (stvec)	53
4.1.5	Supervisor Interrupt Registers (sip and sie)	54
4.1.6	Supervisor Timers and Performance Counters	55
4.1.7	Counter-Enable Register (scounteren)	56
4.1.8	Supervisor Scratch Register (sscratch)	56
4.1.9	Supervisor Exception Program Counter (sepc)	56
4.1.10	Supervisor Cause Register (scause)	57
4.1.11	Supervisor Trap Value (stval) Register	57
4.1.12	Supervisor Address Translation and Protection (satp) Register	59
4.2	Supervisor Instructions	60
4.2.1	Supervisor Memory-Management Fence Instruction	61
4.3	Sv32: Page-Based 32-bit Virtual-Memory Systems	62
4.3.1	Addressing and Memory Protection	62
4.3.2	Virtual Address Translation Process	64
4.4	Sv39: Page-Based 39-bit Virtual-Memory System	65
4.4.1	Addressing and Memory Protection	66
4.5	Sv48: Page-Based 48-bit Virtual-Memory System	67
4.5.1	Addressing and Memory Protection	67
5	Hypervisor Extension, Version 0.2-draft	69
5.1	Privilege Modes	70

5.2	Hypervisor CSRs	70
5.2.1	Hypervisor Status Register (hstatus)	71
5.2.2	Hypervisor Trap Delegation Registers (hedeleg and hideleg)	71
5.2.3	Hypervisor Guest Address Translation and Protection Register (hgatp) . . .	72
5.2.4	Background Supervisor Status Register (bsstatus)	74
5.2.5	Background Supervisor Interrupt Registers (bsip and bsie)	75
5.2.6	Background Supervisor Trap Vector Base Address Register (bstvec)	75
5.2.7	Background Supervisor Scratch Register (bsscratch)	76
5.2.8	Background Supervisor Exception Program Counter (bsepc)	76
5.2.9	Background Supervisor Cause Register (bscause)	76
5.2.10	Background Supervisor Trap Value Register (bstval)	77
5.2.11	Background Supervisor Address Translation and Protection Register (bsatp)	77
5.3	Hypervisor Instructions	78
5.3.1	Hypervisor Memory-Management Fence Instructions	78
5.4	Machine-Level CSRs	79
5.4.1	Machine Status Register (mstatus)	79
5.5	Two-Level Address Translation	81
5.5.1	Guest Physical Address Translation	81
5.5.2	Memory-Management Fences	83
5.6	Base ISA Control	84
5.7	Traps	84
5.8	Trap Return	86
6	RISC-V Privileged Instruction Set Listings	87
7	Platform-Level Interrupt Controller (PLIC)	89
7.1	PLIC Overview	89
7.2	Interrupt Sources	89
7.2.1	Local Interrupt Sources	90

7.2.2	Global Interrupt Sources	91
7.3	Interrupt Targets and Hart Contexts	91
7.4	Interrupt Gateways	91
7.5	Interrupt Identifiers (IDs)	92
7.6	Interrupt Priorities	92
7.7	Interrupt Enables	93
7.8	Interrupt Priority Thresholds	93
7.9	Interrupt Notifications	94
7.10	Interrupt Claims	94
7.11	Interrupt Completion	95
7.12	Interrupt Flow	95
7.13	PLIC Core Specification	96
7.14	Controlling Access to the PLIC	96
8	Machine Configuration Description	97
8.1	Configuration String Search Procedure	97
9	History	99
9.1	Research Funding at UC Berkeley	99

Chapter 1

Introduction

This is a draft of the privileged architecture description document for RISC-V. Feedback welcome. Changes will occur before the final release.

This document describes the RISC-V privileged architecture, which covers all aspects of RISC-V systems beyond the user-level ISA, including privileged instructions as well as additional functionality required for running operating systems and attaching external devices.

Commentary on our design decisions is formatted as in this paragraph, and can be skipped if the reader is only interested in the specification itself.

We briefly note that the entire privileged-level design described in this document could be replaced with an entirely different privileged-level design without changing the user-level ISA, and possibly without even changing the ABI. In particular, this privileged specification was designed to run existing popular operating systems, and so embodies the conventional level-based protection model. Alternate privileged specifications could embody other more flexible protection-domain models.

1.1 RISC-V Privileged Software Stack Terminology

This section describes the terminology we use to describe components of the wide range of possible privileged software stacks for RISC-V.

Figure 1.1 shows some of the possible software stacks that can be supported by the RISC-V architecture. The left-hand side shows a simple system that supports only a single application running on an application execution environment (AEE). The application is coded to run with a particular application binary interface (ABI). The ABI includes the supported user-level ISA plus a set of ABI calls to interact with the AEE. The ABI hides details of the AEE from the application to allow greater flexibility in implementing the AEE. The same ABI could be implemented natively on multiple different host OSs, or could be supported by a user-mode emulation environment running on a machine with a different native ISA.

Our graphical convention represents abstract interfaces using black boxes with white text, to separate them from concrete instances of components implementing the interfaces.

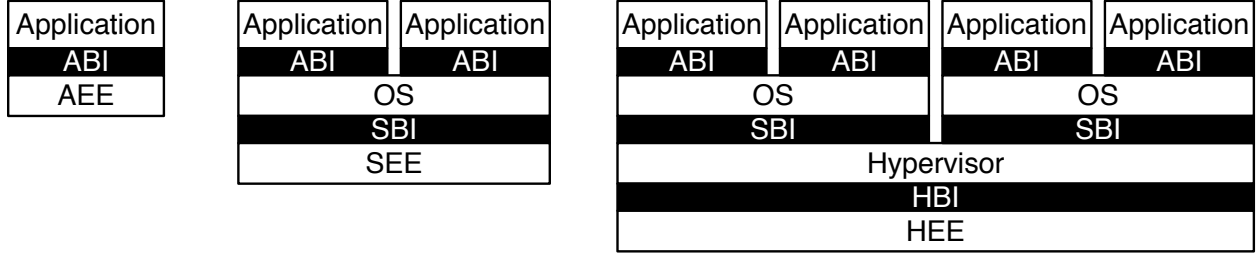


Figure 1.1: Different implementation stacks supporting various forms of privileged execution.

The middle configuration shows a conventional operating system (OS) that can support multiprogrammed execution of multiple applications. Each application communicates over an ABI with the OS, which provides the AEE. Just as applications interface with an AEE via an ABI, RISC-V operating systems interface with a supervisor execution environment (SEE) via a supervisor binary interface (SBI). An SBI comprises the user-level and supervisor-level ISA together with a set of SBI function calls. Using a single SBI across all SEE implementations allows a single OS binary image to run on any SEE. The SEE can be a simple boot loader and BIOS-style IO system in a low-end hardware platform, or a hypervisor-provided virtual machine in a high-end server, or a thin translation layer over a host operating system in an architecture simulation environment.

Most supervisor-level ISA definitions do not separate the SBI from the execution environment and/or the hardware platform, complicating virtualization and bring-up of new hardware platforms.

The rightmost configuration shows a virtual machine monitor configuration where multiple multiprogrammed OSs are supported by a single hypervisor. Each OS communicates via an SBI with the hypervisor, which provides the SEE. The hypervisor communicates with the hypervisor execution environment (HEE) using a hypervisor binary interface (HBI), to isolate the hypervisor from details of the hardware platform.

The ABI, SBI, and HBI are still a work-in-progress, but we are now prioritizing support for Type-2 hypervisors where the SBI is provided recursively by an S-mode OS.

Hardware implementations of the RISC-V ISA will generally require additional features beyond the privileged ISA to support the various execution environments (AEE, SEE, or HEE).

1.2 Privilege Levels

At any time, a RISC-V hardware thread (*hart*) is running at some privilege level encoded as a mode in one or more CSRs (control and status registers). Three RISC-V privilege levels are currently defined as shown in Table 1.1.

Privilege levels are used to provide protection between different components of the software stack, and attempts to perform operations not permitted by the current privilege mode will cause an exception to be raised. These exceptions will normally cause traps into an underlying execution environment.

Level	Encoding	Name	Abbreviation
0	00	User/Application	U
1	01	Supervisor	S
2	10	<i>Reserved</i>	
3	11	Machine	M

Table 1.1: RISC-V privilege levels.

The machine level has the highest privileges and is the only mandatory privilege level for a RISC-V hardware platform. Code run in machine-mode (M-mode) is usually inherently trusted, as it has low-level access to the machine implementation. M-mode can be used to manage secure execution environments on RISC-V. User-mode (U-mode) and supervisor-mode (S-mode) are intended for conventional application and operating system usage respectively.

The previous Hypervisor mode (H-mode) designed to support Type-1 hypervisors has been removed and the encoding space reserved as we are focusing on hypervisor support via an extended S mode suitable for both Type-1 and Type-2 hypervisors as described in Chapter 5. The encoding space for H is reserved for future use and to avoid backwards incompatible changes in bit positions in various status registers. The bit positions might be reused in the future for different Type-1 hypervisor support or possibly additional secure execution modes.

Each privilege level has a core set of privileged ISA extensions with optional extensions and variants. For example, machine-mode supports several optional standard variants for address translation and memory protection. Also, supervisor-mode can be extended to support Type-2 hypervisor execution as described in Chapter 5.

Implementations might provide anywhere from 1 to 3 privilege modes trading off reduced isolation for lower implementation cost, as shown in Table 1.2.

In the description, we try to separate the privilege level for which code is written, from the privilege mode in which it runs, although the two are often tied. For example, a supervisor-level operating system can run in supervisor-mode on a system with three privilege modes, but can also run in user-mode under a classic virtual machine monitor on systems with two or more privilege modes. In both cases, the same supervisor-level operating system binary code can be used, coded to a supervisor-level SBI and hence expecting to be able to use supervisor-level privileged instructions and CSRs. When running a guest OS in user mode, all supervisor-level actions will be trapped and emulated by the SEE running in the higher-privilege level.

Number of levels	Supported Modes	Intended Usage
1	M	Simple embedded systems
2	M, U	Secure embedded systems
3	M, S, U	Systems running Unix-like operating systems

Table 1.2: Supported combinations of privilege modes.

All hardware implementations must provide M-mode, as this is the only mode that has unfettered access to the whole machine. The simplest RISC-V implementations may provide only M-mode, though this will provide no protection against incorrect or malicious application code.

The lock feature of the optional PMP facility can provide some limited protection even with only M-mode implemented.

Many RISC-V implementations will also support at least user mode (U-mode) to protect the rest of the system from application code. Supervisor mode (S-mode) can be added to provide isolation between a supervisor-level operating system and the SEE.

A hart normally runs application code in U-mode until some trap (e.g., a supervisor call or a timer interrupt) forces a switch to a trap handler, which usually runs in a more privileged mode. The hart will then execute the trap handler, which will eventually resume execution at or after the original trapped instruction in U-mode. Traps that increase privilege level are termed *vertical* traps, while traps that remain at the same privilege level are termed *horizontal* traps. The RISC-V privileged architecture provides flexible routing of traps to different privilege layers.

Horizontal traps can be implemented as vertical traps that return control to a horizontal trap handler in the less-privileged mode.

1.3 Debug Mode

Implementations may also include a debug mode to support off-chip debugging and/or manufacturing test. Debug mode (D-mode) can be considered an additional privilege mode, with even more access than M-mode. The separate debug specification proposal describes operation of a RISC-V hart in debug mode. Debug mode reserves a few CSR addresses that are only accessible in D-mode, and may also reserve some portions of the physical memory space on a platform.

Chapter 2

Control and Status Registers (CSRs)

The SYSTEM major opcode is used to encode all privileged instructions in the RISC-V ISA. These can be divided into two main classes: those that atomically read-modify-write control and status registers (CSRs), and all other privileged instructions. In addition to the user-level state described in Volume I of this manual, an implementation may contain additional CSRs, accessible by some subset of the privilege levels using the CSR instructions described in the user-level manual. In this chapter, we map out the CSR address space. The following chapters describe the function of each of the CSRs according to privilege level, as well as the other privileged instructions which are generally closely associated with a particular privilege level. Note that although CSRs and instructions are associated with one privilege level, they are also accessible at all higher privilege levels.

2.1 CSR Address Mapping Conventions

The standard RISC-V ISA sets aside a 12-bit encoding space (`csr[11:0]`) for up to 4,096 CSRs. By convention, the upper 4 bits of the CSR address (`csr[11:8]`) are used to encode the read and write accessibility of the CSRs according to privilege level as shown in Table 2.1. The top two bits (`csr[11:10]`) indicate whether the register is read/write (00, 01, or 10) or read-only (11). The next two bits (`csr[9:8]`) encode the lowest privilege level that can access the CSR.

The CSR address convention uses the upper bits of the CSR address to encode default access privileges. This simplifies error checking in the hardware and provides a larger CSR space, but does constrain the mapping of CSRs into the address space.

Implementations might allow a more-privileged level to trap otherwise permitted CSR accesses by a less-privileged level to allow these accesses to be intercepted. This change should be transparent to the less-privileged software.

Attempts to access a non-existent CSR raise an illegal instruction exception. Attempts to access a CSR without appropriate privilege level or to write a read-only register also raise illegal instruction exceptions. A read/write register might also contain some bits that are read-only, in which case writes to the read-only bits are ignored.

CSR Address			Hex	Use and Accessibility
[11:10]	[9:8]	[7:6]		
User CSRs				
00	00	XX	0x000–0x0FF	Standard read/write
01	00	XX	0x400–0x4FF	Standard read/write
10	00	XX	0x800–0x8FF	Non-standard read/write
11	00	00–10	0xC00–0xCBF	Standard read-only
11	00	11	0xCC0–0xCFF	Non-standard read-only
Supervisor CSRs				
00	01	XX	0x100–0x1FF	Standard read/write
01	01	00–10	0x500–0x5BF	Standard read/write
01	01	11	0x5C0–0x5FF	Non-standard read/write
10	01	00–10	0x900–0x9BF	Standard read/write
10	01	11	0x9C0–0x9FF	Non-standard read/write
11	01	00–10	0xD00–0xDBF	Standard read-only
11	01	11	0xDC0–0xDFF	Non-standard read-only
Reserved CSRs				
XX	10	XX	Reserved	
Machine CSRs				
00	11	XX	0x300–0x3FF	Standard read/write
01	11	00–10	0x700–0x79F	Standard read/write
01	11	10	0x7A0–0x7AF	Standard read/write debug CSRs
01	11	10	0x7B0–0x7BF	Debug-mode-only CSRs
01	11	11	0x7C0–0x7FF	Non-standard read/write
10	11	00–10	0xB00–0xBBF	Standard read/write
10	11	11	0xBC0–0xBFF	Non-standard read/write
11	11	00–10	0xF00–0xFBF	Standard read-only
11	11	11	0xFC0–0xFFF	Non-standard read-only

Table 2.1: Allocation of RISC-V CSR address ranges.

Table 2.1 also indicates the convention to allocate CSR addresses between standard and non-standard uses. The CSR addresses reserved for non-standard uses will not be redefined by future standard extensions.

We have dropped the explicit allocation of CSR space for shadow CSRs to leave more flexibility for allocated other CSRs. Shadow CSRs can still be added in the appropriate R/W space. The counters are the only shadowed CSRs in the current spec.

Shadows CSRs provide a read-write address via which a higher privilege level can modify a register that is read-only at a lower privilege level. Note that if one privilege level has already allocated a read/write shadow address, then any higher privilege level can use the same CSR address for read/write access to the same register.

Effective virtualization requires that as many instructions run natively as possible inside a virtualized environment, while any privileged accesses trap to the virtual machine monitor [1]. CSRs that are read-only at some lower privilege level are shadowed into separate CSR addresses if they are made read-write at a higher privilege level. This avoids trapping permitted lower-privilege accesses while still causing traps on illegal accesses.

Machine-mode standard read-write CSRs 0x7A0–0x7BF are reserved for use by the debug system. Implementations should raise illegal instruction exceptions on machine-mode access to these registers.

2.2 CSR Listing

Tables 2.2–2.5 list the CSRs that have currently been allocated CSR addresses. The timers, counters, and floating-point CSRs are standard user-level CSRs, as well as the additional user trap registers added by the N extension. The other registers are used by privileged code, as described in the following chapters. Note that not all registers are required on all implementations.

Number	Privilege	Name	Description
User Trap Setup			
0x000	URW	<code>ustatus</code>	User status register.
0x004	URW	<code>uie</code>	User interrupt-enable register.
0x005	URW	<code>utvec</code>	User trap handler base address.
User Trap Handling			
0x040	URW	<code>uscratch</code>	Scratch register for user trap handlers.
0x041	URW	<code>uepc</code>	User exception program counter.
0x042	URW	<code>ucause</code>	User trap cause.
0x043	URW	<code>utval</code>	User bad address or instruction.
0x044	URW	<code>uip</code>	User interrupt pending.
User Floating-Point CSRs			
0x001	URW	<code>fflags</code>	Floating-Point Accrued Exceptions.
0x002	URW	<code>frm</code>	Floating-Point Dynamic Rounding Mode.
0x003	URW	<code>fcsr</code>	Floating-Point Control and Status Register (<code>frm</code> + <code>fflags</code>).
User Counter/Timers			
0xC00	URO	<code>cycle</code>	Cycle counter for RDCYCLE instruction.
0xC01	URO	<code>time</code>	Timer for RDTIME instruction.
0xC02	URO	<code>instret</code>	Instructions-retired counter for RDINSTRET instruction.
0xC03	URO	<code>hpmcounter3</code>	Performance-monitoring counter.
0xC04	URO	<code>hpmcounter4</code>	Performance-monitoring counter.
		\vdots	
0xC1F	URO	<code>hpmcounter31</code>	Performance-monitoring counter.
0xC80	URO	<code>cycleh</code>	Upper 32 bits of <code>cycle</code> , RV32I only.
0xC81	URO	<code>timeh</code>	Upper 32 bits of <code>time</code> , RV32I only.
0xC82	URO	<code>instreth</code>	Upper 32 bits of <code>instret</code> , RV32I only.
0xC83	URO	<code>hpmcounter3h</code>	Upper 32 bits of <code>hpmcounter3</code> , RV32I only.
0xC84	URO	<code>hpmcounter4h</code>	Upper 32 bits of <code>hpmcounter4</code> , RV32I only.
		\vdots	
0xC9F	URO	<code>hpmcounter31h</code>	Upper 32 bits of <code>hpmcounter31</code> , RV32I only.

Table 2.2: Currently allocated RISC-V user-level CSR addresses.

Number	Privilege	Name	Description
Supervisor Trap Setup			
0x100	SRW	sstatus	Supervisor status register.
0x102	SRW	sedeleg	Supervisor exception delegation register.
0x103	SRW	sideleg	Supervisor interrupt delegation register.
0x104	SRW	sie	Supervisor interrupt-enable register.
0x105	SRW	stvec	Supervisor trap handler base address.
0x106	SRW	scounteren	Supervisor counter enable.
Supervisor Trap Handling			
0x140	SRW	sscratch	Scratch register for supervisor trap handlers.
0x141	SRW	sepc	Supervisor exception program counter.
0x142	SRW	scause	Supervisor trap cause.
0x143	SRW	stval	Supervisor bad address or instruction.
0x144	SRW	sip	Supervisor interrupt pending.
Supervisor Protection and Translation			
0x180	SRW	satp	Supervisor address translation and protection.

Table 2.3: Currently allocated RISC-V supervisor-level CSR addresses.

Number	Privilege	Name	Description
Machine Information Registers			
0xF11	MRO	<code>mvendorid</code>	Vendor ID.
0xF12	MRO	<code>marchid</code>	Architecture ID.
0xF13	MRO	<code>mimpid</code>	Implementation ID.
0xF14	MRO	<code>mhartid</code>	Hardware thread ID.
Machine Trap Setup			
0x300	MRW	<code>mstatus</code>	Machine status register.
0x301	MRW	<code>misa</code>	ISA and extensions
0x302	MRW	<code>medeleg</code>	Machine exception delegation register.
0x303	MRW	<code>mideleg</code>	Machine interrupt delegation register.
0x304	MRW	<code>mie</code>	Machine interrupt-enable register.
0x305	MRW	<code>mtvec</code>	Machine trap-handler base address.
0x306	MRW	<code>mcOUNTERen</code>	Machine counter enable.
Machine Trap Handling			
0x340	MRW	<code>mscratch</code>	Scratch register for machine trap handlers.
0x341	MRW	<code>mepc</code>	Machine exception program counter.
0x342	MRW	<code>mcause</code>	Machine trap cause.
0x343	MRW	<code>mtval</code>	Machine bad address or instruction.
0x344	MRW	<code>mip</code>	Machine interrupt pending.
Machine Protection and Translation			
0x3A0	MRW	<code>pmpcfg0</code>	Physical memory protection configuration.
0x3A1	MRW	<code>pmpcfg1</code>	Physical memory protection configuration, RV32 only.
0x3A2	MRW	<code>pmpcfg2</code>	Physical memory protection configuration.
0x3A3	MRW	<code>pmpcfg3</code>	Physical memory protection configuration, RV32 only.
0x3B0	MRW	<code>pmpaddr0</code>	Physical memory protection address register.
0x3B1	MRW	<code>pmpaddr1</code>	Physical memory protection address register.
		<code>⋮</code>	
0x3BF	MRW	<code>pmpaddr15</code>	Physical memory protection address register.

Table 2.4: Currently allocated RISC-V machine-level CSR addresses.

Number	Privilege	Name	Description
Machine Counter/Timers			
0xB00	MRW	<code>mcycle</code>	Machine cycle counter.
0xB02	MRW	<code>minstret</code>	Machine instructions-retired counter.
0xB03	MRW	<code>mhpmcounter3</code>	Machine performance-monitoring counter.
0xB04	MRW	<code>mhpmcounter4</code>	Machine performance-monitoring counter.
		\vdots	
0xB1F	MRW	<code>mhpmcounter31</code>	Machine performance-monitoring counter.
0xB80	MRW	<code>mcycleh</code>	Upper 32 bits of <code>mcycle</code> , RV32I only.
0xB82	MRW	<code>minstreth</code>	Upper 32 bits of <code>minstret</code> , RV32I only.
0xB83	MRW	<code>mhpmcounter3h</code>	Upper 32 bits of <code>mhpmcounter3</code> , RV32I only.
0xB84	MRW	<code>mhpmcounter4h</code>	Upper 32 bits of <code>mhpmcounter4</code> , RV32I only.
		\vdots	
0xB9F	MRW	<code>mhpmcounter31h</code>	Upper 32 bits of <code>mhpmcounter31</code> , RV32I only.
Machine Counter Setup			
0x323	MRW	<code>mhpmevent3</code>	Machine performance-monitoring event selector.
0x324	MRW	<code>mhpmevent4</code>	Machine performance-monitoring event selector.
		\vdots	
0x33F	MRW	<code>mhpmevent31</code>	Machine performance-monitoring event selector.
Debug/Trace Registers (shared with Debug Mode)			
0x7A0	MRW	<code>tselect</code>	Debug/Trace trigger register select.
0x7A1	MRW	<code>tdata1</code>	First Debug/Trace trigger data register.
0x7A2	MRW	<code>tdata2</code>	Second Debug/Trace trigger data register.
0x7A3	MRW	<code>tdata3</code>	Third Debug/Trace trigger data register.
Debug Mode Registers			
0x7B0	DRW	<code>dcsr</code>	Debug control and status register.
0x7B1	DRW	<code>dpc</code>	Debug PC.
0x7B2	DRW	<code>dscratch</code>	Debug scratch register.

Table 2.5: Currently allocated RISC-V machine-level CSR addresses.

2.3 CSR Field Specifications

The following definitions and abbreviations are used in specifying the behavior of fields within the CSRs.

Reserved Writes Preserve Values, Reads Ignore Values (WPRI)

Some whole read/write fields are reserved for future use. Software should ignore the values read from these fields, and should preserve the values held in these fields when writing values to other fields of the same register. For forward compatibility, implementations that do not furnish these fields must hardwire them to zero. These fields are labeled **WPRI** in the register descriptions.

To simplify the software model, any backward-compatible future definition of previously reserved fields within a CSR must cope with the possibility that a non-atomic read/modify/write sequence is used to update other fields in the CSR. Alternatively, the original CSR definition must specify that subfields can only be updated atomically, which may require a two-instruction clear bit/set bit sequence in general that can be problematic if intermediate values are not legal.

Write/Read Only Legal Values (WLRL)

Some read/write CSR fields specify behavior for only a subset of possible bit encodings, with other bit encodings reserved. Software should not write anything other than legal values to such a field, and should not assume a read will return a legal value unless the last write was of a legal value, or the register has not been written since another operation (e.g., reset) set the register to a legal value. These fields are labeled **WLRL** in the register descriptions.

Hardware implementations need only implement enough state bits to differentiate between the supported values, but must always return the complete specified bit-encoding of any supported value when read.

Implementations are permitted but not required to raise an illegal instruction exception if an instruction attempts to write a non-supported value to a CSR field. Hardware implementations can return arbitrary bit patterns on the read of a CSR field when the last write was of an illegal value, but the value returned should deterministically depend on the illegal written value.

Write Any Values, Reads Legal Values (WARL)

Some read/write CSR fields are only defined for a subset of bit encodings, but allow any value to be written while guaranteeing to return a legal value whenever read. Assuming that writing the CSR has no other side effects, the range of supported values can be determined by attempting to write a desired setting then reading to see if the value was retained. These fields are labeled **WARL** in the register descriptions.

Implementations will not raise an exception on writes of unsupported values to an **WARL** field. Implementations must always deterministically return the same legal value after a given illegal value is written.

2.4 CSR Width Modulation

If the width of a CSR is changed (for example, by changing MXLEN or UXLEN, as described in Section 3.1.8), the values of the *writable* fields and bits of the new-width CSR are, unless specified otherwise, determined from the previous-width CSR as though by this algorithm:

1. The value of the previous-width CSR is copied to a temporary register of the same width.
2. For the read-only bits of the previous-width CSR, the bits at the same positions in the temporary register are set to zeros.
3. The width of the temporary register is changed to the new width. If the new width W is narrower than the previous width, the least-significant W bits of the temporary register are retained and the more-significant bits are discarded. If the new width is wider than the previous width, the temporary register is zero-extended to the wider width.
4. Each writable field of the new-width CSR takes the value of the bits at the same positions in the temporary register.

Changing the width of a CSR is not a read or write of the CSR and thus does not trigger any side effects.

Chapter 3

Machine-Level ISA, version 1.11-draft

This chapter describes the machine-level operations available in machine-mode (M-mode), which is the highest privilege mode in a RISC-V system. M-mode is the only mandatory privilege mode in a RISC-V hardware implementation. M-mode is used for low-level access to a hardware platform and is the first mode entered at reset. M-mode can also be used to implement features that are too difficult or expensive to implement in hardware directly. The RISC-V machine-level ISA contains a common core that is extended depending on which other privilege levels are supported and other details of the hardware implementation.

3.1 Machine-Level CSRs

In addition to the machine-level CSRs described in this section, M-mode code can access all CSRs at lower privilege levels.

3.1.1 Machine ISA Register `misa`

The `misa` CSR is a **WARL** read-write register reporting the ISA supported by the hart. This register must be readable in any implementation, but a value of zero can be returned to indicate the `misa` register has not been implemented, requiring that CPU capabilities be determined through a separate non-standard mechanism.

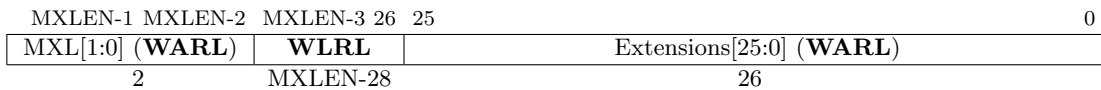


Figure 3.1: Machine ISA register (`misa`).

The MXL (Machine XLEN) field encodes the native base integer ISA width as shown in Table 3.1. The MXL field may be writable in implementations that support multiple base ISA widths. The effective XLEN in M-mode, *MXLEN*, is given by the setting of MXL, or has a fixed value if `misa` is zero. The MXL field is always set to the widest supported ISA variant at reset.

MXL	XLEN
1	32
2	64
3	128

Table 3.1: Encoding of MXL field in `misa`

The `misa` CSR is MXLEN bits wide. If the value read from `misa` is nonzero, field MXL of that value always denotes the current MXLEN. If a write to `misa` causes MXLEN to change, the position of MXL moves to the most-significant two bits of `misa` at the new width.

The base width can be quickly ascertained using branches on the sign of the returned `misa` value, and possibly a shift left by one and a second branch on the sign. These checks can be written in assembly code without knowing the register width (XLEN) of the machine. The base width is given by $XLEN = 2^{MXL+4}$.

The base width can also be found if `misa` is zero, by placing the immediate 4 in a register then shifting the register left by 31 bits at a time. If zero after one shift, then the machine is RV32. If zero after two shifts, then the machine is RV64, else RV128.

When MXL is set to a value less than the widest supported XLEN, all operations must ignore source operand register bits above the configured XLEN, and must sign-extend results to fill the entire widest supported XLEN in the destination register.

We require that operations always fill the entire underlying hardware registers with defined values to avoid implementation-defined behavior.

The Extensions field encodes the presence of the standard extensions, with a single bit per letter of the alphabet (bit 0 encodes presence of extension “A”, bit 1 encodes presence of extension “B”, through to bit 25 which encodes “Z”). The “I” bit will be set for RV32I, RV64I, RV128I base ISAs, and the “E” bit will be set for RV32E. The Extension is a **WARL** field that can contain writable bits where the implementation allows the supported ISA to be modified. At reset, the Extension field should contain the maximal set of supported extensions, and I should be selected over E if both are available.

The “G” bit is used as an escape to allow expansion to a larger space of standard extension names.

G is used to indicate the combination IMAFD, so is redundant in the `misa` CSR, hence we reserve the bit to indicate that additional standard extensions are present.

The “U” and “S” bits will be set if there is support for user and supervisor modes respectively.

The “X” bit will be set if there are any non-standard extensions.

The `misa` CSR exposes a rudimentary catalog of CPU features to machine-mode code. More extensive information can be obtained in machine mode by probing other machine registers, and examining other ROM storage in the system as part of the boot process.

We require that lower privilege levels execute environment calls instead of reading CPU registers to determine features available at each privilege level. This enables virtualization layers to alter the ISA observed at any level, and supports a much richer command interface without burdening hardware designs.

Bit	Character	Description
0	A	Atomic extension
1	B	<i>Tentatively reserved for Bit operations extension</i>
2	C	Compressed extension
3	D	Double-precision floating-point extension
4	E	RV32E base ISA
5	F	Single-precision floating-point extension
6	G	Additional standard extensions present
7	H	Hypervisor extension
8	I	RV32I/64I/128I base ISA
9	J	<i>Tentatively reserved for Dynamically Translated Languages extension</i>
10	K	<i>Reserved</i>
11	L	<i>Tentatively reserved for Decimal Floating-Point extension</i>
12	M	Integer Multiply/Divide extension
13	N	User-level interrupts supported
14	O	<i>Reserved</i>
15	P	<i>Tentatively reserved for Packed-SIMD extension</i>
16	Q	Quad-precision floating-point extension
17	R	<i>Reserved</i>
18	S	Supervisor mode implemented
19	T	<i>Tentatively reserved for Transactional Memory extension</i>
20	U	User mode implemented
21	V	<i>Tentatively reserved for Vector extension</i>
22	W	<i>Reserved</i>
23	X	Non-standard extensions present
24	Y	<i>Reserved</i>
25	Z	<i>Reserved</i>

Table 3.2: Encoding of Extensions field in `misr`. All bits that are reserved for future use must return zero when read.

Writing `misalr` may increase IALIGN, e.g., by disabling the “C” extension. If an instruction that would write `misalr` increases IALIGN, and the subsequent instruction’s address is not IALIGN-bit aligned, the write to `misalr` is suppressed, leaving `misalr` unchanged.

3.1.2 Machine Vendor ID Register `mvendorid`

The `mvendorid` CSR is a 32-bit read-only register providing the JEDEC manufacturer ID of the provider of the core. This register must be readable in any implementation, but a value of 0 can be returned to indicate the field is not implemented or that this is a non-commercial implementation.

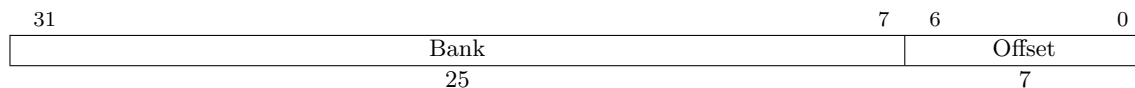


Figure 3.2: Vendor ID register (`mvendorid`).

JEDEC manufacturer IDs are ordinarily encoded as a sequence of one-byte continuation codes `0x7f`, terminated by a one-byte ID not equal to `0x7f`, with an odd parity bit in the most-significant bit of each byte. `mvendorid` encodes the number of one-byte continuation codes in the Bank field, and encodes the final byte in the Offset field, discarding the parity bit. For example, the JEDEC manufacturer ID `0x7f 0x7f 0x7f 0x7f 0x7f 0x7f 0x7f 0x7f 0x7f 0x7f 0x7f 0x7f 0x8a` (twelve continuation codes followed by `0x8a`) would be encoded in the `mvendorid` field as `0x60a`.

Previously the vendor ID was to be a number allocated by the RISC-V Foundation, but this duplicates the work of JEDEC in maintaining a manufacturer ID standard. At time of writing, registering a manufacturer ID with JEDEC has a one-time cost of \$500.

3.1.3 Machine Architecture ID Register `marchid`

The `marchid` CSR is an MXLEN-bit read-only register encoding the base microarchitecture of the hart. This register must be readable in any implementation, but a value of 0 can be returned to indicate the field is not implemented. The combination of `mvendorid` and `marchid` should uniquely identify the type of hart microarchitecture that is implemented.



Figure 3.3: Machine Architecture ID register (`marchid`).

Open-source project architecture IDs are allocated globally by the RISC-V Foundation, and have non-zero architecture IDs with a zero most-significant-bit (MSB). Commercial architecture IDs are allocated by each commercial vendor independently, but must have the MSB set and cannot contain zero in the remaining MXLEN-1 bits.

The intent is for the architecture ID to represent the microarchitecture associated with the repo around which development occurs rather than a particular organization. Commercial fabrications of open-source designs should (and might be required by the license to) retain the original architecture ID. This will aid in reducing fragmentation and tool support costs, as well as provide attribution. Open-source architecture IDs should be administered by the Foundation and should only be allocated to released, functioning open-source projects. Commercial architecture IDs can be managed independently by any registered vendor but are required to have IDs disjoint from

the open-source architecture IDs (MSB set) to prevent collisions if a vendor wishes to use both closed-source and open-source microarchitectures.

The convention adopted within the following Implementation field can be used to segregate branches of the same architecture design, including by organization. The `misra` register also helps distinguish different variants of a design, as does the configuration string if present.

3.1.4 Machine Implementation ID Register `mimpid`

The `mimpid` CSR provides a unique encoding of the version of the processor implementation. This register must be readable in any implementation, but a value of 0 can be returned to indicate that the field is not implemented. The Implementation value should reflect the design of the RISC-V processor itself and not any surrounding system.



Figure 3.4: Machine Implementation ID register (`mimpid`).

The format of this field is left to the provider of the architecture source code, but will be often be printed by standard tools as a hexadecimal string without any leading or trailing zeros, so the Implementation value can be left-justified (i.e., filled in from most-significant nibble down) with subfields aligned on nibble boundaries to ease human readability.

3.1.5 Hart ID Register `mhartid`

The `mhartid` CSR is an MXLEN-bit read-only register containing the integer ID of the hardware thread running the code. This register must be readable in any implementation. Hart IDs might not necessarily be numbered contiguously in a multiprocessor system, but at least one hart must have a hart ID of zero. Hart IDs must be unique.



Figure 3.5: Hart ID register (`mhartid`).

In certain cases, we must ensure exactly one hart runs some code (e.g., at reset), and so require one hart to have a known hart ID of zero.

For efficiency, system implementers should aim to reduce the magnitude of the largest hart ID used in a system.

3.1.6 Machine Status Register (`mstatus`)

The `mstatus` register is an MXLEN-bit read/write register formatted as shown in Figure 3.6 for RV32 and Figure 3.7 for RV64 and RV128. The `mstatus` register keeps track of and controls the hart's current operating state. Restricted views of the `mstatus` register appear as the `sstatus` and `ustatus` registers in the S-level and U-level ISAs respectively.

31	30							23	22	21	20	19	18	17		
SD	WPRI							TSR	TW	TVM	MXR	SUM	MPRV			
1	8							1	1	1	1	1	1	1		
16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
XS[1:0]	FS[1:0]	MPP[1:0]	WPRI	SPP	MPIE	WPRI	SPIE	UPIE	MIE	WPRI	SIE	UIE				
2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1

Figure 3.6: Machine-mode status register (`mstatus`) for RV32.

MXLEN-1	MXLEN-2	36	35	34	33	32	31	23	22	21	20	19	18	17
SD	WPRI		SXL[1:0]	UXL[1:0]	WPRI		TSR	TW	TVM	MXR	SUM	MPRV		
1	MXLEN-37		2	2	9		1	1	1	1	1	1		

16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
XS[1:0]	FS[1:0]	MPP[1:0]	WPRI	SPP	MPIE	WPRI	SPIE	UPIE	MIE	WPRI	SIE	UIE				
2	2	2	2	1	1	1	1	1	1	1	1	1				

Figure 3.7: Machine-mode status register (`mstatus`) for RV64 and RV128.

3.1.7 Privilege and Global Interrupt-Enable Stack in `mstatus` register

Global interrupt-enable bits, MIE, SIE, and UIE, are provided for each privilege mode. These bits are primarily used to guarantee atomicity with respect to interrupt handlers in the current privilege mode.

The global xIE bits are located in the low-order bits of `mstatus`, allowing them to be atomically set or cleared with a single CSR instruction.

When a hart is executing in privilege mode x , interrupts are globally enabled when $xIE=1$ and globally disabled when $xIE=0$. Interrupts for lower-privilege modes, $w < x$, are always globally disabled regardless of the setting of the lower-privilege mode's global wIE bit. Interrupts for higher-privilege modes, $y > x$, are always globally enabled regardless of the setting of the higher-privilege mode's global yIE bit. Higher-privilege-level code can use separate per-interrupt enable bits to disable selected high-privilege-mode interrupts before ceding control to a lower-privilege mode.

A higher-privilege mode y could disable all of its interrupts before ceding control to a lower-privilege mode but this would be unusual as it would leave only a synchronous trap, non-maskable interrupt, or reset as means to regain control of the hart.

To support nested traps, each privilege mode x has a two-level stack of interrupt-enable bits and privilege modes. $xPIE$ holds the value of the interrupt-enable bit active prior to the trap, and xPP holds the previous privilege mode. The xPP fields can only hold privilege modes up to x , so MPP is two bits wide, SPP is one bit wide, and UPP is implicitly zero. When a trap is taken from privilege mode y into privilege mode x , $xPIE$ is set to the value of xIE ; xIE is set to 0; and xPP is set to y .

For lower privilege modes, any trap (synchronous or asynchronous) is usually taken at a higher

privilege mode with interrupts disabled upon entry. The higher-level trap handler will either service the trap and return using the stacked information, or, if not returning immediately to the interrupted context, will save the privilege stack before re-enabling interrupts, so only one entry per stack is required.

The MRET, SRET, or URET instructions are used to return from traps in M-mode, S-mode, or U-mode respectively. When executing an x RET instruction, supposing x PP holds the value y , x IE is set to x PIE; the privilege mode is changed to y ; x PIE is set to 1; and x PP is set to U (or M if user-mode is not supported).

x PP fields are **WARL** fields that can hold only privilege mode x and any implemented privilege mode lower than x . If privilege mode x is not implemented, then x PP must be hardwired to 0.

M-mode software can determine whether a privilege mode is implemented by writing that mode to MPP then reading it back.

If the machine provides only U and M modes, then only a single hardware storage bit is required to represent either 00 or 11 in MPP.

User-level interrupts are an optional extension and have been allocated the ISA extension letter N. If user-level interrupts are omitted, the UIE and UPIE bits are hardwired to zero. For all other supported privilege modes x , the x IE and x PIE must not be hardwired.

User-level interrupts are primarily intended to support secure embedded systems with only M-mode and U-mode present, but can also be supported in systems running Unix-like operating systems to support user-level trap handling.

*Fields that were previously allocated for H-mode support in mstatus have now been reserved as **WPRI** fields. To reduce backwards incompatibility with existing implementations, we did not compact the register after removing these fields.*

3.1.8 Base ISA Control in mstatus Register

For RV64 and RV128 systems, the SXL and UXL fields are **WARL** fields that control the value of XLEN for S-mode and U-mode, respectively. The encoding of these fields is the same as the MXL field of misa, shown in Table 3.1. The effective XLEN in S-mode and U-mode are termed *SXLEN* and *UXLEN*, respectively.

For RV32 systems, the SXL and UXL fields do not exist, and SXLEN=32 and UXLEN=32.

For RV64 and RV128 systems, if S-mode is not supported, then SXL is hardwired to zero. Otherwise, it is a **WARL** field that encodes the current value of SXLEN. In particular, the implementation may hardwire SXL so that SXLEN=MXLEN.

For RV64 and RV128 systems, if U-mode is not supported, then UXL is hardwired to zero. Otherwise, it is a **WARL** field that encodes the current value of UXLEN. In particular, the implementation may hardwire UXL so that UXLEN=MXLEN or UXLEN=SXLEN.

Whenever XLEN in any mode is set to a value less than the widest supported XLEN, all operations must ignore source operand register bits above the configured XLEN, and must sign-extend results to fill the entire widest supported XLEN in the destination register.

To reduce hardware complexity, the architecture imposes no checks that lower-privilege modes have XLEN settings less than or equal to the next-higher privilege mode. In practice, such settings would almost always be an error, but machine operation is well-defined even in this case.

If MXLEN is changed from 32 to a wider width, each of `mstatus` fields SXL and UXL, if not hardwired to a forced value, gets the value corresponding to the widest supported width not wider than the new MXLEN.

3.1.9 Memory Privilege in `mstatus` Register

The MPRV (Modify PRiVilege) bit modifies the privilege level at which loads and stores execute in all privilege modes. When MPRV=0, translation and protection behave as normal. When MPRV=1, load and store memory addresses are translated and protected as though the current privilege mode were set to MPP. Instruction address-translation and protection are unaffected. MPRV is hardwired to 0 if U-mode is not supported.

The MXR (Make eXecutable Readable) bit modifies the privilege with which loads access virtual memory. When MXR=0, only loads from pages marked readable (R=1 in Figure 4.15) will succeed. When MXR=1, loads from pages marked either readable or executable (R=1 or X=1) will succeed. MXR has no effect when page-based virtual memory is not in effect. MXR is hardwired to 0 if S-mode is not supported.

The MPRV and MXR mechanisms were conceived to improve the efficiency of M-mode routines that emulate missing hardware features, e.g., misaligned loads and stores. MPRV obviates the need to perform address translation in software. MXR allows instruction words to be loaded from pages marked execute-only.

For simplicity, MPRV and MXR are in effect regardless of privilege mode, but in normal use will only be enabled for short sequences in machine mode.

The SUM (permit Supervisor User Memory access) bit modifies the privilege with which S-mode loads and stores access virtual memory. When SUM=0, S-mode memory accesses to pages that are accessible by U-mode (U=1 in Figure 4.15) will fault. When SUM=1, these accesses are permitted. SUM has no effect when page-based virtual memory is not in effect. Note that, while SUM is ordinarily ignored when not executing in S-mode, it *is* in effect when MPRV=1 and MPP=S. SUM is hardwired to 0 if S-mode is not supported.

3.1.10 Virtualization Support in `mstatus` Register

The TVM (Trap Virtual Memory) bit supports intercepting supervisor virtual-memory management operations. When TVM=1, attempts to read or write the `satp` CSR or execute the `SFENCE.VMA` instruction while executing in S-mode will raise an illegal instruction exception. When TVM=0, these operations are permitted in S-mode. TVM is hard-wired to 0 when S-mode is not supported.

The TVM mechanism improves virtualization efficiency by permitting guest operating systems to

execute in S-mode, rather than classically virtualizing them in U-mode. This approach obviates the need to trap accesses to most S-mode CSRs.

Trapping `satp` accesses and the `SFENCE.VMA` instruction provides the hooks necessary to lazily populate shadow page tables.

The TW (Timeout Wait) bit supports intercepting the WFI instruction (see Section 3.2.3). When TW=0, the WFI instruction is permitted in S-mode. When TW=1, if WFI is executed in S-mode, and it does not complete within an implementation-specific, bounded time limit, the WFI instruction causes an illegal instruction trap. The time limit may always be 0, in which case WFI always causes an illegal instruction trap in S-mode when TW=1. TW is hard-wired to 0 when S-mode is not supported.

Trapping the WFI instruction can trigger a world switch to another guest OS, rather than wastefully idling in the current guest.

The TSR (Trap SRET) bit supports intercepting the supervisor exception return instruction, SRET. When TSR=1, attempts to execute SRET while executing in S-mode will raise an illegal instruction exception. When TSR=0, this operation is permitted in S-mode. TSR is hard-wired to 0 when S-mode is not supported.

Trapping SRET is necessary to emulate the Augmented Virtualization mechanism (see Chapter 5) on implementations that do not provide it.

3.1.11 Extension Context Status in `mstatus` Register

Supporting substantial extensions is one of the primary goals of RISC-V, and hence we define a standard interface to allow unchanged privileged-mode code, particularly a supervisor-level OS, to support arbitrary user-mode state extensions.

To date, the V extension is the only standard extension that defines additional state beyond the floating-point CSR and data registers.

The FS[1:0] read/write field and the XS[1:0] read-only field are used to reduce the cost of context save and restore by setting and tracking the current state of the floating-point unit and any other user-mode extensions respectively. The FS field encodes the status of the floating-point unit, including the CSR `fcsr` and floating-point data registers `f0–f31`, while the XS field encodes the status of additional user-mode extensions and associated state. These fields can be checked by a context switch routine to quickly determine whether a state save or restore is required. If a save or restore is required, additional instructions and CSRs are typically required to effect and optimize the process.

The design anticipates that most context switches will not need to save/restore state in either or both of the floating-point unit or other extensions, so provides a fast check via the SD bit.

The FS and XS fields use the same status encoding as shown in Table 3.3, with the four possible status values being Off, Initial, Clean, and Dirty.

In systems that do not implement S-mode and do not have a floating-point unit, the FS field is hardwired to zero.

Status	FS Meaning	XS Meaning
0	Off	All off
1	Initial	None dirty or clean, some on
2	Clean	None dirty, some clean
3	Dirty	Some dirty

Table 3.3: Encoding of FS[1:0] and XS[1:0] status fields.

In systems without additional user extensions requiring new state, the XS field is hardwired to zero. Every additional extension with state adds a local field to `mstatus` encoding the equivalent of the XS states. The XS field represents a summary of all extensions' status as shown in Table 3.3.

The XS field effectively reports the maximum status value across all user-extension status fields, though individual extensions can use a different encoding than XS.

The SD bit is a read-only bit that summarizes whether either the FS field or XS field signals the presence of some dirty state that will require saving extended user context to memory. If both XS and FS are hardwired to zero, then SD is also always zero.

When an extension's status is set to Off, any instruction that attempts to read or write the corresponding state will cause an exception. When the status is Initial, the corresponding state should have an initial constant value. When the status is Clean, the corresponding state is potentially different from the initial value, but matches the last value stored on a context swap. When the status is Dirty, the corresponding state has potentially been modified since the last context save.

During a context save, the responsible privileged code need only write out the corresponding state if its status is Dirty, and can then reset the extension's status to Clean. During a context restore, the context need only be loaded from memory if the status is Clean (it should never be Dirty at restore). If the status is Initial, the context must be set to an initial constant value on context restore to avoid a security hole, but this can be done without accessing memory. For example, the floating-point registers can all be initialized to the immediate value 0.

The FS and XS fields are read by the privileged code before saving the context. The FS field is set directly by privileged code when resuming a user context, while the XS field is set indirectly by writing to the status register of the individual extensions. The status fields will also be updated during execution of instructions, regardless of privilege mode.

Extensions to the user-mode ISA often include additional user-mode state, and this state can be considerably larger than the base integer registers. The extensions might only be used for some applications, or might only be needed for short phases within a single application. To improve performance, the user-mode extension can define additional instructions to allow user-mode software to return the unit to an initial state or even to turn off the unit.

For example, a coprocessor might require to be configured before use and can be "unconfigured" after use. The unconfigured state would be represented as the Initial state for context save. If the same application remains running between the unconfigure and the next configure (which would set status to Dirty), there is no need to actually reinitialize the state at the unconfigure instruction, as all state is local to the user process, i.e., the Initial state may only cause the coprocessor state to be initialized to a constant value at context restore, not at every unconfigure.

Executing a user-mode instruction to disable a unit and place it into the Off state will cause an illegal instruction exception to be raised if any subsequent instruction tries to use the unit before it is turned back on. A user-mode instruction to turn a unit on must also ensure the unit's state is properly initialized, as the unit might have been used by another context meantime.

Changing the setting of FS has no effect on the contents of the floating-point register state. In particular, setting FS=Off does not destroy the state, nor does setting FS=Initial clear the contents. Other extensions might not preserve state when set to Off.

Implementations may choose to track the dirtiness of the floating-point register state imprecisely by reporting the state to be dirty even when it has not been modified. On some implementations, some instructions that do not mutate the floating-point state may cause the state to transition from Initial or Clean to Dirty. On other implementations, dirtiness might not be tracked at all, in which case the valid FS states are Off and Dirty, and an attempt to set FS to Initial or Clean causes it to be set to Dirty.

This definition of FS does not disallow setting FS to Dirty as a result of errant speculation. Some platforms may choose to disallow speculatively writing FS to close a potential side channel.

Table 3.4 shows all the possible state transitions for the FS or XS status bits. Note that the standard floating-point extensions do not support user-mode unconfigure or disable/enable instructions.

Current State Action	Off	Initial	Clean	Dirty
At context save in privileged code				
Save state?	No	No	No	Yes
Next state	Off	Initial	Clean	Clean
At context restore in privileged code				
Restore state?	No	Yes, to initial	Yes, from memory	N/A
Next state	Off	Initial	Clean	N/A
Execute instruction to read state				
Action?	Exception	Execute	Execute	Execute
Next state	Off	Initial	Clean	Dirty
Execute instruction to modify state, including configuration				
Action?	Exception	Execute	Execute	Execute
Next state	Off	Dirty	Dirty	Dirty
Execute instruction to unconfigure unit				
Action?	Exception	Execute	Execute	Execute
Next state	Off	Initial	Initial	Initial
Execute instruction to disable unit				
Action?	Execute	Execute	Execute	Execute
Next state	Off	Off	Off	Off
Execute instruction to enable unit				
Action?	Execute	Execute	Execute	Execute
Next state	Initial	Initial	Initial	Initial

Table 3.4: FS and XS state transitions.

Standard privileged instructions to initialize, save, and restore extension state are provided to insulate privileged code from details of the added extension state by treating the state as an opaque object.

Many coprocessor extensions are only used in limited contexts that allows software to safely unconfigure or even disable units when done. This reduces the context-switch overhead of large stateful coprocessors.

We separate out floating-point state from other extension state, as when a floating-point unit is present the floating-point registers are part of the standard calling convention, and so user-mode software cannot know when it is safe to disable the floating-point unit.

The XS field provides a summary of all added extension state, but additional microarchitectural bits might be maintained in the extension to further reduce context save and restore overhead.

The SD bit is read-only and is set when either the FS or XS bits encode a Dirty state (i.e., $SD = ((FS == 11) \text{ OR } (XS == 11))$). This allows privileged code to quickly determine when no additional context save is required beyond the integer register set and PC.

The floating-point unit state is always initialized, saved, and restored using standard instructions (F, D, and/or Q), and privileged code must be aware of FLEN to determine the appropriate space to reserve for each **f** register.

In a supervisor-level OS, any additional user-mode state should be initialized, saved, and restored using SBI calls that treats the additional context as an opaque object of a fixed maximum size. The implementation of the SBI initialize, save, and restore calls might require additional implementation-dependent privileged instructions to initialize, save, and restore microarchitectural state inside a coprocessor.

All privileged modes share a single copy of the FS and XS bits. In a system with more than one privileged mode, supervisor mode would normally use the FS and XS bits directly to record the status with respect to the supervisor-level saved context. Other more-privileged active modes must be more conservative in saving and restoring the extension state in their corresponding version of the context.

In any reasonable use case, the number of context switches between user and supervisor level should far outweigh the number of context switches to other privilege levels. Note that coprocessors should not require their context to be saved and restored to service asynchronous interrupts, unless the interrupt results in a user-level context swap.

3.1.12 Machine Trap-Vector Base-Address Register (**mtvec**)

The **mtvec** register is an MXLEN-bit read/write register that holds trap vector configuration, consisting of a vector base address (BASE) and a vector mode (MODE).

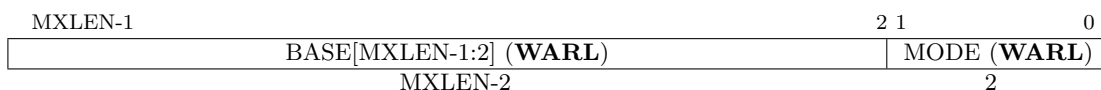


Figure 3.8: Machine trap-vector base-address register (**mtvec**).

The `mtvec` register must always be implemented, but can contain a hardwired read-only value. If `mtvec` is writable, the set of values the register may hold can vary by implementation. The value in the BASE field must always be aligned on a 4-byte boundary, and the MODE setting may impose additional alignment constraints on the value in the BASE field.

We allow for considerable flexibility in implementation of the trap vector base address. On the one hand, we do not wish to burden low-end implementations with a large number of state bits, but on the other hand, we wish to allow flexibility for larger systems.

Value	Name	Description
0	Direct	All exceptions set <code>pc</code> to BASE.
1	Vectored	Asynchronous interrupts set <code>pc</code> to BASE+4×cause.
≥2	—	<i>Reserved</i>

Table 3.5: Encoding of `mtvec` MODE field.

The encoding of the MODE field is shown in Table 3.5. When MODE=Direct, all traps into machine mode cause the `pc` to be set to the address in the BASE field. When MODE=Vectored, all synchronous exceptions into machine mode cause the `pc` to be set to the address in the BASE field, whereas interrupts cause the `pc` to be set to the address in the BASE field plus four times the interrupt cause number. For example, a machine-mode timer interrupt (see Table 3.6) causes the `pc` to be set to BASE+0x1c.

When vectored interrupts are enabled, interrupt cause 0, which corresponds to user-mode software interrupts, are vectored to the same location as synchronous exceptions. This ambiguity does not arise in practice, since user-mode software interrupts are either disabled or delegated to a less-privileged mode.

An implementation may have different alignment constraints for different modes. In particular, MODE=Vectored may have stricter alignment constraints than MODE=Direct.

Allowing coarser alignments in Vectored mode enables vectoring to be implemented without a hardware adder circuit.

Reset and NMI vector locations are given in a platform specification.

3.1.13 Machine Trap Delegation Registers (`medeleg` and `mideleg`)

By default, all traps at any privilege level are handled in machine mode, though a machine-mode handler can redirect traps back to the appropriate level with the MRET instruction (Section 3.2.2). To increase performance, implementations can provide individual read/write bits within `medeleg` and `mideleg` to indicate that certain exceptions and interrupts should be processed directly by a lower privilege level. The machine exception delegation register (`medeleg`) and machine interrupt delegation register (`mideleg`) are MXLEN-bit read/write registers.

In systems with all three privilege modes (M/S/U), setting a bit in `medeleg` or `mideleg` will delegate the corresponding trap in S-mode or U-mode to the S-mode trap handler. If U-mode traps

are supported, S-mode may in turn set corresponding bits in the `sdeleg` and `sideleg` registers to delegate traps that occur in U-mode to the U-mode trap handler.

In systems with two privilege modes (M/U) and support for U-mode traps, setting a bit in `medeleg` or `mideleg` will delegate the corresponding trap in U-mode to the U-mode trap handler.

In systems with only M-mode, or with both M-mode and U-mode but without U-mode trap support, the `medeleg` and `mideleg` registers should not exist.

In versions 1.9.1 and earlier, these registers existed but were hardwired to zero in M-mode only, or M/U without N systems. There is no reason to require they return zero in those cases, as the `misr` register indicates whether they exist.

When a trap is delegated to a less-privileged mode x , the x **cause** register is written with the trap cause; the x **epc** register is written with the virtual address of the instruction that took the trap; the x PP field of **mstatus** is written with the active privilege mode at the time of the trap; the x PIE field of **mstatus** is written with the value of the x IE field at the time of the trap; and the x IE field of **mstatus** is cleared. The **mcause** and **mepc** registers and the MPP and MPIE fields of **mstatus** are not written.

An implementation shall not hardwire any delegation bits to one, i.e., any trap that can be delegated must support not being delegated. An implementation can choose to subset the delegatable traps, with the supported delegatable bits found by writing one to every bit location, then reading back the value in `medeleg` or `mideleg` to see which bit positions hold a one.

Traps never transition from a more-privileged mode to a less-privileged mode. For example, if M-mode has delegated illegal instruction exceptions to S-mode, and M-mode software later executes an illegal instruction, the trap is taken in M-mode, rather than being delegated to S-mode. By contrast, traps may be taken horizontally. Using the same example, if M-mode has delegated illegal instruction exceptions to S-mode, and S-mode software later executes an illegal instruction, the trap is taken in S-mode.

Delegated interrupts result in the interrupt being masked at the delegator privilege level. For example, if the supervisor timer interrupt (STI) is delegated to S-mode by setting `mideleg[5]`, STIs will not be taken when executing in M-mode. By contrast, if `mideleg[5]` is clear, STIs can be taken in any mode and regardless of current mode will transfer control to M-mode.

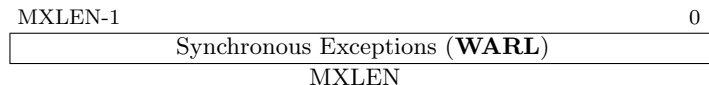
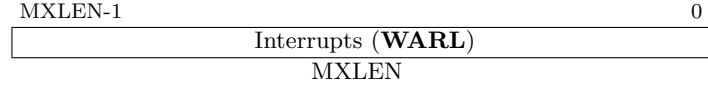


Figure 3.9: Machine Exception Delegation Register `medeleg`.

`mdeleg` has a bit position allocated for every synchronous exception shown in Table 3.6, with the index of the bit position equal to the value returned in the `mcause` register (i.e., setting bit 8 allows user-mode environment calls to be delegated to a lower-privilege trap handler).

midleg holds trap delegation bits for individual interrupts, with the layout of bits matching those in the **mip** register (i.e., STIP interrupt delegation control is located in bit 5).

Some exceptions cannot occur at less privileged modes, and corresponding `xedeleg` bits should be hardwired to zero. In particular, `medeleg`^[11] and `sedeleg`^[11:9] are all hardwired to zero.

Figure 3.10: Machine Interrupt Delegation Register `mideleg`.

3.1.14 Machine Interrupt Registers (`mip` and `mie`)

The `mip` register is an `MXLEN`-bit read/write register containing information on pending interrupts, while `mie` is the corresponding `MXLEN`-bit read/write register containing interrupt enable bits. Only the bits corresponding to lower-privilege software interrupts (`USIP`, `SSIP`), timer interrupts (`UTIP`, `STIP`), and external interrupts (`UEIP`, `SEIP`) in `mip` are writable through this CSR address; the remaining bits are read-only.

The machine-level interrupt registers handle a few root interrupt sources which are assigned a fixed service priority for simplicity, while separate external interrupt controllers can implement a more complex prioritization scheme over a much larger set of interrupts that are then muxed into the machine-level interrupt sources.

Restricted views of the `mip` and `mie` registers appear as the `sip/sie`, and `uip/uie` registers in S-mode and U-mode respectively. If an interrupt is delegated to privilege mode x by setting a bit in the `mideleg` register, it becomes visible in the xip register and is maskable using the xie register. Otherwise, the corresponding bits in xip and xie appear to be hardwired to zero.

MXLEN-1	12	11	10	9	8	7	6	5	4	3	2	1	0
WPRI	MEIP	WPRI	SEIP	UEIP	MTIP	WPRI	STIP	UTIP	MSIP	WPRI	SSIP	USIP	
MXLEN-12	1	1	1	1	1	1	1	1	1	1	1	1	1

Figure 3.11: Machine interrupt-pending register (`mip`).

MXLEN-1	12	11	10	9	8	7	6	5	4	3	2	1	0
WPRI	MEIE	WPRI	SEIE	UEIE	MTIE	WPRI	STIE	UTIE	MSIE	WPRI	SSIE	USIE	
MXLEN-12	1	1	1	1	1	1	1	1	1	1	1	1	1

Figure 3.12: Machine interrupt-enable register (`mie`).

The `MTIP`, `STIP`, `UTIP` bits correspond to timer interrupt-pending bits for machine, supervisor, and user timer interrupts, respectively. The `MTIP` bit is read-only and is cleared by writing to the memory-mapped machine-mode timer compare register. The `UTIP` and `STIP` bits may be written by M-mode software to deliver timer interrupts to lower privilege levels. User and supervisor software may clear the `UTIP` and `STIP` bits with calls to the `AEE` and `SEE` respectively.

There is a separate timer interrupt-enable bit, named `MTIE`, `STIE`, and `UTIE` for M-mode, S-mode, and U-mode timer interrupts respectively.

Each lower privilege level has a separate software interrupt-pending bit (`SSIP`, `USIP`), which can be both read and written by CSR accesses from code running on the local hart at the associated or any higher privilege level. The machine-level `MSIP` bits are written by accesses to memory-mapped control registers, which are used by remote harts to provide machine-mode interprocessor interrupts. Interprocessor interrupts for lower privilege levels are implemented through ABI and

SBI calls to the AEE or SEE respectively, which might ultimately result in a machine-mode write to the receiving hart's MSIP bit. A hart can write its own MSIP bit using the same memory-mapped control register.

The MSIE, SSIE, and USIE fields in the `mie` CSR enable M-mode software interrupts, S-mode software interrupts, and U-mode software interrupts, respectively.

We only allow a hart to directly write its own SSIP or USIP bits when running in the appropriate mode, as other harts might be virtualized and possibly descheduled by higher privilege levels. We rely on ABI and SBI calls to provide interprocessor interrupts for this reason. Machine-mode harts are not virtualized and can directly interrupt other harts by setting their MSIP bits, typically using uncached I/O writes to memory-mapped control registers depending on the platform specification.

The MEIP field in `mip` is a read-only bit that indicates a machine-mode external interrupt is pending. MEIP is set and cleared by a platform-specific interrupt controller, such as the standard platform-level interrupt controller specified in Chapter 7. The MEIE field in `mie` enables machine external interrupts when set.

The SEIP field in `mip` contains a single read-write bit. SEIP may be written by M-mode software to indicate to S-mode that an external interrupt is pending. Additionally, the platform-level interrupt controller may generate supervisor-level external interrupts. The logical-OR of the software-writable bit and the signal from the external interrupt controller is used to generate external interrupts to the supervisor. When the SEIP bit is read with a CSRRW, CSRRS, or CSRRC instruction, the value returned in the `rd` destination register contains the logical-OR of the software-writable bit and the interrupt signal from the interrupt controller. However, the value used in the read-modify-write sequence of a CSRRS or CSRRC instruction is only the software-writable SEIP bit, ignoring the interrupt value from the external interrupt controller.

The SEIP field behavior is designed to allow a higher privilege layer to mimic external interrupts cleanly, without losing any real external interrupts. The behavior of the CSR instructions is slightly modified from regular CSR accesses as a result.

The UEIP field in `mip` provides user-mode external interrupts when the N extension for user-mode interrupts is implemented. It is defined analogously to SEIP.

The MEIE, SEIE, and UEIE fields in the `mie` CSR enable M-mode external interrupts, S-mode external interrupts, and U-mode external interrupts, respectively.

The non-maskable interrupt is not made visible via the `mip` register as its presence is implicitly known when executing the NMI trap handler.

For all the various interrupt types (software, timer, and external), if a privilege level is not supported, or if U-mode is supported but the N extension is not supported, then the associated pending and interrupt-enable bits are hardwired to zero in the `mip` and `mie` registers respectively. Hence, these are all effectively **WARL** fields.

Implementations may add additional platform-specific machine-level interrupt sources to bits 16 and above of the `mip` and `mie` registers. The other unallocated interrupt sources (15–12, 10, 6, and 2) are reserved for future standard use.

An interrupt i will be taken if bit i is set in both `mip` and `mie`, and if interrupts are globally enabled. By default, M-mode interrupts are globally enabled if the hart's current privilege mode is less than M, or if the current privilege mode is M and the MIE bit in the `mstatus` register is set. If bit i in `mideleg` is set, however, interrupts are considered to be globally enabled if the hart's current privilege mode equals the delegated privilege mode (S or U) and that mode's interrupt enable bit (SIE or UIE in `mstatus`) is set, or if the current privilege mode is less than the delegated privilege mode.

Multiple simultaneous interrupts destined for different privilege modes are handled in decreasing order of destined privilege mode. Multiple simultaneous interrupts destined for the same privilege mode are handled in the following decreasing priority order: MEI, MSI, MTI, SEI, SSI, STI, UEI, USI, UTI. Synchronous exceptions are of lower priority than all interrupts.

The machine-level interrupt fixed-priority ordering rules were developed with the following rationale.

Interrupts for higher privilege modes must be serviced before interrupts for lower privilege modes to support pre-emption.

The platform-specific machine-level interrupt sources in bits 16 and above have platform-specific priority, but are typically chosen to have the highest service priority to support very fast local vectored interrupts.

External interrupts are handled before internal (timer/software) interrupts as external interrupts are usually generated by devices that might require low interrupt service times.

Software interrupts are handled before internal timer interrupts, because internal timer interrupts are usually intended for time slicing, where time precision is less important, whereas software interrupts are used for inter-processor messaging. Software interrupts can be avoided when high-precision timing is required, or high-precision timer interrupts can be routed via a different interrupt path. Software interrupts are located in the lowest four bits of `mip` as these are often written by software, and this position allows the use of a single CSR instruction with a five-bit immediate.

Synchronous exceptions are given the lowest priority to minimize worst-case interrupt latency.

3.1.15 Machine Timer Registers (`mtime` and `mtimecmp`)

Platforms provide a real-time counter, exposed as a memory-mapped machine-mode register, `mtime`. `mtime` must run at constant frequency, and the platform must provide a mechanism for determining the timebase of `mtime`.

The `mtime` register has a 64-bit precision on all RV32, RV64, and RV128 systems. Platforms provide a 64-bit memory-mapped machine-mode timer compare register (`mtimecmp`), which causes a timer interrupt to be posted when the `mtime` register contains a value greater than or equal to the value in the `mtimecmp` register. The interrupt remains posted until it is cleared by writing the `mtimecmp` register. The interrupt will only be taken if interrupts are enabled and the MTIE bit is set in the `mie` register.

The timer facility is defined to use wall-clock time rather than a cycle counter to support modern processors that run with a highly variable clock frequency to save energy through dynamic voltage and frequency scaling.

Accurate real-time clocks (RTCs) are relatively expensive to provide (requiring a crystal or MEMS oscillator) and have to run even when the rest of system is powered down, and so

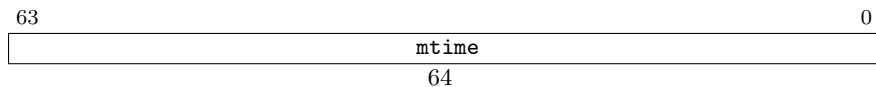


Figure 3.13: Machine time register (memory-mapped control register).

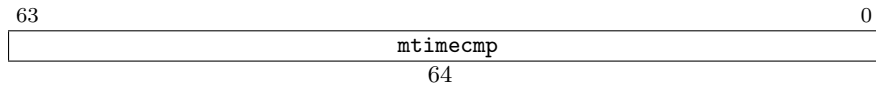


Figure 3.14: Machine time compare register (memory-mapped control register).

there is usually only one in a system located in a different frequency/voltage domain from the processors. Hence, the RTC must be shared by all the harts in a system and accesses to the RTC will potentially incur the penalty of a voltage-level-shifter and clock-domain crossing. It is thus more natural to expose `mtime` as a memory-mapped register than as a CSR.

Lower privilege levels do not have their own `timecmp` registers. Instead, machine-mode software can implement any number of virtual timers on a hart by multiplexing the next timer interrupt into the `mtimecmp` register.

Simple fixed-frequency systems can use a single clock for both cycle counting and wall-clock time.

In RV32, memory-mapped writes to `mtimecmp` modify only one 32-bit part of the register. The following code sequence sets a 64-bit `mtimecmp` value without spuriously generating a timer interrupt due to the intermediate value of the comparand:

```
# New comparand is in a1:a0.
li t0, -1
sw t0, mtimecmp    # No smaller than old value.
sw a1, mtimecmp+4  # No smaller than new value.
sw a0, mtimecmp    # New value.
```

Figure 3.15: Sample code for setting the 64-bit time comparand in RV32 assuming the registers live in a strongly ordered I/O region.

3.1.16 Hardware Performance Monitor

M-mode includes a basic hardware performance-monitoring facility. The `mcycle` CSR counts the number of cycles the hart has executed. The `minstret` CSR counts the number of instructions the hart has retired. The `mcycle` and `minstret` registers have 64-bit precision on all RV32, RV64, and RV128 systems.

The counter registers have an arbitrary value after system reset, and can be written with a given value. Any CSR write takes effect after the writing instruction has otherwise completed.

The hardware performance monitor includes 29 additional 64-bit event counters, `mhpcounter3–mhpcounter31`. The event selector CSRs, `mhpmevent3–mhpmevent31`, are MXLEN-bit **WARL** registers that control which event causes the corresponding counter to increment. The meaning of these events is defined by the platform, but event 0 is defined to mean “no event.” All coun-

ters should be implemented, but a legal implementation is to hard-wire both the counter and its corresponding event selector to 0.

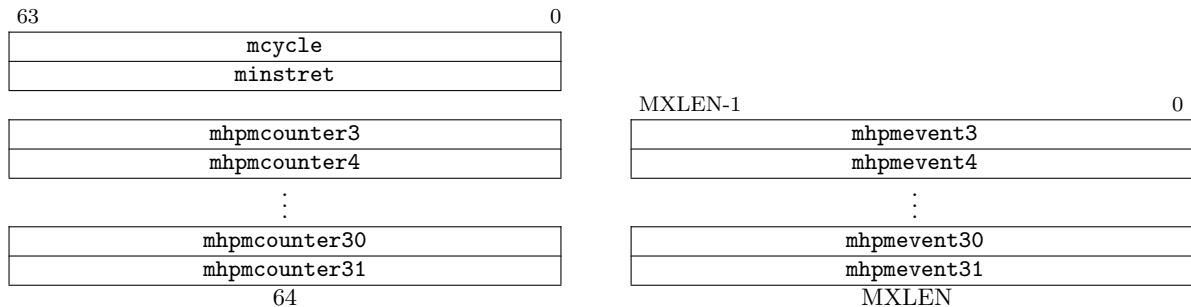


Figure 3.16: Hardware performance monitor counters.

All of these counters have 64-bit precision on RV32, RV64, and RV128.

On RV32 only, reads of the `mcycle`, `minstret`, and `mhpmpcounter n` CSRs return the low 32 bits, while reads of the `mcycleh`, `minstreth`, and `mhpmpcounter n h` CSRs return bits 63–32 of the corresponding counter.

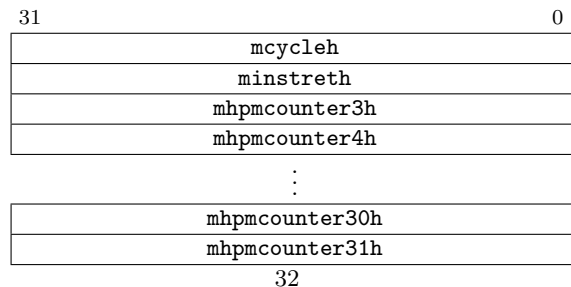


Figure 3.17: Upper 32 bits of hardware performance monitor counters, RV32 only.

On RV128 systems, the 64-bit values in `mcycle`, `minstret`, and `mhpmpcounter n` are sign-extended to 128-bits when read.

On RV128 systems, both signed and unsigned 64-bit values are held in a canonical form with bit 63 repeated in all higher bit positions. The counters are 64-bit values even in RV128, and so the counter CSR reads preserve the sign-extension invariant. Implementations may choose to implement fewer bits of the counters, provided software would be unlikely to experience wraparound (e.g., 2^{63} instructions executed) and thereby avoid having to actually implement the sign-extension circuitry.

3.1.17 Counter-Enable Registers (`[m|s]counteren`)

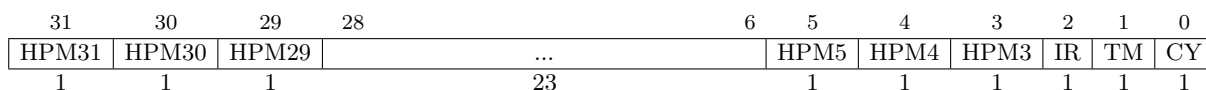


Figure 3.18: Counter-enable registers (`mcounteren` and `scounteren`).

The counter-enable registers `mcounteren` and `scounteren` are 32-bit registers that control the availability of the hardware performance monitoring counters to the next-lowest privileged mode.

The settings in these registers only control accessibility. The act of reading or writing these registers, does not affect the underlying counters, which continue to increment even when not accessible.

When the CY, TM, IR, or HPM n bit in the `mcounter n` register is clear, attempts to read the `cycle`, `time`, `instret`, or `hpmcounter n` register while executing in S-mode or U-mode will cause an illegal instruction exception. When one of these bits is set, access to the corresponding register is permitted in the next implemented privilege mode (S-mode if implemented, otherwise U-mode).

If S-mode is implemented, the same bit positions in the `scounteren` register analogously control access to these registers while executing in U-mode. If S-mode is permitted to access a counter register and the corresponding bit is set in `scounteren`, then U-mode is also permitted to access that register.

Registers `mcounteren` and `scounteren` are **WARL** registers that must be implemented if U-mode and S-mode are implemented. Any of the bits may contain a hardwired value of zero, indicating reads to the corresponding counter will cause an exception when executing in a less-privileged mode.

The counter-enable bits support two common use cases with minimal hardware. For systems that do not need high-performance timers and counters, machine-mode software can trap accesses and implement all features in software. For systems that need high-performance timers and counters but are not concerned with obfuscating the underlying hardware counters, the counters can be directly exposed to lower privilege modes.

The `cycle`, `instret`, and `hpmcounter n` CSRs are read-only shadows of `mcycle`, `minstret`, and `mhpmcounter n` , respectively. The `time` CSR is a read-only shadow of the memory-mapped `mtime` register.

Implementations can convert reads of the `time CSR` into loads to the memory-mapped `mtime` register, or hard-wire the `TM` bit in `mcounteren` to 0 and emulate this functionality in `M-mode` software.

3.1.18 Machine Scratch Register (mscratch)

The `mscratch` register is an MXLEN-bit read/write register dedicated for use by machine mode. Typically, it is used to hold a pointer to a machine-mode hart-local context space and swapped with a user register upon entry to an M-mode trap handler.

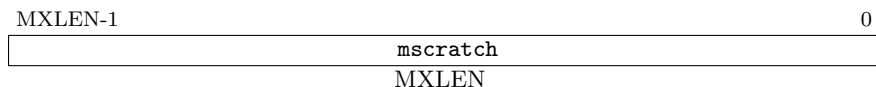


Figure 3.19: Machine-mode scratch register.

The MIPS ISA allocated two user registers (`k0/k1`) for use by the operating system. Although the MIPS scheme provides a fast and simple implementation, it also reduces available user registers, and does not scale to further privilege levels, or nested traps. It can also require both

registers are cleared before returning to user level to avoid a potential security hole and to provide deterministic debugging behavior.

The RISC-V user ISA was designed to support many possible privileged system environments and so we did not want to infect the user-level ISA with any OS-dependent features. The RISC-V CSR swap instructions can quickly save/restore values to the `mscratch` register. Unlike the MIPS design, the OS can rely on holding a value in the `mscratch` register while the user context is running.

3.1.19 Machine Exception Program Counter (`mepc`)

`mepc` is an MXLEN-bit read/write register formatted as shown in Figure 3.20. The low bit of `mepc` (`mepc[0]`) is always zero. On implementations that support only IALIGN=32, the two low bits (`mepc[1:0]`) are always zero.

If an implementation allows IALIGN to be either 16 or 32 (by changing CSR `misa`, for example), then, whenever IALIGN=32, bit `mepc[1]` is masked on reads so that it appears to be 0. This masking occurs also for the implicit read by the MRET instruction. Though masked, `mepc[1]` remains writable when IALIGN=32.

`mepc` is a **WARL** register that must be able to hold all valid physical and virtual addresses. It need not be capable of holding all possible invalid addresses. Implementations may convert some invalid address patterns into other invalid addresses prior to writing them to `mepc`.

When a trap is taken into M-mode, `mepc` is written with the virtual address of the instruction that was interrupted or that encountered the exception. Otherwise, `mepc` is never written by the implementation, though it may be explicitly written by software.

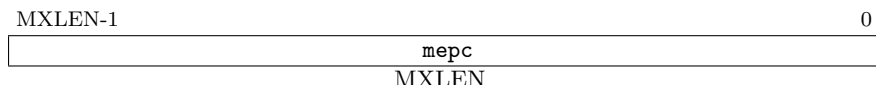


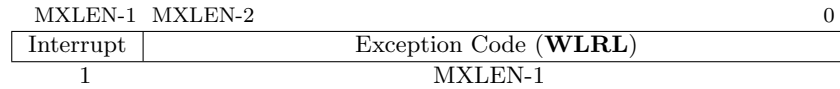
Figure 3.20: Machine exception program counter register.

3.1.20 Machine Cause Register (`mcause`)

The `mcause` register is an MXLEN-bit read-write register formatted as shown in Figure 3.21. When a trap is taken into M-mode, `mcause` is written with a code indicating the event that caused the trap. Otherwise, `mcause` is never written by the implementation, though it may be explicitly written by software.

The Interrupt bit in the `mcause` register is set if the trap was caused by an interrupt. The Exception Code field contains a code identifying the last exception. Table 3.6 lists the possible machine-level exception codes. The Exception Code is an **WLRL** field, so is only guaranteed to hold supported exception codes.

We do not distinguish privileged instruction exceptions from illegal opcode exceptions. This simplifies the architecture and also hides details of which higher-privilege instructions are supported by an implementation. The privilege level servicing the trap can implement a policy on whether

Figure 3.21: Machine Cause register `mcause`.

Interrupt	Exception Code	Description
1	0	User software interrupt
1	1	Supervisor software interrupt
1	2	<i>Reserved</i>
1	3	Machine software interrupt
1	4	User timer interrupt
1	5	Supervisor timer interrupt
1	6	<i>Reserved</i>
1	7	Machine timer interrupt
1	8	User external interrupt
1	9	Supervisor external interrupt
1	10	<i>Reserved</i>
1	11	Machine external interrupt
1	≥ 12	<i>Reserved</i>
0	0	Instruction address misaligned
0	1	Instruction access fault
0	2	Illegal instruction
0	3	Breakpoint
0	4	Load address misaligned
0	5	Load access fault
0	6	Store/AMO address misaligned
0	7	Store/AMO access fault
0	8	Environment call from U-mode
0	9	Environment call from S-mode
0	10	<i>Reserved</i>
0	11	Environment call from M-mode
0	12	Instruction page fault
0	13	Load page fault
0	14	<i>Reserved</i>
0	15	Store/AMO page fault
0	≥ 16	<i>Reserved</i>

Table 3.6: Machine cause register (`mcause`) values after trap.

these need to be distinguished, and if so, whether a given opcode should be treated as illegal or privileged.

Interrupts can be separated from other traps with a single branch on the sign of the `mcause` register value. A shift left can remove the interrupt bit and scale the exception codes to index into a trap vector table.

cache misses or slow uncached accesses when loads are used to fetch the instruction into a data register. There is also a problem of atomicity if another agent is manipulating the instruction memory, as might occur in a dynamic translation system.

A requirement is that the entire instruction (or at least the first $XLEN$ bits) are fetched into `mtval` before taking the trap. This should not constrain implementations, which would typically fetch the entire instruction before attempting to decode the instruction, and avoids complicating software handlers.

A value of zero in `mtval` signifies either that the feature is not supported, or an illegal zero instruction was fetched. A load from the instruction memory pointed to by `mepc` can be used to distinguish these two cases (or alternatively, the system configuration information can be interrogated to install the appropriate trap handling before runtime).

If the hardware platform specifies that no exceptions set `mtval` to a nonzero value, then it may be hardwired to zero. Otherwise, `mtval` is a **WARL** register that must be able to hold all valid physical and virtual addresses and the value 0. It need not be capable of holding all possible invalid addresses. Implementations may convert some invalid address patterns into other invalid addresses prior to writing them to `mtval`. If the feature to return the faulting instruction bits is implemented, `mtval` must also be able to hold all values less than 2^N , where N is the smaller of $XLEN$ and the width of the longest supported instruction.

3.2 Machine-Mode Privileged Instructions

3.2.1 Environment Call and Breakpoint

31	20 19	15 14	12 11	7 6	0
funct12	rs1	funct3	rd	opcode	
12	5	3	5	7	
ECALL	0	PRIV	0	SYSTEM	
EBREAK	0	PRIV	0	SYSTEM	

The ECALL instruction is used to make a request to the supporting execution environment. When executed in U-mode, S-mode, or M-mode, it generates an environment-call-from-U-mode exception, environment-call-from-S-mode exception, or environment-call-from-M-mode exception, respectively, and performs no other operation.

ECALL generates a different exception for each originating privilege mode so that environment call exceptions can be selectively delegated. A typical use case for Unix-like operating systems is to delegate to S-mode the environment-call-from-U-mode exception but not the others.

The EBREAK instruction is used by debuggers to cause control to be transferred back to a debugging environment. It generates a breakpoint exception and performs no other operation.

As described in the “C” Standard Extension for Compressed Instructions in Volume I of this manual, the C.EBREAK instruction performs the same operation as the EBREAK instruction.

ECALL and EBREAK cause the receiving privilege mode’s `epc` register to be set to the address of the ECALL or EBREAK instruction itself, *not* the address of the following instruction.

3.2.2 Trap-Return Instructions

Instructions to return from trap are encoded under the PRIV minor opcode.

31	20 19	15 14	12 11	7 6	0
funct12	rs1	funct3	rd	opcode	
12	5	3	5	7	
MRET/SRET/URET	0	PRIV	0	SYSTEM	

To return after handling a trap, there are separate trap return instructions per privilege level: MRET, SRET, and URET. MRET is always provided. SRET must be provided if supervisor mode is supported, and should raise an illegal instruction exception otherwise. SRET should also raise an illegal instruction exception when TSR=1 in `mstatus`, as described in Section 3.1.10. URET is only provided if user-mode traps are supported, and should raise an illegal instruction otherwise. An `xRET` instruction can be executed in privilege mode `x` or higher, where executing a lower-privilege `xRET` instruction will pop the relevant lower-privilege interrupt enable and privilege mode stack. In addition to manipulating the privilege stack as described in Section 3.1.7, `xRET` sets the `pc` to the value stored in the `xepc` register.

Previously, there was only a single ERET instruction (which was also earlier known as SRET). To support the addition of user-level interrupts, we needed to add a separate URET instruction to continue to allow classic virtualization of OS code using the ERET instruction. It then became more orthogonal to support a different xRET instruction per privilege level.

If the A extension is supported, the `xRET` instruction is allowed to clear any outstanding LR address reservation but is not required to. Trap handlers should explicitly clear the reservation if required (e.g., by using a dummy SC) before executing the `xRET`.

If xRET instructions always cleared LR reservations, it would be impossible to single-step through LR/SC sequences using a debugger.

3.2.3 Wait for Interrupt

The Wait for Interrupt instruction (WFI) provides a hint to the implementation that the current hart can be stalled until an interrupt might need servicing. Execution of the WFI instruction can also be used to inform the hardware platform that suitable interrupts should preferentially be routed to this hart. WFI is available in all of the supported S and M privilege modes, and optionally available to U-mode for implementations that support U-mode interrupts. This instruction may raise an illegal instruction exception when TW=1 in `mstatus`, as described in Section 3.1.10.

31	20 19	15 14	12 11	7 6	0
funct12	rs1	funct3	rd	opcode	
12	5	3	5	7	
WFI	0	PRIV	0	SYSTEM	

If an enabled interrupt is present or later becomes present while the hart is stalled, the interrupt exception will be taken on the following instruction, i.e., execution resumes in the trap handler and $\text{mepc} = \text{pc} + 4$.

The following instruction takes the interrupt exception and trap, so that a simple return from the trap handler will execute code after the WFI instruction.

The WFI instruction is just a hint, and a legal implementation is to implement WFI as a NOP.

If the implementation does not stall the hart on execution of the instruction, then the interrupt will be taken on some instruction in the idle loop containing the WFI, and on a simple return from the handler, the idle loop will resume execution.

We have removed the earlier requirement that implementations ignore the rs1 and rd fields, so non-zero values in these fields should now raise illegal instruction exceptions.

The WFI instruction can also be executed when interrupts are disabled. The operation of WFI must be unaffected by the global interrupt bits in `mstatus` (MIE/SIE/UIE) and the delegation registers `[m|s]ideleg` (i.e., the hart must resume if a locally enabled interrupt becomes pending, even if it has been delegated to a less-privileged mode), but should honor the individual interrupt enables (e.g., MTIE) (i.e., implementations should avoid resuming the hart if the interrupt is pending but not individually enabled). WFI is also required to resume execution for locally enabled interrupts pending at any privilege level, regardless of the global interrupt enable at each privilege level.

If the event that causes the hart to resume execution does not cause an interrupt to be taken, execution will resume at $\text{pc} + 4$, and software must determine what action to take, including looping back to repeat the WFI if there was no actionable event.

By allowing wakeup when interrupts are disabled, an alternate entry point to an interrupt handler can be called that does not require saving the current context, as the current context can be saved or discarded before the WFI is executed.

As implementations are free to implement WFI as a NOP, software must explicitly check for any relevant pending but disabled interrupts in the code following an WFI, and should loop back to the WFI if no suitable interrupt was detected. The `mip`, `sip`, or `uip` registers can be interrogated to determine the presence of any interrupt in machine, supervisor, or user mode respectively.

The operation of WFI is unaffected by the delegation register settings.

WFI is defined so that an implementation can trap into a higher privilege mode, either immediately on encountering the WFI or after some interval to initiate a machine-mode transition to a lower power state, for example.

The same “wait-for-event” template might be used for possible future extensions that wait on memory locations changing, or message arrival.

3.3 Reset

Upon reset, a hart’s privilege mode is set to M. The `mstatus` fields MIE and MPRV are reset to 0. The `pc` is set to an implementation-defined reset vector. The `mcause` register is set to a value indicating the cause of the reset. All other hart state is undefined.

The **mcause** values after reset have implementation-specific interpretation, but the value 0 should be returned on implementations that do not distinguish different reset conditions. Implementations that distinguish different reset conditions should only use 0 to indicate the most complete reset (e.g., hard reset).

Some designs may have multiple causes of reset (e.g., power-on reset, external hard reset, brownout detected, watchdog timer elapse, sleep-mode wakeup), which machine-mode software and debuggers may wish to distinguish.

***mcause** reset values may alias **mcause** values following synchronous exceptions. There should be no ambiguity in this overlap, since on reset the **pc** is typically set to a different value than on other traps.*

3.4 Non-Maskable Interrupts

Non-maskable interrupts (NMIs) are only used for hardware error conditions, and cause an immediate jump to an implementation-defined NMI vector running in M-mode regardless of the state of a hart's interrupt enable bits. The **mepc** register is written with the address of the next instruction to be executed at the time the NMI was taken, and **mcause** is set to a value indicating the source of the NMI. The NMI can thus overwrite state in an active machine-mode interrupt handler.

The values written to **mcause** on an NMI are implementation-defined, but a value of 0 is reserved to mean “unknown cause” and implementations that do not distinguish sources of NMIs via the **mcause** register should return 0.

Unlike resets, NMIs do not reset processor state, enabling diagnosis, reporting, and possible containment of the hardware error.

3.5 Physical Memory Attributes

The physical memory map for a complete system includes various address ranges, some corresponding to memory regions, some to memory-mapped control registers, and some to empty holes in the address space. Some memory regions might not support reads, writes, or execution; some might not support subword or subblock accesses; some might not support atomic operations; and some might not support cache coherence or might have different memory models. Similarly, memory-mapped control registers vary in their supported access widths, support for atomic operations, and whether read and write accesses have associated side effects. In RISC-V systems, these properties and capabilities of each region of the machine's physical address space are termed *physical memory attributes* (PMAs). This section describes RISC-V PMA terminology and how RISC-V systems implement and check PMAs.

PMAs are inherent properties of the underlying hardware and rarely change during system operation. Unlike physical memory protection values described in Section 3.6, PMAs do not vary by execution context. The PMAs of some memory regions are fixed at chip design time—for example, for an on-chip ROM. Others are fixed at board design time, depending, for example, on which other chips are connected to off-chip buses. Off-chip buses might also support devices that could be changed on every power cycle (cold pluggable) or dynamically while the system is running (hot

pluggable). Some devices might be configurable at run time to support different uses that imply different PMAs—for example, an on-chip scratchpad RAM might be cached privately by one core in one end-application, or accessed as a shared non-cached memory in another end-application.

Most systems will require that at least some PMAs are dynamically checked in hardware later in the execution pipeline after the physical address is known, as some operations will not be supported at all physical memory addresses, and some operations require knowing the current setting of a configurable PMA attribute. While many other architectures specify some PMAs in the virtual memory page tables and use the TLB to inform the pipeline of these properties, this approach injects platform-specific information into a virtualized layer and can cause system errors unless attributes are correctly initialized in each page-table entry for each physical memory region. In addition, the available page sizes might not be optimal for specifying attributes in the physical memory space, leading to address-space fragmentation and inefficient use of expensive TLB entries.

For RISC-V, we separate out specification and checking of PMAs into a separate hardware structure, the *PMA checker*. In many cases, the attributes are known at system design time for each physical address region, and can be hardwired into the PMA checker. Where the attributes are run-time configurable, platform-specific memory-mapped control registers can be provided to specify these attributes at a granularity appropriate to each region on the platform (e.g., for an on-chip SRAM that can be flexibly divided between cacheable and uncacheable uses). PMAs are checked for any access to physical memory, including accesses that have undergone virtual to physical memory translation. To aid in system debugging, we strongly recommend that, where possible, RISC-V processors precisely trap physical memory accesses that fail PMA checks. Precisely trapped PMA violations manifest as load, store, or instruction-fetch access exceptions, distinct from virtual-memory page-fault exceptions. Precise PMA traps might not always be possible, for example, when probing a legacy bus architecture that uses access failures as part of the discovery mechanism. In this case, error responses from slave devices will be reported as imprecise bus-error interrupts.

PMAs must also be readable by software to correctly access certain devices or to correctly configure other hardware components that access memory, such as DMA engines. As PMAs are tightly tied to a given physical platform’s organization, many details are inherently platform-specific, as is the means by which software can learn the PMA values for a platform. The configuration string (Chapter 8) can encode PMAs for on-chip devices and might also describe on-chip controllers for off-chip buses that can be dynamically interrogated to discover attached device PMAs. Some devices, particularly legacy buses, do not support discovery of PMAs and so will give error responses or time out if an unsupported access is attempted. Typically, platform-specific machine-mode code will extract PMAs and ultimately present this information to higher-level less-privileged software using some standard representation.

Where platforms support dynamic reconfiguration of PMAs, an interface will be provided to set the attributes by passing requests to a machine-mode driver that can correctly reconfigure the platform. For example, switching cacheability attributes on some memory regions might involve platform-specific operations, such as cache flushes, that are available only to machine-mode.

3.5.1 Main Memory versus I/O versus Empty Regions

The most important characterization of a given memory address range is whether it holds regular main memory, or I/O devices, or is empty. Regular main memory is required to have a number of properties, specified below, whereas I/O devices can have a much broader range of attributes. Memory regions that do not fit into regular main memory, for example, device scratchpad RAMs, are categorized as I/O regions. Empty regions are also classified as I/O regions but with attributes specifying that no accesses are supported.

3.5.2 Supported Access Type PMAs

Access types specify which access widths, from 8-bit byte to long multi-word burst, are supported, and also whether misaligned accesses are supported for each access width.

Although software running on a RISC-V hart cannot directly generate bursts to memory, software might have to program DMA engines to access I/O devices and might therefore need to know which access sizes are supported.

Main memory regions always support read, write, and execute of all access widths required by the attached devices.

In some cases, the design of a processor or device accessing main memory might support other widths, but must be able to function with the types supported by the main memory.

I/O regions can specify which combinations of read, write, or execute accesses to which data widths are supported.

3.5.3 Atomicity PMAs

Atomicity PMAs describes which atomic instructions are supported in this address region. Main memory regions must support the atomic operations required by the processors attached. I/O regions may only support a subset or none of the processor-supported atomic operations.

Support for atomic instructions is divided into two categories: *LR/SC* and *AMOs*. Within AMOs, there are four levels of support: *AMONone*, *AMOSwap*, *AMOLogical*, and *AMOArithmetic*. *AMONone* indicates that no AMO operations are supported. *AMOSwap* indicates that only `amoswap` instructions are supported in this address range. *AMOLogical* indicates that swap instructions plus all the logical AMOs (`amoand`, `amoor`, `amoxor`) are supported. *AMOArithmetic* indicates that all RISC-V AMOs are supported. For each level of support, naturally aligned AMOs of a given width are supported if the underlying memory region supports reads and writes of that width.

We recommend providing at least AMOLogical support for I/O regions where possible. Most I/O regions will not support LR/SC accesses, as these are most conveniently built on top of a cache-coherence scheme.

AMO Class	Supported Operations
AMONone	<i>None</i>
AMOSwap	amoswap
AMOLogical	above + amoand , amoor , amoxor
AMOArithmetic	above + amoadd , amomin , amomax , amominu , amomaxu

Table 3.7: Classes of AMOs supported by I/O regions. Main memory regions must always support all AMOs required by the processor.

Memory regions that support LR/SC will define the conditions under which LR/SC sequences must succeed or must fail. Main memory must guarantee the eventual success of any LR/SC sequence that meets the requirements described in the user ISA specification.

Memory regions that support aligned LR/SC or aligned AMOs might also support misaligned LR/SC or misaligned AMOs for some addresses and access widths. If, for a given address and access width, a misaligned LR/SC or AMO generates a misaligned address exception, then *all* loads, stores, LR/SCs, and AMOs using that address and access width must generate misaligned address exceptions.

The standard “A” extension does not support misaligned AMOs or LR/SC pairs. Support for misaligned AMOs is provided by the standard “Zam” extension. Support for misaligned LR/SC sequences is not currently standardized, so LR and SC to misaligned addresses must raise an exception.

Mandating that misaligned loads and stores trap wherever misaligned AMOs trap permits the emulation of misaligned AMOs in an M-mode trap handler. The handler guarantees atomicity by acquiring a global mutex and emulating the access within the critical section. Provided that the handler for misaligned loads and stores uses the same mutex, all accesses to a given address that use the same word size will be mutually atomic.

3.5.4 Memory-Ordering PMAs

Regions of the address space are classified as either *main memory* or *I/O* for the purposes of ordering by the FENCE instruction and atomic-instruction ordering bits.

Accesses by one hart to main memory regions are observable not only by other harts but also by other devices with the capability to initiate requests in the main memory system (e.g., DMA engines). Main memory regions always have the standard RISC-V relaxed memory model.

Accesses by one hart to the I/O space are observable not only by other harts and bus mastering devices, but also by targeted slave I/O devices. Within I/O, regions may further be classified as implementing either *relaxed* or *strong* ordering. A relaxed I/O region has no ordering guarantees on how memory accesses made by one hart are observable by different harts or I/O devices beyond those enforced by FENCE and AMO instructions. A strongly ordered I/O region ensures that all accesses made by a hart to that region are only observable in program order by all other harts or I/O devices.

Each strongly ordered I/O region specifies a numbered ordering channel, which is a mechanism by which ordering guarantees can be provided between different I/O regions. Channel 0 is used to

indicate point-to-point strong ordering only, where only accesses by the hart to the single associated I/O region are strongly ordered.

Channel 1 is used to provide global strong ordering across all I/O regions. Any accesses by a hart to any I/O region associated with channel 1 can only be observed to have occurred in program order by all other harts and I/O devices, including relative to accesses made by that hart to relaxed I/O regions or strongly ordered I/O regions with different channel numbers. In other words, any access to a region in channel 1 is equivalent to executing a `fence io,io` instruction before and after the instruction.

Other larger channel numbers provide program ordering to accesses by that hart across any regions with the same channel number.

Systems might support dynamic configuration of ordering properties on each memory region.

Strong ordering can be used to improve compatibility with legacy device driver code, or to enable increased performance compared to insertion of explicit ordering instructions when the implementation is known to not reorder accesses.

Local strong ordering (channel 0) is the default form of strong ordering as it is often straightforward to provide if there is only a single in-order communication path between the hart and the I/O device.

Generally, different strongly ordered I/O regions can share the same ordering channel without additional ordering hardware if they share the same interconnect path and the path does not reorder requests.

3.5.5 Coherence and Cacheability PMAs

Coherence is a property defined for a single physical address, and indicates that writes to that address by one agent will eventually be made visible to other agents in the system. Coherence is not to be confused with the memory consistency model of a system, which defines what values a memory read can return given the previous history of reads and writes to the entire memory system. In RISC-V platforms, the use of hardware-incoherent regions is discouraged due to software complexity, performance, and energy impacts.

The cacheability of a memory region should not affect the software view of the region except for differences reflected in other PMAs, such as main memory versus I/O classification, memory ordering, supported accesses and atomic operations, and coherence. For this reason, we treat cacheability as a platform-level setting managed by machine-mode software only.

Where a platform supports configurable cacheability settings for a memory region, a platform-specific machine-mode routine will change the settings and flush caches if necessary, so the system is only incoherent during the transition between cacheability settings. This transitory state should not be visible to lower privilege levels.

We categorize RISC-V caches into three types: master-private, shared, and slave-private. Master-private caches are attached to a single master agent, i.e., one that issues read/write requests to the memory system. Shared caches are located between masters and slaves and may be hierarchically organized. Slave-private caches do not impact coherence, as they are local to a single slave and do not affect other PMAs at a master, so are not considered further here.

We use private cache to mean a master-private cache in the following section, unless explicitly stated otherwise.

Coherence is straightforward to provide for a shared memory region that is not cached by any agent. The PMA for such a region would simply indicate it should not be cached in a private or shared cache.

Coherence is also straightforward for read-only regions, which can be safely cached by multiple agents without requiring a cache-coherence scheme. The PMA for this region would indicate that it can be cached, but that writes are not supported.

Some read-write regions might only be accessed by a single agent, in which case they can be cached privately by that agent without requiring a coherence scheme. The PMA for such regions would indicate they can be cached. The data can also be cached in a shared cache, as other agents should not access the region.

If an agent can cache a read-write region that is accessible by other agents, whether caching or non-caching, a cache-coherence scheme is required to avoid use of stale values. In regions lacking hardware cache coherence (hardware-incoherent regions), cache coherence can be implemented entirely in software, but software coherence schemes are notoriously difficult to implement correctly and often have severe performance impacts due to the need for conservative software-directed cache-flushing. Hardware cache-coherence schemes require more complex hardware and can impact performance due to the cache-coherence probes, but are otherwise invisible to software.

For each hardware cache-coherent region, the PMA would indicate that the region is coherent and which hardware coherence controller to use if the system has multiple coherence controllers. For some systems, the coherence controller might be an outer-level shared cache, which might itself access further outer-level cache-coherence controllers hierarchically.

Most memory regions within a platform will be coherent to software, because they will be fixed as either uncached, read-only, hardware cache-coherent, or only accessed by one agent.

3.5.6 Idempotency PMAs

Idempotency PMAs describe whether reads and writes to an address region are idempotent. Main memory regions are assumed to be idempotent. For I/O regions, idempotency on reads and writes can be specified separately (e.g., reads are idempotent but writes are not). If accesses are non-idempotent, i.e., there is potentially a side effect on any read or write access, then speculative or redundant accesses must be avoided.

For the purposes of defining the idempotency PMAs, changes in observed memory ordering created by redundant accesses are not considered a side effect.

While hardware should always be designed to avoid speculative or redundant accesses to memory regions marked as non-idempotent, it is also necessary to ensure software or compiler optimizations do not generate spurious accesses to non-idempotent memory regions.

Non-idempotent regions might not support misaligned accesses, in which case software might emulate misaligned accesses as a sequence of aligned accesses, each possibly causing side effects. Therefore, portable software should not issue misaligned accesses to non-idempotent regions.

3.6 Physical Memory Protection

To support secure processing and contain faults, it is desirable to limit the physical addresses accessible by software running on a hart. An optional physical memory protection (PMP) unit

provides per-hart machine-mode control registers to allow physical memory access privileges (read, write, execute) to be specified for each physical memory region. The PMP values are checked in parallel with the PMA checks described in Section 3.5.

The granularity of PMP access control settings are platform-specific and within a platform may vary by physical memory region, but the standard PMP encoding supports regions as small as four bytes. Certain regions' privileges can be hardwired—for example, some regions might only ever be visible in machine mode but in no lower-privilege layers.

Platforms vary widely in demands for physical memory protection, and some platforms may provide other PMP structures in addition to or instead of the scheme described in this section.

PMP checks are applied to all accesses when the hart is running in S or U modes, and for loads and stores when the MPRV bit is set in the `mstatus` register and the MPP field in the `mstatus` register contains S or U. PMP checks are also applied to page-table accesses for virtual-address translation, for which the effective privilege mode is S. Optionally, PMP checks may additionally apply to M-mode accesses, in which case the PMP registers themselves are locked, so that even M-mode software cannot change them without a system reset. In effect, PMP can *grant* permissions to S and U modes, which by default have none, and can *revoke* permissions from M-mode, which by default has full permissions.

PMP violations are always trapped precisely at the processor.

3.6.1 Physical Memory Protection CSRs

PMP entries are described by an 8-bit configuration register and one MXLEN-bit address register. Some PMP settings additionally use the address register associated with the preceding PMP entry. Up to 16 PMP entries are supported. If any PMP entries are implemented, then all PMP CSRs must be implemented, but all PMP CSR fields are **WARL** and may be hardwired to zero. PMP CSRs are only accessible to M-mode.

The PMP configuration registers are densely packed into CSRs to minimize context-switch time. For RV32, four CSRs, `pmpcfg0`–`pmpcfg3`, hold the configurations `pmp0cfg`–`pmp15cfg` for the 16 PMP entries, as shown in Figure 3.23. For RV64, `pmpcfg0` and `pmpcfg2` hold the configurations for the 16 PMP entries, as shown in Figure 3.24; `pmpcfg1` and `pmpcfg3` are illegal.

RV64 systems use `pmpcfg2`, rather than `pmpcfg1`, to hold configurations for PMP entries 8–15. This design reduces the cost of supporting multiple MXLEN values, since the configurations for PMP entries 8–11 appear in `pmpcfg2[31:0]` for both RV32 and RV64.

The PMP address registers are CSRs named `pmpaddr0`–`pmpaddr15`. Each PMP address register encodes bits 33–2 of a 34-bit physical address for RV32, as shown in Figure 3.25. For RV64, each PMP address register encodes bits 55–2 of a 56-bit physical address, as shown in Figure 3.26. Not all physical address bits may be implemented, and so the `pmpaddr` registers are **WARL**.

The Sv32 page-based virtual-memory scheme described in Section 4.3 supports 34-bit physical addresses for RV32, so the PMP scheme must support addresses wider than XLEN for RV32. The Sv39 and Sv48 page-based virtual-memory schemes described in Sections 4.4 and 4.5 support a 56-bit physical address space, so the RV64 PMP address registers impose the same limit.

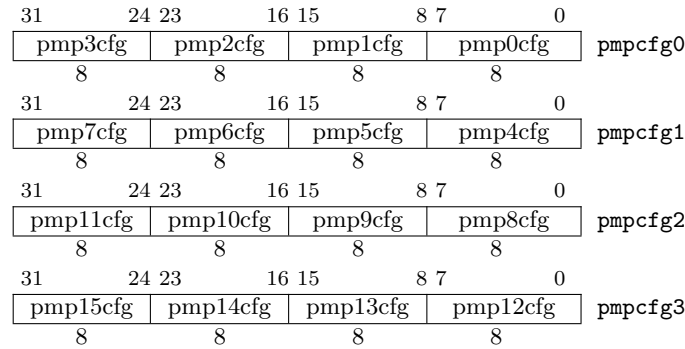


Figure 3.23: RV32 PMP configuration CSR layout.

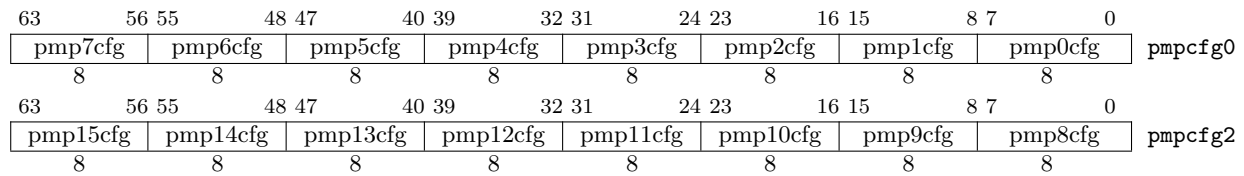


Figure 3.24: RV64 PMP configuration CSR layout.

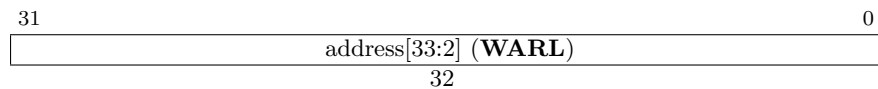


Figure 3.25: PMP address register format, RV32.

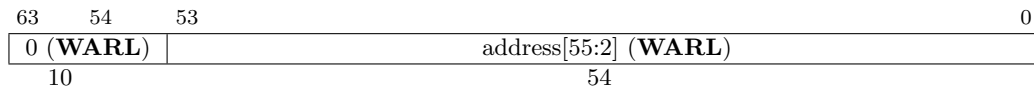


Figure 3.26: PMP address register format, RV64.

Figure 3.27 shows the layout of a PMP configuration register. The R, W, and X bits, when set, indicate that the PMP entry permits read, write, and instruction execution, respectively. When one of these bits is clear, the corresponding access type is denied. The combination R=0 and W=1 is reserved for future use. The remaining two fields, A and L, are described in the following sections.

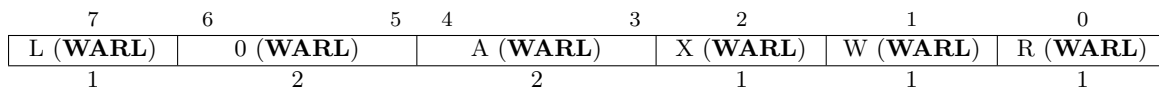


Figure 3.27: PMP configuration register format.

Attempting to fetch an instruction from a PMP region that does not have execute permissions raises a fetch access exception. Attempting to execute a load or load-reserved instruction whose effective address lies within a PMP region without read permissions raises a load access exception. Attempting to execute a store, store-conditional (regardless of success), or AMO instruction whose effective address lies within a PMP region without write permissions raises a store access exception.

Address Matching

The A field in a PMP entry’s configuration register encodes the address-matching mode of the associated PMP address register. The encoding of this field is shown in Table 3.8. When A=0, this PMP entry is disabled and matches no addresses. Two other address-matching modes are supported: naturally aligned power-of-2 regions (NAPOT), including the special case of naturally aligned four-byte regions (NA4); and the top boundary of an arbitrary range (TOR). These modes support four-byte granularity.

A	Name	Description
0	OFF	Null region (disabled)
1	TOR	Top of range
2	NA4	Naturally aligned four-byte region
3	NAPOT	Naturally aligned power-of-two region, ≥ 8 bytes

Table 3.8: Encoding of A field in PMP configuration registers.

NAPOT ranges make use of the low-order bits of the associated address register to encode the size of the range, as shown in Table 3.9.

pmpaddr	pmpcfg.A	Match type and size
yyyy...yyy	NA4	4-byte NAPOT range
yyyy...yyy0	NAPOT	8-byte NAPOT range
yyyy...yy01	NAPOT	16-byte NAPOT range
yyyy...y011	NAPOT	32-byte NAPOT range
...
yy01...1111	NAPOT	2^{XLEN} -byte NAPOT range
y011...1111	NAPOT	2^{XLEN+1} -byte NAPOT range
0111...1111	NAPOT	2^{XLEN+2} -byte NAPOT range
1111...1111	NAPOT	<i>Reserved</i>

Table 3.9: NAPOT range encoding in PMP address and configuration registers.

If TOR is selected, the associated address register forms the top of the address range, and the preceding PMP address register forms the bottom of the address range. If PMP entry i ’s A field is set to TOR, the entry matches any address y such that $\text{pmpaddr}_{i-1} \leq y < \text{pmpaddr}_i$. If PMP entry 0’s A field is set to TOR, zero is used for the lower bound, and so it matches any address $y < \text{pmpaddr}_0$.

Although the PMP mechanism supports regions as small as four bytes, platforms may specify coarser PMP regions. In general, the PMP grain is 2^{G+2} bytes and must be the same across all PMP regions. When $G \geq 1$, the NA4 mode may not be selected. When $G \geq 2$ and $\text{pmpcfg}_i.A[1]$ is set, i.e. the mode is NAPOT, then bits $\text{pmpaddr}_i[G-2:0]$ read as all ones. When $G \geq 1$ and $\text{pmpcfg}_i.A[1]$ is clear, i.e. the mode is OFF or TOR, then bits $\text{pmpaddr}_i[G-1:0]$ read as all zeros. Bits $\text{pmpaddr}_i[G-1:0]$ do not affect the TOR address-matching logic.

Software may determine the PMP granularity by writing zero to `pmp0cfg`, then writing all ones

to `pmpaddr0`, then reading back `pmpaddr0`. If G is the index of the least-significant bit set, the PMP granularity is 2^{G+2} bytes.

Locking and Privilege Mode

The L bit indicates that the PMP entry is locked, i.e., writes to the configuration register and associated address registers are ignored. Locked PMP entries may only be unlocked with a system reset. If PMP entry i is locked, writes to `pmpicfg` and `pmpaddr i` are ignored. Additionally, if `pmpicfg.A` is set to TOR, writes to `pmpaddr $i-1$` are ignored.

In addition to locking the PMP entry, the L bit indicates whether the R/W/X permissions are enforced on M-mode accesses. When the L bit is set, these permissions are enforced for all privilege modes. When the L bit is clear, any M-mode access matching the PMP entry will succeed; the R/W/X permissions apply only to S and U modes.

Priority and Matching Logic

PMP entries are statically prioritized. The lowest-numbered PMP entry that matches any byte of an access determines whether that access succeeds or fails. The matching PMP entry must match all bytes of an access, or the access fails, irrespective of the L, R, W, and X bits. For example, if a PMP entry is configured to match the four-byte range `0xC-0xF`, then an 8-byte access to the range `0x8-0xF` will fail, assuming that PMP entry is the highest-priority entry that matches those addresses.

If a PMP entry matches all bytes of an access, then the L, R, W, and X bits determine whether the access succeeds or fails. If the L bit is clear and the privilege mode of the access is M, the access succeeds. Otherwise, if the L bit is set or the privilege mode of the access is S or U, then the access succeeds only if the R, W, or X bit corresponding to the access type is set.

If no PMP entry matches an M-mode access, the access succeeds. If no PMP entry matches an S-mode or U-mode access, but at least one PMP entry is implemented, the access fails.

Failed accesses generate a load, store, or instruction access exception. Note that a single instruction may generate multiple accesses, which may not be mutually atomic. An access exception is generated if at least one access generated by an instruction fails, though other accesses generated by that instruction may succeed with visible side effects. Notably, instructions that reference virtual memory are decomposed into multiple accesses.

On some implementations, misaligned loads, stores, and instruction fetches may also be decomposed into multiple accesses, some of which may succeed before an access exception occurs. In particular, a portion of a misaligned store that passes the PMP check may become visible, even if another portion fails the PMP check. The same behavior may manifest for floating-point stores wider than XLEN bits (e.g., the FSD instruction in RV32D), even when the store address is naturally aligned.

3.6.2 Physical Memory Protection and Paging

The Physical Memory Protection mechanism is designed to compose with the page-based virtual memory systems described in Chapter 4. When paging is enabled, instructions that access virtual memory may result in multiple physical-memory accesses, including implicit references to the page tables. The PMP checks apply to all of these accesses. The effective privilege mode for implicit page-table accesses is S.

Implementations with virtual memory are permitted to perform address translations speculatively and earlier than required by an explicit virtual-memory access. The PMP settings for the resulting physical address may be checked at any point between the address translation and the explicit virtual-memory access. Hence, when the PMP settings are modified in a manner that affects either the physical memory that holds the page tables or the physical memory to which the page tables point, M-mode software must synchronize the PMP settings with the virtual memory system. This is accomplished by executing an SFENCE.VMA instruction with $rs1=x0$ and $rs2=x0$, after the PMP CSRs are written.

If page-based virtual memory is not implemented, or when it is disabled, memory accesses check the PMP settings synchronously, so no fence is needed.

Chapter 4

Supervisor-Level ISA, Version 1.11-draft

This chapter describes the RISC-V supervisor-level architecture, which contains a common core that is used with various supervisor-level address translation and protection schemes. Supervisor-level code relies on a supervisor execution environment to initialize the environment and enter the supervisor code at an entry point defined by the system binary interface (SBI). The SBI also defines function entry points that provide supervisor environment services for supervisor-level code.

Supervisor mode is deliberately restricted in terms of interactions with underlying physical hardware, such as physical memory and device interrupts, to support clean virtualization.

4.1 Supervisor CSRs

A number of CSRs are provided for the supervisor.

The supervisor should only view CSR state that should be visible to a supervisor-level operating system. In particular, there is no information about the existence (or non-existence) of higher privilege levels (hypervisor or machine) visible in the CSRs accessible by the supervisor.

Many supervisor CSRs are a subset of the equivalent machine-mode CSR, and the machine-mode chapter should be read first to help understand the supervisor-level CSR descriptions.

4.1.1 Supervisor Status Register (`sstatus`)

The `sstatus` register is an SXLEN-bit read/write register formatted as shown in Figure 4.1 for RV32 and Figure 4.2 for RV64 and RV128. The `sstatus` register keeps track of the processor's current operating state.

The SPP bit indicates the privilege level at which a hart was executing before entering supervisor mode. When a trap is taken, SPP is set to 0 if the trap originated from user mode, or 1 otherwise. When an SRET instruction (see Section 3.2.2) is executed to return from the trap handler, the

31	30	20	19	18	17	16	15	14	13	12	9	8	7	6	5	4	3	2	1	0
SD	WPRI	MXR	SUM	WPRI	XS[1:0]	FS[1:0]	WPRI	SPP	WPRI	SPIE	UPIE	WPRI	SIE	UIE						
1	11	1	1	1	2	2	4	1	2	1	1	2	1	1						

Figure 4.1: Supervisor-mode status register (`sstatus`) for RV32.

SXLEN-1	SXLEN-2	34	33	32	31	20	19	18	17											
SD	WPRI	UXL	WPRI	MXR	SUM	WPRI														
1	SXLEN-35	2	12	1	1	1														

16	15	14	13	12	9	8	7	6	5	4	3	2	1	0
XS[1:0]	FS[1:0]	WPRI	SPP	WPRI	SPIE	UPIE	WPRI	SIE	UIE					
2	2	4	1	2	1	1	2	1	1					

Figure 4.2: Supervisor-mode status register (`sstatus`) for RV64 and RV128.

privilege level is set to user mode if the SPP bit is 0, or supervisor mode if the SPP bit is 1; SPP is then set to 0.

The SIE bit enables or disables all interrupts in supervisor mode. When SIE is clear, interrupts are not taken while in supervisor mode. When the hart is running in user-mode, the value in SIE is ignored, and supervisor-level interrupts are enabled. The supervisor can disable individual interrupt sources using the `sie` register.

The SPIE bit indicates whether supervisor interrupts were enabled prior to trapping into supervisor mode. When a trap is taken into supervisor mode, SPIE is set to SIE, and SIE is set to 0. When an SRET instruction is executed, SIE is set to SPIE, then SPIE is set to 1.

The UIE bit enables or disables user-mode interrupts. User-level interrupts are enabled only if UIE is set and the hart is running in user-mode. The UPIE bit indicates whether user-level interrupts were enabled prior to taking a user-level trap. When a URET instruction is executed, UIE is set to UPIE, and UPIE is set to 1. User-level interrupts are optional. If omitted, the UIE and UPIE bits are hardwired to zero.

The `sstatus` register is a subset of the `mstatus` register. In a straightforward implementation, reading or writing any field in `sstatus` is equivalent to reading or writing the homonymous field in `mstatus`.

4.1.2 Base ISA Control in `sstatus` Register

The UXL field controls the value of XLEN for U-mode, termed *UXLEN*, which may differ from the value of XLEN for S-mode, termed *SXLEN*. The encoding of UXL is the same as that of the MXL field of `misa`, shown in Table 3.1.

For RV32 systems, the UXL field does not exist, and UXLEN=32. For RV64 and RV128 systems, it is a **WARL** field that encodes the current value of UXLEN. In particular, the implementation may hardwire UXL so that UXLEN=SXLEN.

If $UXLEN \neq SXLEN$, instructions executed in the narrower mode must ignore source register operand bits above the configured $XLEN$, and must sign-extend results to fill the widest supported $XLEN$ in the destination register.

If $UXLEN < SXLEN$, user-mode instruction-fetch addresses and load and store effective addresses are taken modulo 2^{UXLEN} . For example, when $UXLEN=32$ and $SXLEN=64$, user-mode memory accesses reference the lowest 4 GiB of the address space.

4.1.3 Memory Privilege in sstatus Register

The MXR (Make eXecutable Readable) bit modifies the privilege with which loads access virtual memory. When $MXR=0$, only loads from pages marked readable ($R=1$ in Figure 4.15) will succeed. When $MXR=1$, loads from pages marked either readable or executable ($R=1$ or $X=1$) will succeed. MXR has no effect when page-based virtual memory is not in effect.

The SUM (permit Supervisor User Memory access) bit modifies the privilege with which S-mode loads and stores access virtual memory. When $SUM=0$, S-mode memory accesses to pages that are accessible by U-mode ($U=1$ in Figure 4.15) will fault. When $SUM=1$, these accesses are permitted. SUM has no effect when page-based virtual memory is not in effect, nor when executing in U-mode. Note that S-mode can never execute instructions from user pages, regardless of the state of SUM.

The SUM mechanism prevents supervisor software from inadvertently accessing user memory. Operating systems can execute the majority of code with SUM clear; the few code segments that should access user memory can temporarily set SUM.

The SUM mechanism does not avail S-mode software of permission to execute instructions in user code pages. Legitimate uses cases for execution from user memory in supervisor context are rare in general and nonexistent in POSIX environments. However, bugs in supervisors that lead to arbitrary code execution are much easier to exploit if the supervisor exploit code can be stored in a user buffer at a virtual address chosen by an attacker.

Some non-POSIX single address space operating systems do allow certain privileged software to partially execute in supervisor mode, while most programs run in user mode, all in a shared address space. This use case can be realized by mapping the physical code pages at multiple virtual addresses with different permissions, possibly with the assistance of the instruction page-fault handler to direct supervisor software to use the alternate mapping.

4.1.4 Supervisor Trap Vector Base Address Register (stvec)

The **stvec** register is an $SXLEN$ -bit read/write register that holds trap vector configuration, consisting of a vector base address (BASE) and a vector mode (MODE).



Figure 4.3: Supervisor trap vector base address register (**stvec**).

The BASE field in **stvec** is a **WARL** field that can hold any valid virtual or physical address, subject to the following alignment constraints: the address must always be at least 4-byte aligned,

and the MODE setting may impose additional alignment constraints on the value in the BASE field.

Value	Name	Description
0	Direct	All exceptions set <code>pc</code> to BASE.
1	Vectored	Asynchronous interrupts set <code>pc</code> to BASE+4×cause.
≥2	—	<i>Reserved</i>

Table 4.1: Encoding of `stvec` MODE field.

The encoding of the MODE field is shown in Table 4.1. When MODE=Direct, all traps into supervisor mode cause the `pc` to be set to the address in the BASE field. When MODE=Vectored, all synchronous exceptions into supervisor mode cause the `pc` to be set to the address in the BASE field, whereas interrupts cause the `pc` to be set to the address in the BASE field plus four times the interrupt cause number. For example, a supervisor-mode timer interrupt (see Table 4.2) causes the `pc` to be set to BASE+0x14. Setting MODE=Vectored may impose an additional alignment constraint on BASE, requiring up to 4×SXLEN-byte alignment.

When vectored interrupts are enabled, interrupt cause 0, which corresponds to user-mode software interrupts, are vectored to the same location as synchronous exceptions. This ambiguity does not arise in practice for supervisor software, since user-mode software interrupts are either disabled or delegated to user mode.

4.1.5 Supervisor Interrupt Registers (`sip` and `sie`)

The `sip` register is an SXLEN-bit read/write register containing information on pending interrupts, while `sie` is the corresponding SXLEN-bit read/write register containing interrupt enable bits.

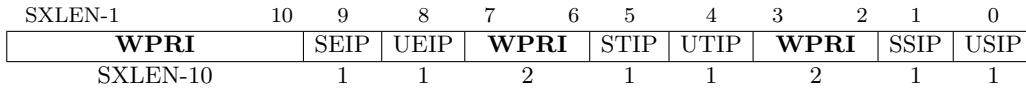


Figure 4.4: Supervisor interrupt-pending register (`sip`).

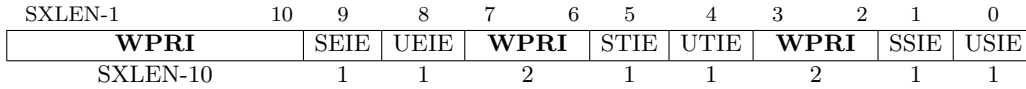


Figure 4.5: Supervisor interrupt-enable register (`sie`).

Three types of interrupts are defined: software interrupts, timer interrupts, and external interrupts. A supervisor-level software interrupt is triggered on the current hart by writing 1 to its supervisor software interrupt-pending (SSIP) bit in the `sip` register. A pending supervisor-level software interrupt can be cleared by writing 0 to the SSIP bit in `sip`. Supervisor-level software interrupts are disabled when the SSIE bit in the `sie` register is clear.

Interprocessor interrupts are sent to other harts by means of SBI calls, which will ultimately cause the SSIP bit to be set in the recipient hart’s `sip` register.

A user-level software interrupt is triggered on the current hart by writing 1 to its user software interrupt-pending (USIP) bit in the `sip` register. A pending user-level software interrupt can be cleared by writing 0 to the USIP bit in `sip`. User-level software interrupts are disabled when the USIE bit in the `sie` register is clear. If user-level interrupts are not supported, USIP and USIE are hardwired to zero.

All bits besides SSIP, USIP, and UEIP in the `sip` register are read-only.

A supervisor-level timer interrupt is pending if the STIP bit in the `sip` register is set. Supervisor-level timer interrupts are disabled when the STIE bit in the `sie` register is clear. An SBI call to the SEE may be used to clear the pending timer interrupt.

A user-level timer interrupt is pending if the UTIP bit in the `sip` register is set. User-level timer interrupts are disabled when the UTIE bit in the `sie` register is clear. If user-level interrupts are supported, the ABI should provide a facility for scheduling timer interrupts in terms of real-time counter values. If user-level interrupts are not supported, UTIP and UTIE are hardwired to zero.

A supervisor-level external interrupt is pending if the SEIP bit in the `sip` register is set. Supervisor-level external interrupts are disabled when the SEIE bit in the `sie` register is clear. The SBI should provide facilities to mask, unmask, and query the cause of external interrupts.

The UEIP field in `sip` contains a single read-write bit. UEIP may be written by S-mode software to indicate to U-mode that an external interrupt is pending. Additionally, the platform-level interrupt controller may generate user-level external interrupts. The logical-OR of the software-writable bit and the signal from the external interrupt controller are used to generate external interrupts for user mode. When the UEIP bit is read with a CSRRW, CSRRS, or CSRRC instruction, the value returned in the `rd` destination register contains the logical-OR of the software-writable bit and the interrupt signal from the interrupt controller. However, the value used in the read-modify-write sequence of a CSRWS or CSRRC instruction is only the software-writable UEIP bit, ignoring the interrupt value from the external interrupt controller.

Analogous to SEIP, the UEIP field behavior is designed to allow a higher privilege layer to mimic external interrupts without losing any real external interrupts.

User-level external interrupts are disabled when the UEIE bit in the `sie` register is clear. If the N extension for user-level interrupts is not implemented, UEIP and UEIE are hardwired to zero.

The `sip` and `sie` registers are subsets of the `mip` and `mie` registers. Reading any field, or writing any writable field, of `sip/sie` effects a read or write of the homonymous field of `mip/mie`.

4.1.6 Supervisor Timers and Performance Counters

Supervisor software uses the same hardware performance monitoring facility as user-mode software, including the `time`, `cycle`, and `instret` CSRs. The SBI should provide a mechanism to modify the counter values.

The SBI must provide a facility for scheduling timer interrupts in terms of the real-time counter, `time`.

4.1.7 Counter-Enable Register (`scounteren`)

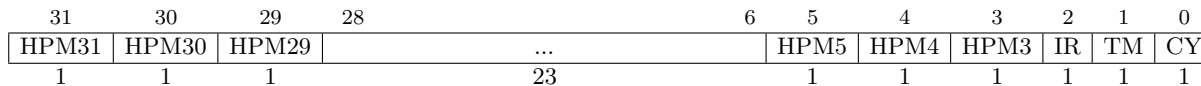


Figure 4.6: Counter-enable register (`scounteren`).

The counter-enable register `scounteren` is a 32-bit register that controls the availability of the hardware performance monitoring counters to U-mode.

When the CY, TM, IR, or HPM n bit in the `scounteren` register is clear, attempts to read the `cycle`, `time`, `instret`, or `hpmcountern` register while executing in U-mode will cause an illegal instruction exception. When one of these bits is set, access to the corresponding register is permitted.

`scounteren` must be implemented. However, any of the bits may contain a hardwired value of zero, indicating reads to the corresponding counter will cause an exception when executing in U-mode. Hence, they are effectively **WARL** fields.

4.1.8 Supervisor Scratch Register (`sscratch`)

The `sscratch` register is an SXLEN-bit read/write register, dedicated for use by the supervisor. Typically, `sscratch` is used to hold a pointer to the hart-local supervisor context while the hart is executing user code. At the beginning of a trap handler, `sscratch` is swapped with a user register to provide an initial working register.

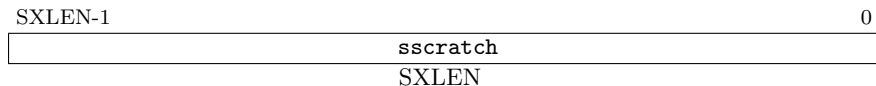


Figure 4.7: Supervisor Scratch Register.

4.1.9 Supervisor Exception Program Counter (`sepc`)

`sepc` is an SXLEN-bit read/write register formatted as shown in Figure 4.8. The low bit of `sepc` (`sepc[0]`) is always zero. On implementations that support only IALIGN=32, the two low bits (`sepc[1:0]`) are always zero.

If an implementation allows IALIGN to be either 16 or 32 (by changing CSR `misa`, for example), then, whenever IALIGN=32, bit `sepc[1]` is masked on reads so that it appears to be 0. This masking occurs also for the implicit read by the SRET instruction. Though masked, `sepc[1]` remains writable when IALIGN=32.

`sepc` is a **WARL** register that must be able to hold all valid physical and virtual addresses. It need not be capable of holding all possible invalid addresses. Implementations may convert some invalid address patterns into other invalid addresses prior to writing them to `sepc`.

When a trap is taken into S-mode, **sepc** is written with the virtual address of the instruction that was interrupted or that encountered the exception. Otherwise, **sepc** is never written by the implementation, though it may be explicitly written by software.

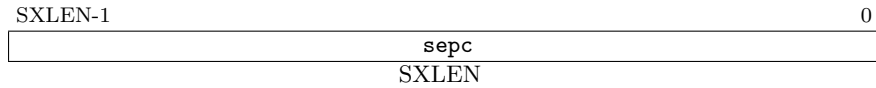


Figure 4.8: Supervisor exception program counter register.

4.1.10 Supervisor Cause Register (**scause**)

The **scause** register is an SXLEN-bit read-write register formatted as shown in Figure 4.9. When a trap is taken into S-mode, **scause** is written with a code indicating the event that caused the trap. Otherwise, **scause** is never written by the implementation, though it may be explicitly written by software.

The Interrupt bit in the **scause** register is set if the trap was caused by an interrupt. The Exception Code field contains a code identifying the last exception. Table 4.2 lists the possible exception codes for the current supervisor ISAs, in descending order of priority. The Exception Code is an **WLRL** field, so is only guaranteed to hold supported exception codes.

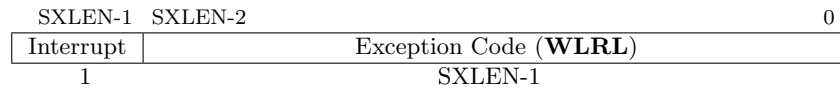


Figure 4.9: Supervisor Cause register **scause**.

4.1.11 Supervisor Trap Value (**stval**) Register

The **stval** register is an SXLEN-bit read-write register formatted as shown in Figure 4.10. When a trap is taken into S-mode, **stval** is written with exception-specific information to assist software in handling the trap. Otherwise, **stval** is never written by the implementation, though it may be explicitly written by software. The hardware platform will specify which exceptions must set **stval** informatively and which may unconditionally set it to zero.

When a hardware breakpoint is triggered, or an instruction-fetch, load, or store address-misaligned, access, or page-fault exception occurs, **stval** is written with the faulting virtual address. On an illegal instruction trap, **stval** may be written with the first XLEN or ILEN bits of the faulting instruction as described below. For other exceptions, **stval** is set to zero, but a future standard may redefine **stval**'s setting for other exceptions.

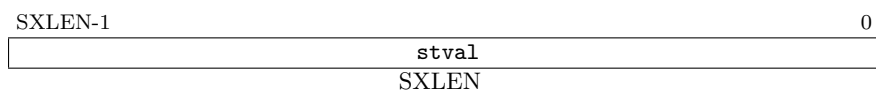


Figure 4.10: Supervisor Trap Value register.

Interrupt	Exception Code	Description
1	0	User software interrupt
1	1	Supervisor software interrupt
1	2–3	<i>Reserved</i>
1	4	User timer interrupt
1	5	Supervisor timer interrupt
1	6–7	<i>Reserved</i>
1	8	User external interrupt
1	9	Supervisor external interrupt
1	≥ 10	<i>Reserved</i>
0	0	Instruction address misaligned
0	1	Instruction access fault
0	2	Illegal instruction
0	3	Breakpoint
0	4	Load address misaligned
0	5	Load access fault
0	6	Store/AMO address misaligned
0	7	Store/AMO access fault
0	8	Environment call from U-mode
0	9	Environment call from S-mode
0	10–11	<i>Reserved</i>
0	12	Instruction page fault
0	13	Load page fault
0	14	<i>Reserved</i>
0	15	Store/AMO page fault
0	≥ 16	<i>Reserved</i>

Table 4.2: Supervisor cause register (**scause**) values after trap.

For misaligned loads and stores that cause access or page-fault exceptions, **stval** will contain the virtual address of the portion of the access that caused the fault. For instruction-fetch access or page-fault exceptions on systems with variable-length instructions, **stval** will contain the virtual address of the portion of the instruction that caused the fault while **sepc** will point to the beginning of the instruction.

The **stval** register can optionally also be used to return the faulting instruction bits on an illegal instruction exception (**sepc** points to the faulting instruction in memory).

If this feature is not provided, then **stval** is set to zero on an illegal instruction fault.

If this feature is provided, after an illegal instruction trap, **stval** will contain the shortest of:

- the actual faulting instruction
- the first ILEN bits of the faulting instruction
- the first XLEN bits of the faulting instruction

The value loaded into **stval** is right-justified and all unused upper bits are cleared to zero.

stval is a **WARL** register that must be able to hold all valid physical and virtual addresses and the value 0. It need not be capable of holding all possible invalid addresses. Implementations may convert some invalid address patterns into other invalid addresses prior to writing them to **stval**. If the feature to return the faulting instruction bits is implemented, **stval** must also be able to hold all values less than 2^N , where N is the smaller of **XLEN** and the width of the longest supported instruction.

4.1.12 Supervisor Address Translation and Protection (**satp**) Register

The **satp** register is an **SXLEN**-bit read/write register, formatted as shown in Figure 4.11 for **SXLEN**=32 and Figure 4.12 for **SXLEN**=64, which controls supervisor-mode address translation and protection. This register holds the physical page number (PPN) of the root page table, i.e., its supervisor physical address divided by 4 KiB; an address space identifier (ASID), which facilitates address-translation fences on a per-address-space basis; and the **MODE** field, which selects the current address-translation scheme. Further details on the access to this register are described in Section 3.1.10.

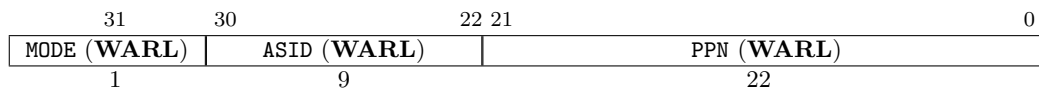


Figure 4.11: RV32 Supervisor address translation and protection register **satp**.

*Storing a PPN in **satp**, rather than a physical address, supports a physical address space larger than 4 GiB for RV32.*

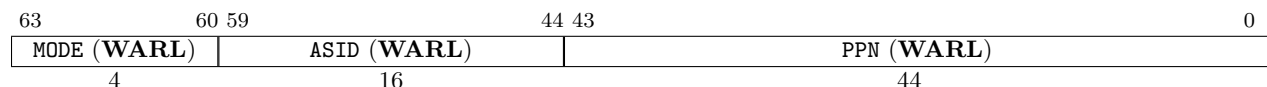


Figure 4.12: RV64 Supervisor address translation and protection register **satp**, for **MODE** values Bare, Sv39, and Sv48.

We store the ASID and the page table base address in the same CSR to allow the pair to be changed atomically on a context switch. Swapping them non-atomically could pollute the old virtual address space with new translations, or vice-versa. This approach also slightly reduces the cost of a context switch.

Table 4.3 shows the encodings of the **MODE** field for RV32 and RV64. When **MODE**=Bare, supervisor virtual addresses are equal to supervisor physical addresses, and there is no additional memory protection beyond the physical memory protection scheme described in Section 3.6. In this case, the remaining fields in **satp** have no effect.

For RV32, the only other valid setting for **MODE** is Sv32, a paged virtual-memory scheme described in Section 4.3.

For RV64, two paged virtual-memory schemes are defined: Sv39 and Sv48, described in Sections 4.4 and 4.5, respectively. Two additional schemes, Sv57 and Sv64, will be defined in a later version of this specification. The remaining **MODE** settings are reserved for future use and may define different interpretations of the other fields in **satp**.

Implementations are not required to support all MODE settings, and if **satp** is written with an unsupported MODE, the entire write has no effect; no fields in **satp** are modified.

RV32		
Value	Name	Description
0	Bare	No translation or protection.
1	Sv32	Page-based 32-bit virtual addressing.
RV64		
Value	Name	Description
0	Bare	No translation or protection.
1–7	—	<i>Reserved</i>
8	Sv39	Page-based 39-bit virtual addressing.
9	Sv48	Page-based 48-bit virtual addressing.
10	<i>Sv57</i>	<i>Reserved for page-based 57-bit virtual addressing.</i>
11	<i>Sv64</i>	<i>Reserved for page-based 64-bit virtual addressing.</i>
12–15	—	<i>Reserved</i>

Table 4.3: Encoding of **satp** MODE field.

The number of supervisor physical address bits is implementation-defined; any unimplemented address bits are hardwired to zero in the **satp** register. The number of ASID bits is also implementation-defined and may be zero. The number of implemented ASID bits, termed *ASIDLEN*, may be determined by writing one to every bit position in the ASID field, then reading back the value in **satp** to see which bit positions in the ASID field hold a one. The least-significant bits of ASID are implemented first: that is, if $ASIDLEN > 0$, $ASID[ASIDLEN-1:0]$ is writable. The maximal value of *ASIDLEN*, termed *ASIDMAX*, is 9 for Sv32 or 16 for Sv39 and Sv48.

For many applications, the choice of page size has a substantial performance impact. A large page size increases TLB reach and loosens the associativity constraints on virtually-indexed, physically-tagged caches. At the same time, large pages exacerbate internal fragmentation, wasting physical memory and possibly cache capacity.

After much deliberation, we have settled on a conventional page size of 4 KiB for both RV32 and RV64. We expect this decision to ease the porting of low-level runtime software and device drivers. The TLB reach problem is ameliorated by transparent superpage support in modern operating systems [2]. Additionally, multi-level TLB hierarchies are quite inexpensive relative to the multi-level cache hierarchies whose address space they map.

Note that writing **satp** does not imply any ordering constraints between page-table updates and subsequent address translations. If the new address space’s page tables have been modified, it may be necessary to execute an SFENCE.VMA instruction (see Section 4.2.1) prior to writing **satp**.

*Not imposing upon implementations to flush address-translation caches upon **satp** writes reduces the cost of context switches, provided a sufficiently large ASID space.*

4.2 Supervisor Instructions

In addition to the SRET instruction defined in Section 3.2.2, one new supervisor-level instruction is provided.

4.2.1 Supervisor Memory-Management Fence Instruction

31	25 24	20 19	15 14	12 11	7 6	0
funct7	rs2	rs1	funct3	rd	opcode	
7	5	5	3	5	7	
SFENCE.VMA	asid	vaddr	PRIV	0	SYSTEM	

The supervisor memory-management fence instruction SFENCE.VMA is used to synchronize updates to in-memory memory-management data structures with current execution. Instruction execution causes implicit reads and writes to these data structures; however, these implicit references are ordinarily not ordered with respect to loads and stores in the instruction stream. Executing an SFENCE.VMA instruction guarantees that any stores in the instruction stream prior to the SFENCE.VMA are ordered before all implicit references subsequent to the SFENCE.VMA. Further details on the behaviour of this instruction are described in Section 3.1.10 and Section 3.6.2.

The SFENCE.VMA is used to flush any local hardware caches related to address translation. It is specified as a fence rather than a TLB flush to provide cleaner semantics with respect to which instructions are affected by the flush operation and to support a wider variety of dynamic caching structures and memory-management schemes. SFENCE.VMA is also used by higher privilege levels to synchronize page table writes and the address translation hardware.

Note the instruction has no effect on the translations of other harts, which must be notified separately. One approach is to use 1) a local data fence to ensure local writes are visible globally, then 2) an interprocessor interrupt to the other thread, then 3) a local SFENCE.VMA in the interrupt handler of the remote thread, and finally 4) signal back to originating thread that operation is complete. This is, of course, the RISC-V analog to a TLB shutdown. Alternatively, implementations might provide direct hardware support for remote TLB invalidation. TLB shutdowns are handled by an SBI call to hide implementation details.

For the common case that the translation data structures have only been modified for a single address mapping (i.e., one page or superpage), *rs1* can specify a virtual address within that mapping to effect a translation fence for that mapping only. Furthermore, for the common case that the translation data structures have only been modified for a single address-space identifier, *rs2* can specify the address space. The behavior of SFENCE.VMA depends on *rs1* and *rs2* as follows:

- If *rs1*=x0 and *rs2*=x0, the fence orders all reads and writes made to any level of the page tables, for all address spaces.
- If *rs1*=x0 and *rs2*≠x0, the fence orders all reads and writes made to any level of the page tables, but only for the address space identified by integer register *rs2*. Accesses to *global* mappings (see Section 4.3.1) are not ordered.
- If *rs1*≠x0 and *rs2*=x0, the fence orders only reads and writes made to the leaf page table entry corresponding to the virtual address in *rs1*, for all address spaces.
- If *rs1*≠x0 and *rs2*≠x0, the fence orders only reads and writes made to the leaf page table entry corresponding to the virtual address in *rs1*, for the address space identified by integer register *rs2*. Accesses to global mappings are not ordered.

When $rs2 \neq x0$, bits SXLEN-1:ASIDMAX of the value held in $rs2$ are reserved for future use and should be zeroed by software and ignored by current implementations. Furthermore, if ASIDLEN < ASIDMAX, the implementation shall ignore bits ASIDMAX-1:ASIDLEN of the value held in $rs2$.

Simpler implementations can ignore the virtual address in $rs1$ and the ASID value in $rs2$ and always perform a global fence.

4.3 Sv32: Page-Based 32-bit Virtual-Memory Systems

When Sv32 is written to the MODE field in the `satp` register (see Section 4.1.12), the supervisor operates in a 32-bit paged virtual-memory system. Sv32 is supported on RV32 systems and is designed to include mechanisms sufficient for supporting modern Unix-based operating systems.

The initial RISC-V paged virtual-memory architectures have been designed as straightforward implementations to support existing operating systems. We have architected page table layouts to support a hardware page-table walker. Software TLB refills are a performance bottleneck on high-performance systems, and are especially troublesome with decoupled specialized coprocessors. An implementation can choose to implement software TLB refills using a machine-mode trap handler as an extension to M-mode.

4.3.1 Addressing and Memory Protection

Sv32 implementations support a 32-bit virtual address space, divided into 4 KiB pages. An Sv32 virtual address is partitioned into a virtual page number (VPN) and page offset, as shown in Figure 4.13. When Sv32 virtual memory mode is selected in the MODE field of the `satp` register, supervisor virtual addresses are translated into supervisor physical addresses via a two-level page table. The 20-bit VPN is translated into a 22-bit physical page number (PPN), while the 12-bit page offset is untranslated. The resulting supervisor-level physical addresses are then checked using any physical memory protection structures (Sections 3.6), before being directly converted to machine-level physical addresses.

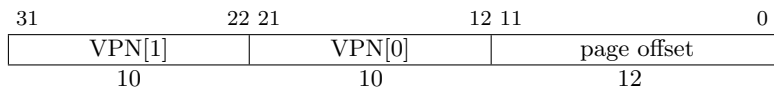


Figure 4.13: Sv32 virtual address.

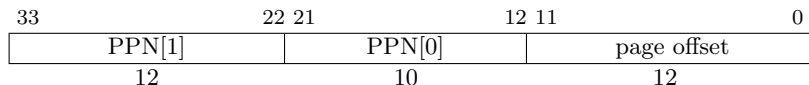


Figure 4.14: Sv32 physical address.

Sv32 page tables consist of 2^{10} page-table entries (PTEs), each of four bytes. A page table is exactly the size of a page and must always be aligned to a page boundary. The physical page number of the root page table is stored in the `satp` register.

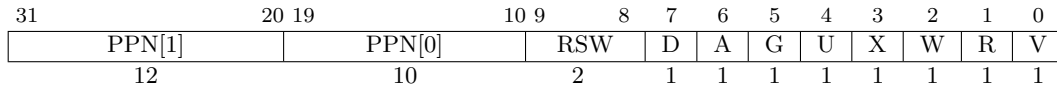


Figure 4.15: Sv32 page table entry.

The PTE format for Sv32 is shown in Figures 4.15. The V bit indicates whether the PTE is valid; if it is 0, all other bits in the PTE are don't-cares and may be used freely by software. The permission bits, R, W, and X, indicate whether the page is readable, writable, and executable, respectively. When all three are zero, the PTE is a pointer to the next level of the page table; otherwise, it is a leaf PTE. Writable pages must also be marked readable; the contrary combinations are reserved for future use. Table 4.4 summarizes the encoding of the permission bits.

X	W	R	Meaning
0	0	0	Pointer to next level of page table.
0	0	1	Read-only page.
0	1	0	<i>Reserved for future use.</i>
0	1	1	Read-write page.
1	0	0	Execute-only page.
1	0	1	Read-execute page.
1	1	0	<i>Reserved for future use.</i>
1	1	1	Read-write-execute page.

Table 4.4: Encoding of PTE R/W/X fields.

Attempting to fetch an instruction from a page that does not have execute permissions raises a fetch page-fault exception. Attempting to execute a load or load-reserved instruction whose effective address lies within a page without read permissions raises a load page-fault exception. Attempting to execute a store, store-conditional (regardless of success), or AMO instruction whose effective address lies within a page without write permissions raises a store page-fault exception.

AMOs never raise load page-fault exceptions. Since any unreadable page is also unwritable, attempting to perform an AMO on an unreadable page always raises a store page-fault exception.

The U bit indicates whether the page is accessible to user mode. U-mode software may only access the page when U=1. If the SUM bit in the `sstatus` register is set, supervisor mode software may also access pages with U=1. However, supervisor code normally operates with the SUM bit clear, in which case, supervisor code will fault on accesses to user-mode pages. Irrespective of SUM, the supervisor may not execute code on pages with U=1.

An alternative PTE format would support different permissions for supervisor and user. We omitted this feature because it would be largely redundant with the SUM mechanism (see Section 4.1.3) and would require more encoding space in the PTE.

The G bit designates a *global* mapping. Global mappings are those that exist in all address spaces. For non-leaf PTEs, the global setting implies that all mappings in the subsequent levels of the page table are global. Note that failing to mark a global mapping as global merely reduces performance, whereas marking a non-global mapping as global is an error.

Global mappings need not be stored redundantly in address-translation caches for multiple ASIDs. Additionally, they need not be flushed from local address-translation caches when an SFENCE.VMA instruction is executed with rs2≠x0.

The RSW field is reserved for use by supervisor software; the implementation shall ignore this field.

Each leaf PTE contains an accessed (A) and dirty (D) bit. The A bit indicates the virtual page has been read, written, or fetched from since the last time the A bit was cleared. The D bit indicates the virtual page has been written since the last time the D bit was cleared.

Two schemes to manage the A and D bits are permitted:

- When a virtual page is accessed and the A bit is clear, or is written and the D bit is clear, the implementation sets the corresponding bit in the PTE. The PTE update must be atomic with respect to other accesses to the PTE, and must atomically check that the PTE is valid and grants sufficient permissions. The PTE update must be exact (i.e., not speculative), and observed in program order by the local hart. The ordering on loads and stores provided by FENCE instructions and the acquire/release bits on atomic instructions also orders the PTE updates associated with those loads and stores as observed by remote harts.
- When a virtual page is accessed and the A bit is clear, or is written and the D bit is clear, a page-fault exception is raised.

All harts in a system must employ the same PTE-update scheme as each other.

Mandating that the PTE updates to be exact, atomic, and in program order simplifies the specification, and makes the feature more useful for system software. Simple implementations may instead generate page-fault exceptions.

The A and D bits are never cleared by the implementation. If the supervisor software does not rely on accessed and/or dirty bits, e.g. if it does not swap memory pages to secondary storage or if the pages are being used to map I/O space, it should always set them to 1 in the PTE to improve performance.

Any level of PTE may be a leaf PTE, so in addition to 4 KiB pages, Sv32 supports 4 MiB *megapages*. A megapage must be virtually and physically aligned to a 4 MiB boundary; a page-fault exception is raised if the physical address is insufficiently aligned.

For non-leaf PTEs, the D, A, and U bits are reserved for future use and must be cleared by software for forward compatibility.

4.3.2 Virtual Address Translation Process

A virtual address va is translated into a physical address pa as follows:

1. Let a be $\text{satp.ppn} \times \text{PAGESIZE}$, and let $i = \text{LEVELS} - 1$. (For Sv32, $\text{PAGESIZE}=2^{12}$ and $\text{LEVELS}=2$.)

2. Let pte be the value of the PTE at address $a + va.vpn[i] \times PTESIZE$. (For Sv32, $PTESIZE=4$.) If accessing pte violates a PMA or PMP check, raise an access exception corresponding to the original access type.
3. If $pte.v = 0$, or if $pte.r = 0$ and $pte.w = 1$, stop and raise a page-fault exception corresponding to the original access type.
4. Otherwise, the PTE is valid. If $pte.r = 1$ or $pte.x = 1$, go to step 5. Otherwise, this PTE is a pointer to the next level of the page table. Let $i = i - 1$. If $i < 0$, stop and raise a page-fault exception corresponding to the original access type. Otherwise, let $a = pte.ppn \times PAGE_SIZE$ and go to step 2.
5. A leaf PTE has been found. Determine if the requested memory access is allowed by the $pte.r$, $pte.w$, $pte.x$, and $pte.u$ bits, given the current privilege mode and the value of the SUM and MXR fields of the `mstatus` register. If not, stop and raise a page-fault exception corresponding to the original access type.
6. If $i > 0$ and $pte.ppn[i - 1 : 0] \neq 0$, this is a misaligned superpage; stop and raise a page-fault exception corresponding to the original access type.
7. If $pte.a = 0$, or if the memory access is a store and $pte.d = 0$, either raise a page-fault exception corresponding to the original access type, or:
 - Set $pte.a$ to 1 and, if the memory access is a store, also set $pte.d$ to 1.
 - If this access violates a PMA or PMP check, raise an access exception corresponding to the original access type.
 - This update and the loading of pte in step 2 must be atomic; in particular, no intervening store to the PTE may be perceived to have occurred in-between.
8. The translation is successful. The translated physical address is given as follows:
 - $pa.pgoff = va.pgoff$.
 - If $i > 0$, then this is a superpage translation and $pa.ppn[i - 1 : 0] = va.vpn[i - 1 : 0]$.
 - $pa.ppn[LEVELS - 1 : i] = pte.ppn[LEVELS - 1 : i]$.

4.4 Sv39: Page-Based 39-bit Virtual-Memory System

This section describes a simple paged virtual-memory system designed for RV64 systems, which supports 39-bit virtual address spaces. The design of Sv39 follows the overall scheme of Sv32, and this section details only the differences between the schemes.

We specified multiple virtual memory systems for RV64 to relieve the tension between providing a large address space and minimizing address-translation cost. For many systems, 512 GiB of virtual-address space is ample, and so Sv39 suffices. Sv48 increases the virtual address space to 256 TiB, but increases the physical memory capacity dedicated to page tables, the latency of page-table traversals, and the size of hardware structures that store virtual addresses.

4.4.1 Addressing and Memory Protection

Sv39 implementations support a 39-bit virtual address space, divided into 4 KiB pages. An Sv39 address is partitioned as shown in Figure 4.16. Instruction fetch addresses and load and store effective addresses, which are 64 bits, **must have bits 63–39 all equal to bit 38, or else a page-fault exception will occur**. The 27-bit VPN is translated into a 44-bit PPN via a three-level page table, while the 12-bit page offset is untranslated.

When mapping between narrower and wider addresses, RISC-V usually zero-extends a narrower address to a wider size. The mapping between 64-bit virtual addresses and the 39-bit usable address space of Sv39 is not based on zero-extension but instead follows an entrenched convention that allows an OS to use one or a few of the most-significant bits of a full-size (64-bit) virtual address to quickly distinguish user and supervisor address regions.

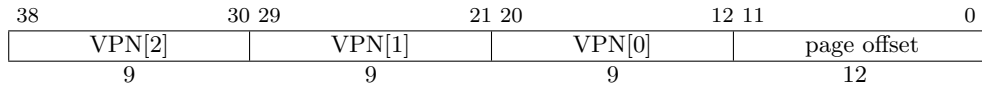


Figure 4.16: Sv39 virtual address.

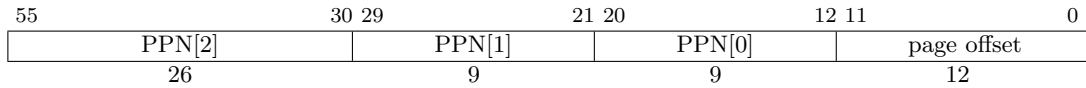


Figure 4.17: Sv39 physical address.

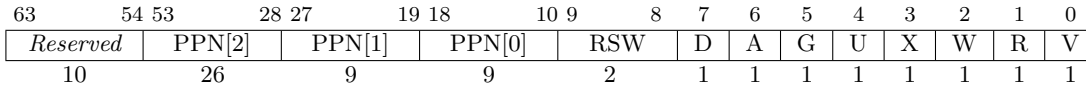


Figure 4.18: Sv39 page table entry.

Sv39 page tables contain 2^9 page table entries (PTEs), eight bytes each. A page table is exactly the size of a page and must always be aligned to a page boundary. The physical page number of the root page table is stored in the **satp** register’s PPN field.

The PTE format for Sv39 is shown in Figure 4.18. Bits 9–0 have the same meaning as for Sv32. Bits 63–54 are reserved for future use and must be zeroed by software for forward compatibility.

We reserved several PTE bits for a possible extension that improves support for sparse address spaces by allowing page-table levels to be skipped, reducing memory usage and TLB refill latency. These reserved bits may also be used to facilitate research experimentation. The cost is reducing the physical address space, but 64 PiB is presently ample. When it no longer suffices, the reserved bits that remain unallocated could be used to expand the physical address space.

Any level of PTE may be a leaf PTE, so in addition to 4 KiB pages, Sv39 supports 2 MiB *megapages* and 1 GiB *gigapages*, each of which must be virtually and physically aligned to a boundary equal to its size. A page-fault exception is raised if the physical address is insufficiently aligned.

The algorithm for virtual-to-physical address translation is the same as in Section 4.3.2, except LEVELS equals 3 and PTE_SIZE equals 8.

4.5 Sv48: Page-Based 48-bit Virtual-Memory System

This section describes a simple paged virtual-memory system designed for RV64 systems, which supports 48-bit virtual address spaces. Sv48 is intended for systems for which a 39-bit virtual address space is insufficient. It closely follows the design of Sv39, simply adding an additional level of page table, and so this chapter only details the differences between the two schemes.

Implementations that support Sv48 must also support Sv39.

Systems that support Sv48 can also support Sv39 at essentially no cost, and so should do so to maintain compatibility with supervisor software that assumes Sv39.

4.5.1 Addressing and Memory Protection

Sv48 implementations support a 48-bit virtual address space, divided into 4 KiB pages. An Sv48 address is partitioned as shown in Figure 4.19. Instruction fetch addresses and load and store effective addresses, which are 64 bits, must have bits 63–48 all equal to bit 47, or else a page-fault exception will occur. The 36-bit VPN is translated into a 44-bit PPN via a four-level page table, while the 12-bit page offset is untranslated.

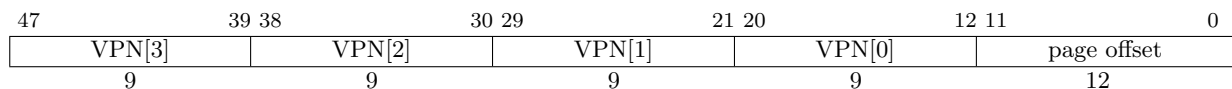


Figure 4.19: Sv48 virtual address.

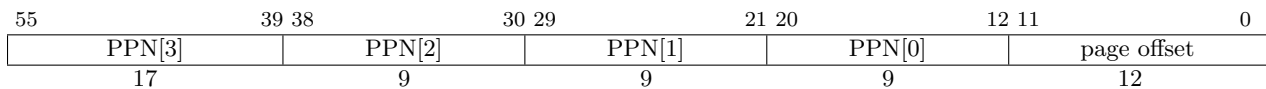


Figure 4.20: Sv48 physical address.

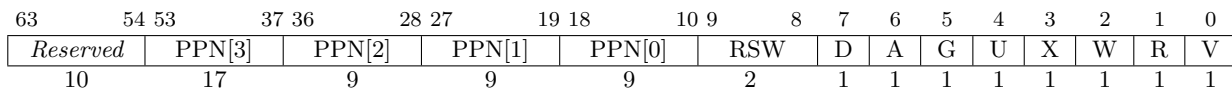


Figure 4.21: Sv48 page table entry.

The PTE format for Sv48 is shown in Figure 4.21. Bits 9–0 have the same meaning as for Sv32. Any level of PTE may be a leaf PTE, so in addition to 4 KiB pages, Sv48 supports 2 MiB *megapages*, 1 GiB *gigapages*, and 512 GiB *terapages*, each of which must be virtually and physically aligned to a boundary equal to its size. A page-fault exception is raised if the physical address is insufficiently aligned.

The algorithm for virtual-to-physical address translation is the same as in Section 4.3.2, except LEVELS equals 4 and PTESIZE equals 8.

Chapter 5

Hypervisor Extension, Version 0.2-draft

This chapter describes the RISC-V hypervisor extension, which virtualizes the supervisor-level architecture to support the efficient hosting of guest operating systems atop a type-1 or type-2 hypervisor. The hypervisor extension changes supervisor mode into *hypervisor-extended supervisor mode* (HS-mode, or *hypervisor mode* for short), where a hypervisor or a hosting-capable operating system runs. The hypervisor extension also adds another level of address translation, from *guest physical addresses* to supervisor physical addresses, to virtualize the memory and memory-mapped I/O subsystems for a guest operating system. HS-mode acts the same as S-mode, but with additional instructions and CSRs that control the new level of address translation and support hosting a guest OS in virtual S-mode (VS-mode). Regular S-mode operating systems can execute without modification either in HS-mode or as VS-mode guests.

In HS-mode, an OS or hypervisor interacts with the machine through the same SBI as an OS normally does from S-mode. An HS-mode hypervisor is expected to implement the SBI for its VS-mode guest.

The hypervisor extension is enabled by setting bit 7 in the `misa` CSR, which corresponds to the letter H. When `misa[7]` is clear, the hart behaves as though this extension were not implemented, and attempts to use hypervisor CSRs or instructions raise an illegal instruction exception. Implementations that include the hypervisor extension are encouraged not to hardwire `misa[7]`, so that the extension may be disabled.

This draft is based on earlier proposals by John Hauser and Paolo Bonzini.

The baseline privileged architecture is designed to simplify the use of classic virtualization techniques, where a guest OS is run at user-level, as the few privileged instructions can be easily detected and trapped. The hypervisor extension improves virtualization performance by reducing the frequency of these traps.

The hypervisor extension has been designed to be efficiently emulable on platforms that do not implement the extension, by running the hypervisor in S-mode and trapping into M-mode for hypervisor CSR accesses and to maintain shadow page tables. The majority of CSR accesses for type-2 hypervisors are valid S-mode accesses so need not be trapped. Hypervisors can support nested virtualization analogously.

5.1 Privilege Modes

The current *virtualization mode*, denoted V , indicates whether the hart is currently executing in a guest. When $V=1$, the hart is either in virtual S-mode (VS-mode), or in virtual U-mode (VU-mode) under a guest OS running in VS-mode. When $V=0$, the hart is either in M-mode, in HS-mode, or in U-mode under an OS running in HS-mode. The virtualization mode also indicates whether two-level address translation is active ($V=1$) or inactive ($V=0$). Table 5.1 lists the possible operating modes of a RISC-V hart with the hypervisor extension.

Virtualization Mode (V)	Privilege Encoding	Abbreviation	Name	Two-Level Translation
0	0	U-mode	User mode	Off
0	1	HS-mode	Hypervisor-extended supervisor mode	Off
0	3	M-mode	Machine mode	Off
1	0	VU-mode	Virtual user mode	On
1	1	VS-mode	Virtual supervisor mode	On

Table 5.1: Operating modes with the hypervisor extension.

5.2 Hypervisor CSRs

An OS or hypervisor running in HS-mode uses the supervisor CSRs to interact with the exception, interrupt, and address-translation subsystems. Additional CSRs are provided to HS-mode, but not to VS-mode, to manage two-level address translation and to control the behavior of a VS-mode guest: `hstatus`, `hedeleg`, `hideleg`, and `hgatep`.

Additionally, several *background* supervisor CSRs are copies of one of the existing *foreground* supervisor CSRs. For example, the `bsstatus` CSR is the background copy of the foreground `sstatus` CSR. When transitioning between virtualization modes ($V=0$ to $V=1$, or vice-versa), the implementation swaps the background supervisor CSRs with their foreground counterparts. When $V=0$, the background supervisor CSRs contain VS-mode’s version of those CSRs, and the foreground supervisor CSRs contain HS-mode’s version. When $V=1$, the background supervisor CSRs contain HS-mode’s version, and the foreground supervisor CSRs contain VS-mode’s version. The background registers are accessible to HS-mode, but not to VS-mode.

The swapping of foreground and background supervisor registers can be implemented either by physically copying bits or by dynamically changing the CSR addresses of hardware registers. The CSR addresses of the background supervisor registers have been aligned with their foreground counterparts to minimize the cost of swapping registers simply by changing their addresses.

In this section, we use the term *HSXLEN* to refer to the effective XLEN when executing in HS-mode.

5.2.1 Hypervisor Status Register (**hstatus**)

The **hstatus** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.1. The **hstatus** register provides facilities analogous to the **mstatus** register that track and control the exception behavior of a VS-mode guest.

HSXLEN-1	23	22	21	20	19	11	10	9	8	7	6	1	0
WPRI	VTSR	VTW	VTVM	WPRI	SPV	STL	SP2P	SP2V	WPRI	SPRV			
HSXLEN-23	1	1	1	9	1	1	1	1	6	1			

Figure 5.1: Hypervisor-mode status register (**hstatus**).

The **hstatus** fields **VTSR**, **VTW**, and **VTVM** are defined analogously to the **mstatus** fields **TSR**, **TW**, and **TVM**, but affect the trapping behavior of the **SRET**, **WFI**, and virtual-memory management instructions only when **V**=1.

The **SPV** bit (Supervisor Previous Virtualization Mode) is written by the implementation whenever a trap is taken into HS-mode. Just as the **SPP** bit in **sstatus** is set to the privilege mode at the time of the trap, the **SPV** bit in **hstatus** is set to the value of the virtualization mode **V** at the time of the trap. When an **SRET** instruction is executed when **V**=0, **V** is set to **SPV**.

When a trap is taken into HS-mode, bits **SP2V** and **SP2P** are set to the values that **SPV** and the HS-level **SPP** had before the trap. (Before the trap, the HS-level **SPP** is **sstatus.SPP** if **V**=0, or **bsstatus.SPP** if **V**=1.) When an **SRET** instruction is executed when **V**=0, the reverse assignments occur: after **SPV** and **sstatus.SPP** have supplied the new virtualization and privilege modes, they are written with **SP2V** and **SP2P**, respectively.

The **STL** bit (Supervisor Translation Level), which indicates which address-translation level caused an access-fault or page-fault exception, is also written by the implementation whenever a trap is taken into HS-mode. On an access or page fault due to guest physical address translation, **STL** is set to 1. For any other trap into HS-mode, **STL** is set to 0.

The **SPRV** bit modifies the privilege with which loads and stores execute when not in M-mode. When **SPRV**=0, translation and protection behave as normal. When **SPRV**=1, load and store memory addresses are translated and protected as though the current virtualization mode were set to **hstatus.SPV** and the current privilege mode were set to the HS-level **SPP** (**sstatus.SPP** when **V**=0, or **bsstatus.SPP** when **V**=1). Table 5.2 enumerates the cases.

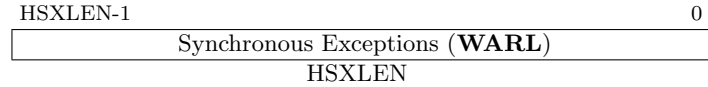
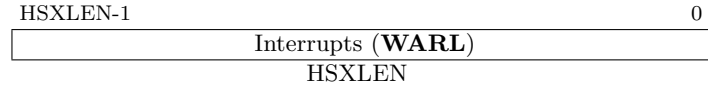
*For simplicity, **SPRV** is in effect even when in U-mode, VS-mode, or VU-mode, but in normal use will only be enabled for short sequences in HS-mode.*

5.2.2 Hypervisor Trap Delegation Registers (**hedeleg** and **hideleg**)

By default, all traps at any privilege level are handled in M-mode, though M-mode usually uses the **medeleg** and **mideleg** CSRs to delegate some traps to HS-mode. The **hedeleg** and **hideleg** CSRs allow these traps to be further delegated to a VS-mode guest; their layout is the same as **medeleg** and **mideleg**.

SPRV	SPV	SPP	Effect
0	–	–	Normal access; current privilege and virtualization modes apply.
1	0	0	U-level access with HS-level translation and protection only.
1	0	1	HS-level access with HS-level translation and protection only.
1	1	0	VU-level access with two-level translation and protection. The HS-level MXR bit makes any executable page readable. The VS-level MXR makes readable those pages marked executable at the VS translation level only if readable at the guest-physical translation level.
1	1	1	VS-level access with two-level translation and protection. The HS-level MXR bit makes any executable page readable. The VS-level MXR makes readable those pages marked executable at the VS translation level only if readable at the guest-physical translation level. The VS-level SUM bit applies instead of HS-level SUM.

Table 5.2: Effect on load and store translation and protection under SPRV.

Figure 5.2: Hypervisor Exception Delegation Register **hedeleg**.Figure 5.3: Hypervisor Interrupt Delegation Register **hideleg**.

The **hedeleg** and **hideleg** registers are only active when $V=1$. When $V=1$, a trap that has been delegated to HS-mode (using **medeleg** or **mideleg**) is further delegated to VS-mode if the corresponding **hedeleg** or **hideleg** bit is set. If the N extension for user-mode interrupts is implemented, the VS-mode guest may further delegate the interrupt to VU-mode by setting the corresponding bit in **sedeleg** or **sideleg**.

When $V=0$ and the N extension for user-mode interrupts is implemented, a trap that has been delegated to HS-mode can be further delegated to U-mode by setting the corresponding bit in **sedeleg** or **sideleg**.

When an access-fault or page-fault exception is caused by guest physical address translation, the trap is not delegated beyond HS-mode, regardless of the setting of **hedeleg**.

5.2.3 Hypervisor Guest Address Translation and Protection Register (**hgap**)

The **hgap** register is an HSXLEN-bit read/write register, formatted as shown in Figure 5.4 for HSXLEN=32 and Figure 5.5 for HSXLEN=64, which controls guest physical address translation and protection. Similar to CSR **satp**, this register holds the physical page number (PPN) of the guest-physical root page table; a virtual machine identifier (VMID), which facilitates address-

translation fences on a per-virtual-machine basis; and the MODE field, which selects the address-translation scheme for guest physical addresses. When `mstatus.TVM=1`, attempts to read or write `hgatp` while executing in HS-mode will raise an illegal instruction exception.

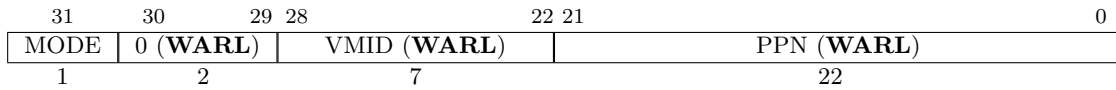


Figure 5.4: RV32 Hypervisor guest address translation and protection register `hgatp`.

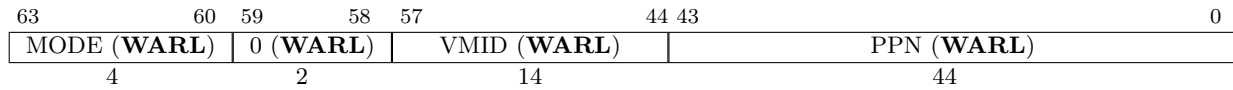


Figure 5.5: RV64 Hypervisor guest address translation and protection register `hgatp`, for MODE values Bare, Sv39x4, and Sv48x4.

Table 5.3 shows the encodings of the MODE field for RV32 and RV64. When MODE=Bare, guest physical addresses are equal to supervisor physical addresses, and there is no further memory protection for a guest virtual machine beyond the physical memory protection scheme described in Section 3.6. In this case, the remaining fields in `hgatp` have no effect.

For RV32, the only other valid setting for MODE is Sv32x4, which is a modification of the usual Sv32 paged virtual-memory scheme, extended to support 34-bit guest physical addresses. For RV64, modes Sv39x4 and Sv48x4 are defined as modifications of the Sv39 and Sv48 paged virtual-memory schemes. All these paged virtual-memory schemes are described in Section 5.5.1. An additional RV64 scheme, Sv57x4, may be defined in a later version of this specification.

The remaining MODE settings for RV64 are reserved for future use and may define different interpretations of the other fields in `hgatp`.

RV64 implementations are not required to support all defined RV64 MODE settings. (However, a write to `hgatp` with an unsupported MODE value is not ignored as it is for `satp`.)

RV32		
Value	Name	Description
0	Bare	No translation or protection.
1	Sv32x4	Page-based 34-bit virtual addressing (2-bit extension of Sv32).
RV64		
Value	Name	Description
0	Bare	No translation or protection.
1–7	—	<i>Reserved</i>
8	Sv39x4	Page-based 41-bit virtual addressing (2-bit extension of Sv39).
9	Sv48x4	Page-based 50-bit virtual addressing (2-bit extension of Sv48).
10	<i>Sv57x4</i>	<i>Reserved for page-based 59-bit virtual addressing.</i>
11–15	—	<i>Reserved</i>

Table 5.3: Encoding of `hgatp` MODE field.

As explained in Section 5.5.1, for the paged virtual-memory schemes (Sv32x4, Sv39x4, and Sv48x4),

the root page table is 16 KiB and must be aligned to a 16-KiB boundary. In these modes, the lowest two bits of the physical page number (PPN) in **hgatp** are ignored. An implementation that supports only the defined paged virtual-memory schemes and/or Bare may hardwire PPN[1:0] to zero.

The number of supervisor physical address bits is implementation-defined; any unimplemented address bits are hardwired to zero in **hgatp.PPN**. The number of VMID bits is also implementation-defined and may be zero. The number of implemented VMID bits, termed *VMIDLEN*, may be determined by writing one to every bit position in the VMID field, then reading back the value in **hgatp** to see which bit positions in the VMID field hold a one. The least-significant bits of VMID are implemented first: that is, if $VMIDLEN > 0$, $VMID[VMIDLEN-1:0]$ is writable. The maximal value of VMIDLEN, termed VMIDMAX, is 7 for Sv32x4 or 14 for Sv39x4 and Sv48x4.

Note that writing **hgatp** does not imply any ordering constraints between page-table updates and subsequent guest physical address translations. If the new virtual machine's guest physical page tables have been modified, it may be necessary to execute an HFENCE.GVMA instruction (see Section 5.3.1) before or after writing **hgatp**.

5.2.4 Background Supervisor Status Register (**bsstatus**)

The **bsstatus** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.6. When $V=0$, the **bsstatus** register holds VS-mode's version of several fields of the **sstatus** register: UXL, MXR, SUM, FS, SPP, SPIE, and SIE. When $V=1$, **bsstatus** holds HS-mode's version of these fields. When transitioning between virtualization modes ($V=0$ to $V=1$, or vice-versa), the implementation swaps these fields in **bsstatus** with their counterparts in **sstatus**. The other fields in **sstatus** are unchanged.

When $V=1$, both **bsstatus.FS** and **sstatus.FS** are in effect. Attempts to execute a floating-point instruction when either field is 0 (Off) raise an illegal-instruction exception. Modifying the floating-point state when $V=1$ causes both fields to be set to 3 (Dirty).

When $V=0$, **bsstatus** does not directly affect the behavior of the machine, unless the MPRV feature in the **mstatus** register or the SPRV feature in the **hstatus** register is used to execute a load or store *as though* $V=1$.

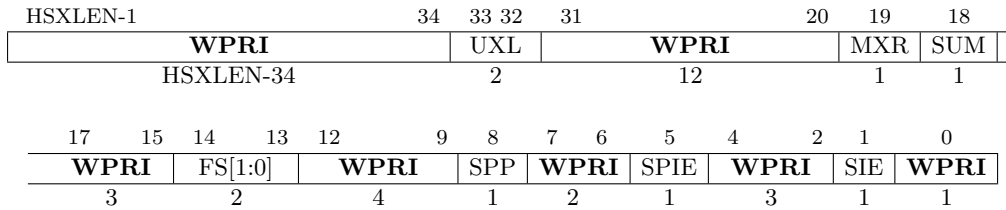


Figure 5.6: Background supervisor status register (**bsstatus**) for RV64 and RV128.

5.2.5 Background Supervisor Interrupt Registers (bsip and bsie)

The **bsip** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.7. When $V=0$, the **bsip** register holds VS-mode's version of the **sip** register. When $V=1$, **bsip** holds HS-mode's version of the **sip** register. When transitioning between virtualization modes ($V=0$ to $V=1$, or vice-versa), the implementation swaps the defined fields of **bsip** with their counterparts in **sip**. The other fields in **sip** are unchanged.

[NOTE: AW: Need to describe how **bsip.SEIP** interacts with PLIC. I think **bsip.SEIP** should purely be a read-write storage bit to emulate the PLIC for VS-mode; the PLIC should not be wired into **bsip.SEIP**.]

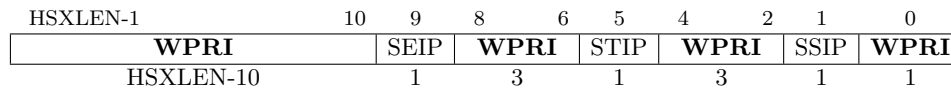


Figure 5.7: Background supervisor interrupt-pending register (**bsip**).

The **bsie** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.8. When $V=0$, the **bsie** register holds VS-mode's version of the **sie** register. When $V=1$, **bsie** holds HS-mode's version of the **sie** register. When transitioning between virtualization modes ($V=0$ to $V=1$, or vice-versa), the implementation swaps the defined fields of **bsie** with their counterparts in **sie**. The other fields in **sie** are unchanged.

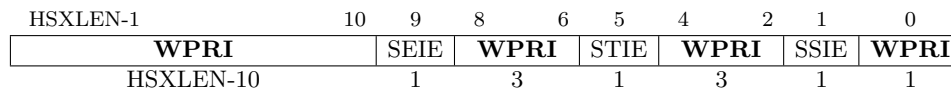


Figure 5.8: Background supervisor interrupt-enable register (**bsie**).

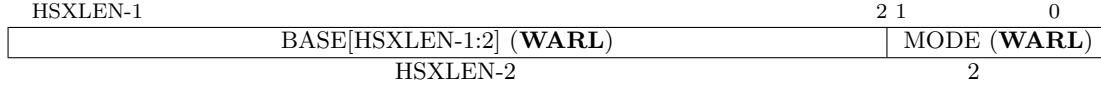
When $V=0$, **bsip** and **bsie** do not affect the behavior of the machine. When $V=1$, they hold the active interrupt-pending and interrupt-enable bits, respectively, for HS-mode; if any bit position holds a 1 in both registers, an interrupt will be taken.

*The **bsip** and **bsie** CSRs do not hold copies of the user-mode interrupt fields. The expectation is that the context-switch code will swap the **uip** and **uie** CSRs along with the other user-mode interrupt registers (**ustatus**, **utvec**, etc.) if that feature is enabled.*

5.2.6 Background Supervisor Trap Vector Base Address Register (bstvec)

The **bstvec** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.9. When $V=0$, the **bstvec** register holds VS-mode's version of the **stvec** register. When $V=1$, **bstvec** holds HS-mode's version of the **stvec** register. When transitioning between virtualization modes ($V=0$ to $V=1$, or vice-versa), the implementation swaps the contents of **bstvec** and **stvec**.

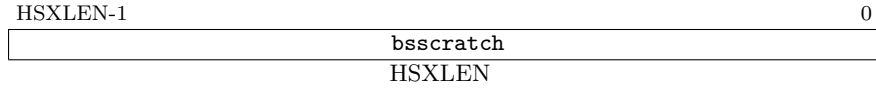
When $V=0$, **bstvec** does not directly affect the behavior of the machine. When $V=1$, it controls the value to which the pc will be set upon a trap into HS-mode.

Figure 5.9: Background supervisor trap vector base address register (**bstvec**).

5.2.7 Background Supervisor Scratch Register (**bsscratch**)

The **bsscratch** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.10. When V=0, the **bsscratch** register holds VS-mode’s version of the **sscratch** register. When V=1, **bsscratch** holds HS-mode’s version of the **sscratch** register. When transitioning between virtualization modes (V=0 to V=1, or vice-versa), the implementation swaps the contents of **bsscratch** and **sscratch**.

Typically, **bsscratch** is used to hold a pointer to the hart-local hypervisor context (when V=1) or supervisor context (when V=0). The contents of **bsscratch** do not directly affect the behavior of the machine.

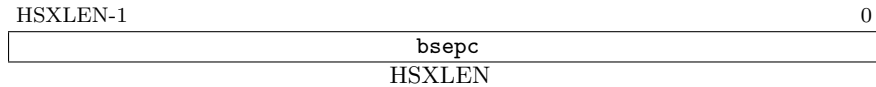
Figure 5.10: Background supervisor scratch register (**bsscratch**).

5.2.8 Background Supervisor Exception Program Counter (**bsepc**)

The **bsepc** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.11. When V=0, the **bsepc** register holds VS-mode’s version of the **sepc** register. When V=1, **bsepc** holds HS-mode’s version of the **sepc** register. When transitioning between virtualization modes (V=0 to V=1, or vice-versa), the implementation swaps the contents of **bsepc** and **sepc**.

The contents of **bsepc** do not directly affect the behavior of the machine.

bsepc is a **WARL** register that must be able to hold the same set of values that **sepc** can hold.

Figure 5.11: Background supervisor exception program counter (**bsepc**).

5.2.9 Background Supervisor Cause Register (**bscause**)

The **bscause** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.12. When V=0, the **bscause** register holds VS-mode’s version of the **scause** register. When V=1, **bscause** holds HS-mode’s version of the **scause** register. When transitioning between virtualization modes (V=0 to V=1, or vice-versa), the implementation swaps the contents of **bscause** and **scause**.

The contents of **bscause** do not directly affect the behavior of the machine.

bscause is a **WLRL** register that must be able to hold the same set of values that **scause** can hold.

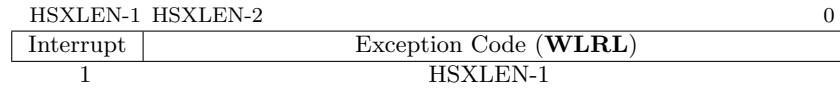


Figure 5.12: Background supervisor cause register (**bscause**).

5.2.10 Background Supervisor Trap Value Register (**bstval**)

The **bstval** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.13. When $V=0$, the **bstval** register holds VS-mode’s version of the **stval** register. When $V=1$, **bstval** holds HS-mode’s version of the **stval** register. When transitioning between virtualization modes ($V=0$ to $V=1$, or vice-versa), the implementation swaps the contents of **bstval** and **stval**.

The contents of **bstval** do not directly affect the behavior of the machine.

bstval is a **WARL** register that must be able to hold the same set of values that **stval** can hold.

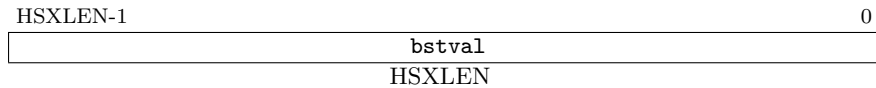


Figure 5.13: Background supervisor trap value register (**bstval**).

5.2.11 Background Supervisor Address Translation and Protection Register (**bsatp**)

The **bsatp** register is an HSXLEN-bit read/write register formatted as shown in Figure 5.14 for RV32 and Figure 5.15 for RV64. When $V=0$, the **bsatp** register holds VS-mode’s version of the **satp** register. When $V=1$, **bsatp** holds HS-mode’s version of the **satp** register. When transitioning between virtualization modes ($V=0$ to $V=1$, or vice-versa), the implementation swaps the contents of **bsatp** and **satp**.

bsatp does not directly affect the behavior of the machine, unless the MPRV feature in the **mstatus** register or the SPRV feature in the **hstatus** register is used to execute a load or store *as though* V is the opposite of its actual setting. The interpretation of the MODE, ASID, and PPN fields is the same as for **satp**.

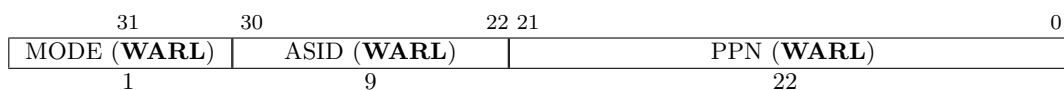


Figure 5.14: RV32 background supervisor address translation and protection register **bsatp**.

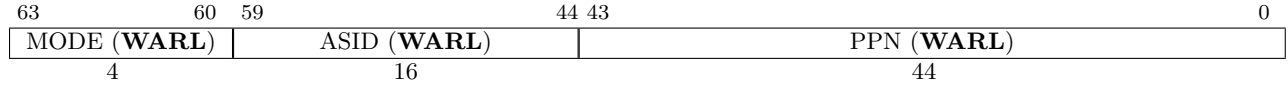


Figure 5.15: RV64 background supervisor address translation and protection register **bsatp**, for MODE values Bare, Sv39, and Sv48.

5.3 Hypervisor Instructions

The hypervisor extension adds two new privileged fence instructions.

5.3.1 Hypervisor Memory-Management Fence Instructions

31	25	24	20	19	15	14	12	11	7	6	0
funct7			rs2		rs1		funct3		rd		opcode
7			5		5		3		5		7
HFENCE.GVMA			vmid		gaddr		PRIV		0		SYSTEM
HFENCE.BVMA			asid		vaddr		PRIV		0		SYSTEM

The hypervisor memory-management fence instructions, HFENCE.GVMA and HFENCE.BVMA, are valid only in HS-mode when **mstatus.TVM**=0, or in M-mode (irrespective of **mstatus.TVM**). These instructions perform a function similar to SFENCE.VMA (Section 4.2.1), except applying to the guest-physical memory-management data structures controlled by CSR **hgatp** (HFENCE.GVMA) or the VS-level memory-management data structures controlled by CSR **bsatp** (HFENCE.BVMA). Instruction SFENCE.VMA applies only to the memory-management data structures controlled by the foreground **satp**.

If an HFENCE.BVMA instruction executes without trapping, its effect is much the same as temporarily entering VS-mode (with the usual swapping of foreground and background supervisor registers) and executing SFENCE.VMA. Executing an HFENCE.BVMA guarantees that any stores in the instruction stream prior to the HFENCE.BVMA are ordered before implicit references to VS-level memory-management data structures when those implicit references

- are subsequent to the HFENCE.BVMA, and
- occur when **hgatp.VMID** has the same setting as it did when HFENCE.BVMA executed.

Implicit references need not be ordered when **hgatp.VMID** is different than at the time HFENCE.BVMA executed. If operand **rs1**≠x0, it specifies a single guest virtual address, and if operand **rs2**≠x0, it specifies a single guest address-space identifier (ASID).

*An HFENCE.BVMA instruction applies only to a single virtual machine, identified by the setting of **hgatp.VMID** when HFENCE.BVMA executes.*

When **rs2**≠x0, bits XLEN-1:ASIDMAX of the value held in **rs2** are reserved for future use and should be zeroed by software and ignored by current implementations. Furthermore, if ASIDLEN < ASIDMAX, the implementation shall ignore bits ASIDMAX-1:ASIDLEN of the value held in **rs2**.

Simpler implementations of HFENCE.BVMA can ignore the guest virtual address in rs1 and the guest ASID value in rs2, as well as hgatp.VMID, and always perform a global fence for the VS-level memory management of all virtual machines, or even a global fence for all memory-management data structures.

Executing an HFENCE.GVMA instruction guarantees that any stores in the instruction stream prior to the HFENCE.GVMA are ordered before all implicit references to guest-physical memory-management data structures subsequent to the HFENCE.GVMA. If operand $rs1 \neq x0$, it specifies a single guest physical address, shifted right by 2 bits, and if operand $rs2 \neq x0$, it specifies a single virtual machine identifier (VMID).

For HFENCE.GVMA, a guest physical address specified in rs1 is shifted right by 2 bits to accommodate addresses wider than the current XLEN. For RV32, the hypervisor extension permits guest physical addresses as wide as 34 bits, and rs1 specifies bits 33:2 of such an address. This shift-by-2 encoding of guest physical addresses matches the encoding of physical addresses in PMP address registers (Section 3.6) and in page table entries (Sections 4.3, 4.4, and 4.5).

When $rs2 \neq x0$, bits XLEN-1:VMIDMAX of the value held in $rs2$ are reserved for future use and should be zeroed by software and ignored by current implementations. Furthermore, if VMIDLEN < VMIDMAX, the implementation shall ignore bits VMIDMAX-1:VMIDLEN of the value held in $rs2$.

Simpler implementations of HFENCE.GVMA can ignore the guest physical address in rs1 and the VMID value in rs2 and always perform a global fence for the guest-physical memory management of all virtual machines, or even a global fence for all memory-management data structures.

5.4 Machine-Level CSRs

The hypervisor extension augments the `mstatus` CSR.

5.4.1 Machine Status Register (`mstatus`)

The hypervisor extension adds two fields to the machine-mode `mstatus` CSR, MPV and MTL, and modifies the behavior of several existing fields. Figure 5.16 shows the `mstatus` register when the hypervisor extension is provided.

MXLEN-1	MXLEN-2	36	35	34	33	32	31	23	22	21	20	19	18	17		
SD	WPRI		SXL[1:0]	UXL[1:0]		WPRI		TSR	TW	TVM	MXR	SUM	MPRV			
1	MXLEN-37		2	2		9		1	1	1	1	1	1			
16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
XS[1:0]	FS[1:0]	MPP[1:0]	MPV	MTL	SPP	MPIE	WPRI	SPIE	UPIE	MIE	WPRI	SIE	UIE			
2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1	

Figure 5.16: Machine-mode status register (`mstatus`) for RV64 and RV128.

The MPV bit (Machine Previous Virtualization Mode) is written by the implementation whenever a trap is taken into M-mode. Just as the MPP bit is set to the privilege mode at the time of the trap, the MPV bit is set to the value of the virtualization mode V at the time of the trap. When an MRET instruction is executed, the virtualization mode V is set to MPV, unless MPP=3, in which case V remains 0.

The MTL bit (Machine Translation Level), which indicates which address-translation level caused an access-fault or page-fault exception, is also written by the implementation whenever a trap is taken into M-mode. On an access or page fault due to guest physical address translation, MTL is set to 1. For any other trap into M-mode, MTL is set to 0.

The SXL field controls the value of XLEN for HS-mode. The UXL field controls the value of XLEN for VS-mode or U-mode when V=0, or for VU-mode when V=1.

The TSR and TVM fields only affect execution in HS-mode.

The TW field affects execution in all modes except M-mode.

The hypervisor extension changes the behavior of the the Modify Privilege field, MPRV. When MPRV=0, translation and protection behave as normal. When MPRV=1, loads and stores are translated and protected as though the current privilege mode were set to MPP and the current virtualization mode were set to MPV. Table 5.4 enumerates the cases.

MPRV	MPV	MPP	Effect
0	–	–	Normal access; current privilege and virtualization modes apply.
1	0	0	U-level access with HS-level translation and protection only.
1	0	1	HS-level access with HS-level translation and protection only.
1	–	3	M-level access with no translation.
1	1	0	VU-level access with two-level translation and protection. The HS-level MXR bit makes any executable page readable. The VS-level MXR makes readable those pages marked executable at the VS translation level only if readable at the guest-physical translation level.
1	1	1	VS-level access with two-level translation and protection. The HS-level MXR bit makes any executable page readable. The VS-level MXR makes readable those pages marked executable at the VS translation level only if readable at the guest-physical translation level. The VS-level SUM bit applies instead of HS-level SUM.

Table 5.4: Effect on load and store translation and protection under MPRV. When MPRV=1, MPP≠3, and `hstatus.SPRV`=1, the effective privilege is further modified: `hstatus.SPV` applies instead of MPV, and the HS-level SPP applies instead of MPP.

The `mstatus` register is a superset of the foreground `sstatus` register; modifying a field in `sstatus` modifies the homonymous field in `mstatus`, and vice-versa.

5.5 Two-Level Address Translation

Whenever the current virtualization mode `V` is 1 (and assuming both `mstatus.MPRV=0` and `hstatus.SPRV=0`), two-level address translation and protection is in effect. For any virtual memory access, the original virtual address is first converted by VS-level address translation, as controlled by the VS-level `satp` register, into a *guest physical address*. The guest physical address is then converted by guest physical address translation, as controlled by the `hgatp` register, into a supervisor physical address. Although there is no option to disable two-level address translation when `V=1`, either level of translation can be effectively disabled by zeroing the corresponding `satp` or `hgatp` register.

The VS-level MXR setting, which makes execute-only pages readable, only overrides VS-level page protection. Setting MXR at VS-level does not override guest-physical page protections. Setting MXR at HS-level, however, overrides both VS-level and guest-physical execute-only permissions.

When `V=1`, memory accesses that would normally bypass address translation are subject to guest physical address translation alone. This includes memory accesses made in support of VS-level address translation, such as reads and writes of VS-level page tables.

Machine-level physical memory protection applies to supervisor physical addresses and is in effect regardless of virtualization mode.

5.5.1 Guest Physical Address Translation

The mapping of guest physical addresses to supervisor physical addresses is controlled by CSR `hgatp` (Section 5.2.3).

When the address translation scheme selected by the `MODE` field of `hgatp` is Bare, guest physical addresses are equal to supervisor physical addresses without modification, and no memory protection applies in the trivial translation of guest physical addresses to supervisor physical addresses.

When `hgatp.MODE` specifies a translation scheme of Sv32x4, Sv39x4, or Sv48x4, guest physical address translation is a variation on the usual page-based virtual address translation scheme of Sv32, Sv39, or Sv48, respectively. In each case, the size of the incoming address is widened by 2 bits (to 34, 41, or 50 bits). To accommodate the 2 extra bits, the root page table (only) is expanded by a factor of four to be 16 KiB instead of the usual 4 KiB. Matching its larger size, the root page table also must be aligned to a 16 KiB boundary instead of the usual 4 KiB page boundary. Except as noted, all other aspects of Sv32, Sv39, or Sv48 are adopted unchanged for guest physical address translation. Non-root page tables and all page table entries (PTEs) have the same formats as documented in Sections 4.3, 4.4, and 4.5.

For Sv32x4, an incoming guest physical address is partitioned into a virtual page number (VPN) and page offset as shown in Figure 5.17. This partitioning is identical to that for an Sv32 virtual address as depicted in Figure 4.13 (page 62), except with 2 more bits at the high end in `VPN[1]`. (Note that the fields of a partitioned guest physical address also correspond one-for-one with the structure that Sv32 assigns to a physical address, depicted in Figure 4.14.)

For Sv39x4, an incoming guest physical address is partitioned as shown in Figure 5.18. This

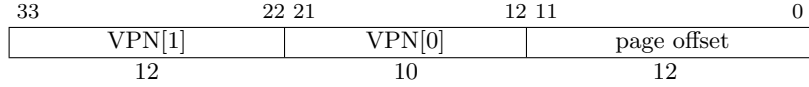


Figure 5.17: Sv32x4 virtual address (guest physical address).

partitioning is identical to that for an Sv39 virtual address as depicted in Figure 4.16 (page 66), except with 2 more bits at the high end in VPN[2]. Address bits 63:41 must all be zeros, or else a page-fault exception occurs, attributed to guest physical address translation.

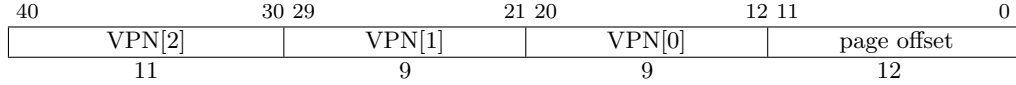


Figure 5.18: Sv39x4 virtual address (guest physical address).

For Sv48x4, an incoming guest physical address is partitioned as shown in Figure 5.19. This partitioning is identical to that for an Sv48 virtual address as depicted in Figure 4.19 (page 67), except with 2 more bits at the high end in VPN[3]. Address bits 63:50 must all be zeros, or else a page-fault exception occurs, attributed to guest physical address translation.

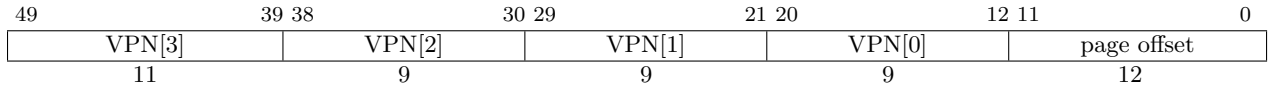


Figure 5.19: Sv48x4 virtual address (guest physical address).

The page-based guest physical address translation scheme for RV32, Sv32x4, is defined to support a 34-bit guest physical address so that an RV32 hypervisor need not be limited in its ability to virtualize real 32-bit RISC-V machines, even those with 33-bit or 34-bit physical addresses. This may include the possibility of a machine virtualizing itself, if it happens to use 33-bit or 34-bit physical addresses. Multiplying the size and alignment of the root page table by a factor of four is the cheapest way to extend Sv32 to cover a 34-bit address. The possible wastage of 12 KiB for an unnecessarily large root page table is expected to be of negligible consequence for most (maybe all) real uses.

A consistent ability to virtualize machines having as much as four times the physical address space as virtual address space is believed to be of some utility also for RV64. For a machine supporting 39-bit virtual addresses (Sv39), for example, this allows the hypervisor extension to support up to a 41-bit physical address space without either necessitating hardware support for 48-bit virtual addresses (Sv48) or falling back to emulating the larger address space with shadow page tables.

The conversion of an Sv32x4, Sv39x4, or Sv48x4 guest physical address is accomplished with the same algorithm used for Sv32, Sv39, or Sv48, as presented in Section 4.3.2, except that:

- in step 1, $a = \text{hgatp.PPN} \times \text{PAGESIZE}$;
- the current privilege mode is always taken to be U-mode; and
- instead of `mstatus.MXR`, the HS-level MXR applies (usually in `bsstatus`).

For guest physical address translation, all memory accesses (including those made to access data structures for VS-level address translation) are considered to be user-level accesses, as though executed in U-mode. Access type permissions—readable, writable, or executable—are checked

during guest physical address translation the same as for VS-level address translation. For a memory access made to support VS-level address translation (such as to read/write a VS-level page table), permissions are checked as though for a load or store, not for the original access type. However, any exception is always reported for the original access type (instruction, load, or store/AMO).

Access faults and page faults raised by guest physical address translation are treated as HS-level exceptions for the purpose of exception delegation, so are not delegated to VS-mode, regardless of the setting of the `hedeleg` register.

5.5.2 Memory-Management Fences

The behavior of the SFENCE.VMA instruction is affected by the current virtualization mode V. When V=0, the virtual-address argument is an HS-level virtual address, and the ASID argument is an HS-level ASID. The instruction orders stores only to HS-level address-translation structures with subsequent HS-level address translations.

When V=1, the virtual-address argument to SFENCE.VMA is a guest virtual address within the current virtual machine, and the ASID argument is a VS-level ASID within the current virtual machine. The current virtual machine is identified by the VMID field of CSR `hgatp`, and the effective ASID can be considered to be the combination of this VMID with the VS-level ASID. The SFENCE.VMA instruction orders stores only to the VS-level address-translation structures with subsequent VS-level address translations for the same virtual machine, i.e., only when `hgatp.VMID` is the same as when the SFENCE.VMA executed.

Hypervisor instructions HFENCE.GVMA and HFENCE.BVMA provide additional memory-management fences to complement SFENCE.VMA. These instructions are described in Section 5.3.1.

Section 3.6.2 discusses the intersection between physical memory protection (PMP) and page-based address translation. It is noted there that, when PMP settings are modified in a manner that affects either the physical memory that holds page tables or the physical memory to which page tables point, M-mode software must synchronize the PMP settings with the virtual memory system. For HS-level address translation, this is accomplished by executing in M-mode an SFENCE.VMA instruction with `rs1=x0` and `rs2=x0`, after the PMP CSRs are written. If guest physical address translation is in use, synchronization with its data structures is also needed. When PMP settings are modified in a manner that affects either the physical memory that holds guest-physical page tables or the physical memory to which guest-physical page tables point, an HFENCE.GVMA instruction with `rs1=x0` and `rs2=x0` must be executed in M-mode after the PMP CSRs are written. An HFENCE.BVMA instruction is not required.

5.6 Base ISA Control

The `mstatus` field `SXL` determines `XLEN` for HS-mode.

When executing in VS-mode, `XLEN` is determined by the `UXL` field of the background register `bsstatus`. Because `bsstatus` is swapped with `sstatus` when transitioning from VS-mode into HS-mode or M-mode, HS-mode and M-mode can modify VS-mode's `XLEN` via the `UXL` field of the foreground register `sstatus`.

When executing in U-mode or VU-mode, `XLEN` is determined by the `UXL` field of the foreground register `sstatus`.

HS-mode controls U-mode's `XLEN` the same way it controls VS-mode's `XLEN`, via `sstatus.UXL`.

5.7 Traps

The hypervisor extension augments the environment-call exception cause encoding. Environment calls from HS-mode use cause 9, whereas environment calls from VS-mode use cause 10. Table 5.5 lists the possible M-mode and HS-mode exception codes when the hypervisor extension is present.

HS-mode and VS-mode ECALLs use different cause values so they can be delegated separately.

When a trap occurs in HS-mode or U-mode, it goes to M-mode, unless delegated by `medeleg` or `mideleg`, in which case it goes to HS-mode. If the N extension for user-mode interrupts is implemented, then U-mode traps destined for HS-mode may be further delegated to U-mode using the `sedeleg` and `sideleg` CSRs.

When a trap occurs in VS-mode or VU-mode, it goes to M-mode, unless delegated by `medeleg` or `mideleg`, in which case it goes to HS-mode, unless further delegated by `hedeleg` or `hideleg`, in which case it goes to VS-mode. If the N extension for user-mode interrupts is implemented, then VU-mode traps destined for VS-mode may be further delegated to VU-mode using the `sedeleg` and `sideleg` CSRs.

When a trap is taken into M-mode, the following occurs: first, if the virtualization mode `V` was 1, the contents of the background supervisor registers are swapped with their foreground counterparts. Then, `mstatus.MPV` and `mstatus.MPP` are set according to Table 5.6.

When a trap is taken into HS-mode, the following occurs: first, if the virtualization mode `V` was 1, the contents of the background supervisor registers are swapped with their foreground counterparts. Then, `hstatus.SP2V` is set to `hstatus.SPV`, `hstatus.SP2P` is set to `sstatus.SPP`, and lastly `hstatus.SPV` and `sstatus.SPP` are set according to Table 5.7.

When a trap is taken into VS-mode, `sstatus.SPP` is set according to Table 5.8. Bits `SP2V`, `SP2P`, and `SPV` of `hstatus` are not modified, and the current virtualization state `V` remains 1.

Interrupt	Exception Code	Description
1	0	User software interrupt
1	1	Supervisor software interrupt
1	2	<i>Reserved</i>
1	3	Machine software interrupt
1	4	User timer interrupt
1	5	Supervisor timer interrupt
1	6	<i>Reserved</i>
1	7	Machine timer interrupt
1	8	User external interrupt
1	9	Supervisor external interrupt
1	10	<i>Reserved</i>
1	11	Machine external interrupt
1	≥ 12	<i>Reserved</i>
0	0	Instruction address misaligned
0	1	Instruction access fault
0	2	Illegal instruction
0	3	Breakpoint
0	4	Load address misaligned
0	5	Load access fault
0	6	Store/AMO address misaligned
0	7	Store/AMO access fault
0	8	Environment call from U-mode or VU-mode
0	9	Environment call from HS-mode
0	10	Environment call from VS-mode
0	11	Environment call from M-mode
0	12	Instruction page fault
0	13	Load page fault
0	14	<i>Reserved</i>
0	15	Store/AMO page fault
0	≥ 16	<i>Reserved</i>

Table 5.5: Supervisor and machine cause register (**scause** and **mcause**) values when the hypervisor extension is enabled.

Previous Mode	MPV	MPP
U-mode	0	0
HS-mode	0	1
M-mode	0	3
VU-mode	1	0
VS-mode	1	1

Table 5.6: Value of **mstatus** fields MPV and MPP after a trap into M-mode. Upon trap return, MPV is ignored when MPP=3.

Previous Mode	SPV	SPP
U-mode	0	0
HS-mode	0	1
VU-mode	1	0
VS-mode	1	1

Table 5.7: Value of `hstatus` field SPV and `sstatus` field SPP after a trap into HS-mode.

Previous Mode	SPP
VU-mode	0
VS-mode	1

Table 5.8: Value of `sstatus` field SPP after a trap into VS-mode.

5.8 Trap Return

The MRET instruction is used to return from a trap taken into M-mode. MRET sets the privilege mode according to the values of MPP and MPV in `mstatus`, as encoded in Table 5.6. MRET then in `mstatus` sets MPV=0, MPP=0, MIE=MPIE, and MPIE=1, and also sets `pc=mepc`. Finally, if the new virtualization mode V=1, the contents of the background supervisor registers are swapped with their foreground counterparts.

The SRET instruction is used to return from a trap taken into HS-mode or VS-mode. Its behavior depends on the current virtualization mode.

When executed in M-mode or HS-mode (i.e., V=0), SRET sets the virtualization and privilege modes according to the values in `hstatus.SP2V` and `sstatus.SPP`, as encoded in Table 5.7. SRET then sets `hstatus.SP2V=hstatus.SP2V`, `sstatus.SPP=hstatus.SP2P`, `hstatus.SP2V=0`, `hstatus.SP2P=0`, `sstatus.SIE=sstatus.SPIE`, `sstatus.SPIE=1`, and `pc=sepc`. Finally, if the new virtualization mode V=1, the contents of the background supervisor registers are swapped with their foreground counterparts.

When executed in VS-mode (i.e., V=1), SRET sets the privilege mode according to Table 5.8, then in `sstatus` sets SPP=0, SIE=SPIE, and SPIE=1, and lastly sets `pc=sepc`.

Chapter 6

RISC-V Privileged Instruction Set Listings

This chapter presents instruction-set listings for all instructions defined in the RISC-V Privileged Architecture.

31	27	26	25	24	20	19	15	14	12	11	7	6	0	
funct7				rs2		rs1		funct3		rd		opcode		R-type
imm[11:0]						rs1		funct3		rd		opcode		I-type
Environment Call and Breakpoint														
000000000000						00000		000		00000		1110011		ECALL
000000000001						00000		000		00000		1110011		EBREAK
Trap-Return Instructions														
0000000				00010		00000		000		00000		1110011		URET
0001000				00010		00000		000		00000		1110011		SRET
0011000				00010		00000		000		00000		1110011		MRET
Interrupt-Management Instructions														
0001000				00101		00000		000		00000		1110011		WFI
Memory-Management Instructions														
0001001				rs2		rs1		000		00000		1110011		SFENCE.VMA

Table 6.1: RISC-V Privileged Instructions

Chapter 7

Platform-Level Interrupt Controller (PLIC)

This chapter describes the general architecture for the RISC-V platform-level interrupt controller (PLIC), which prioritizes and distributes global interrupts in a RISC-V system.

7.1 PLIC Overview

Figure 7.1 provides a quick overview of PLIC operation. The PLIC connects global *interrupt sources*, which are usually I/O devices, to *interrupt targets*, which are usually *hart contexts*. The PLIC contains multiple *interrupt gateways*, one per interrupt source, together with a *PLIC core* that performs interrupt prioritization and routing. Global interrupts are sent from their source to an *interrupt gateway* that processes the interrupt signal from each source and sends a single *interrupt request* to the PLIC core, which latches these in the core interrupt pending bits (IP). Each interrupt source is assigned a separate priority. The PLIC core contains a matrix of interrupt enable (IE) bits to select the interrupts that are enabled for each target. The PLIC core forwards an *interrupt notification* to one or more targets if the targets have any pending interrupts enabled, and the priority of the pending interrupts exceeds a per-target threshold. When the target takes the external interrupt, it sends an *interrupt claim* request to retrieve the identifier of the highest-priority global interrupt source pending for that target from the PLIC core, which then clears the corresponding interrupt source pending bit. After the target has serviced the interrupt, it sends the associated interrupt gateway an *interrupt completion* message and the interrupt gateway can now forward another interrupt request for the same source to the PLIC. The rest of this chapter describes each of these components in detail, though many details are necessarily platform specific.

7.2 Interrupt Sources

RISC-V harts can have both local and global interrupt sources. Only global interrupt sources are handled by the PLIC.

to lower privilege levels, using the corresponding bits in the `mideleg` register. The priority of non-standard local interrupt sources relative to external, timer, and software interrupts is platform-specific.

7.2.2 Global Interrupt Sources

Global interrupt sources are those that are prioritized and distributed by the PLIC. Depending on the platform-specific PLIC implementation, any global interrupt source could be routed to any hart context.

Global interrupt sources can take many forms, including level-triggered, edge-triggered, and message-signalled. Some sources might queue up a number of interrupt requests. All global interrupt sources are converted to a common interrupt request format for the PLIC.

7.3 Interrupt Targets and Hart Contexts

Interrupt targets are usually hart contexts, where a hart context is a given privilege mode on a given hart (though there are other possible interrupt targets, such as DMA engines). Not all hart contexts need be interrupt targets, in particular, if a processor core does not support delegating external interrupts to lower-privilege modes, then the lower-privilege hart contexts will not be interrupt targets. Interrupt notifications generated by the PLIC appear in the `meip/seip/ueip` bits of the `mip/sip/uiip` registers for M/S/U modes, respectively. For the notifications to appear in lower-privilege `xip` registers, the corresponding external interrupts must have been delegated in the higher-privilege `yideleg` registers.

Each processor core must define a policy on how simultaneous active interrupts are taken by multiple hart contexts on the core. For the simple case of a single stack of hart contexts, one for each supported privileged mode, interrupts for higher-privilege contexts can preempt execution of interrupt handlers for lower-privilege contexts. A multithreaded processor core could run multiple independent interrupt handlers on different hart contexts at the same time. A processor core could also provide hart contexts that are only used for interrupt handling to reduce interrupt service latency, and these might preempt interrupt handlers for other harts on the same core.

The PLIC treats each interrupt target independently and does not take into account any interrupt prioritization scheme used by a component that contains multiple interrupt targets. As a result, the PLIC provides no concept of interrupt preemption or nesting so this must be handled by the cores hosting multiple interrupt target contexts.

7.4 Interrupt Gateways

The interrupt gateways are responsible for converting global interrupt signals into a common interrupt request format, and for controlling the flow of interrupt requests to the PLIC core. At most one interrupt request per interrupt source can be pending in the PLIC core at any time, indicated by setting the source's IP bit. The gateway only forwards a new interrupt request to the PLIC

core after receiving notification that the interrupt handler servicing the previous interrupt request from the same source has completed.

If the global interrupt source uses level-sensitive interrupts, the gateway will convert the first assertion of the interrupt level into an interrupt request, but thereafter the gateway will not forward an additional interrupt request until it receives an interrupt completion message. On receiving an interrupt completion message, if the interrupt is level-triggered and the interrupt is still asserted, a new interrupt request will be forwarded to the PLIC core. The gateway does not have the facility to retract an interrupt request once forwarded to the PLIC core. If a level-sensitive interrupt source deasserts the interrupt after the PLIC core accepts the request and before the interrupt is serviced, the interrupt request remains present in the IP bit of the PLIC core and will be serviced by a handler, which will then have to determine that the interrupt device no longer requires service.

If the global interrupt source was edge-triggered, the gateway will convert the first matching signal edge into an interrupt request. Depending on the design of the device and the interrupt handler, between sending an interrupt request and receiving notice of its handler's completion, the gateway might either ignore additional matching edges or increment a counter of pending interrupts. In either case, the next interrupt request will not be forwarded to the PLIC core until the previous completion message has been received. If the gateway has a pending interrupt counter, the counter will be decremented when the interrupt request is accepted by the PLIC core.

Unlike dedicated-wire interrupt signals, message-signalled interrupts (MSIs) are sent over the system interconnect via a message packet that describes which interrupt is being asserted. The message is decoded to select an interrupt gateway, and the relevant gateway then handles the MSI similar to an edge-triggered interrupt.

7.5 Interrupt Identifiers (IDs)

Global interrupt sources are assigned small unsigned integer identifiers, beginning at the value 1. An interrupt ID of 0 is reserved to mean “no interrupt”.

Interrupt identifiers are also used to break ties when two or more interrupt sources have the same assigned priority. Smaller values of interrupt ID take precedence over larger values of interrupt ID.

7.6 Interrupt Priorities

Interrupt priorities are small unsigned integers, with a platform-specific maximum number of supported levels. The priority value 0 is reserved to mean “never interrupt”, and interrupt priority increases with increasing integer values.

Each global interrupt source has an associated interrupt priority held in a platform-specific memory-mapped register. Different interrupt sources need not support the same set of priority values. A valid implementation can hardwire all input priority levels. Interrupt source priority registers should be **WARL** fields to allow software to determine the number and position of read-write bits in each priority specification, if any. To simplify discovery of supported priority values, each

priority register must support any combination of values in the bits that are variable within the register, i.e., if there are two variable bits in the register, all four combinations of values in those bits must operate as valid priority levels.

In the degenerate case, all priorities can be hardwired to the value 1, in which case input priorities are effectively determined by interrupt ID.

The supported priority values can be determined as follows: 1) write all zeros to the priority register then 2) read back the value. Any set bits are hardwired to 1. Next, 3) write all ones to the register, and 4) read back the value. Any clear bits are hardwired to 0. Any set bits that were not found to be hardwired in step 2 are variable. The supported priority levels are the set of values obtained by substituting all combinations of ones and zeros in the variable bits within the priority field.

7.7 Interrupt Enables

Each target has a vector of interrupt enable (IE) bits, one per interrupt source. The target will not receive interrupts from sources that are disabled. The IE bits for a single target should be packed together as a bit vector in platform-specific memory-mapped control registers to support rapid context switching of the IE bits for a target. IE bits are **WARL** fields that can be hardwired to either 0 or 1.

A large number of potential IE bits might be hardwired to zero in cases where some interrupt sources can only be routed to a subset of targets.

A larger number of bits might be wired to 1 for an embedded device with fixed interrupt routing. Interrupt priorities, thresholds, and hart-internal interrupt masking provide considerable flexibility in ignoring external interrupts even if a global interrupt source is always enabled.

7.8 Interrupt Priority Thresholds

Each interrupt target has an associated priority threshold, held in a platform-specific memory-mapped register. Only active interrupts that have a priority strictly greater than the threshold will cause an interrupt notification to be sent to the target. Different interrupt targets need not support the same set of priority threshold values. Interrupt target threshold registers should be **WARL** fields to allow software to determine the supported thresholds. A threshold register should always be able to hold the value zero, in which case, no interrupts are masked. If implemented, the threshold register will usually also be able to hold the maximum priority level, in which case all interrupts are masked.

A simple valid implementation is to hardwire the threshold to zero, in which case it has no effect, and the individual enable bits will have to be saved and restored to attain the same effect. While the function of the threshold can be achieved by changing the interrupt-enable bits, manipulating a single threshold value avoids the target having to consider the individual priority levels of each interrupt source, and saving and restoring all the interrupt enables. Changing the threshold quickly might be especially important for systems that move frequently between power states.

7.9 Interrupt Notifications

Each interrupt target has an *external interrupt pending* (EIP) bit in the PLIC core that indicates that the corresponding target has a pending interrupt waiting for service. The value in EIP can change as a result of changes to state in the PLIC core, brought on by interrupt sources, interrupt targets, or other agents manipulating register values in the PLIC. The value in EIP is communicated to the destination target as an interrupt notification. If the target is a RISC-V hart context, the interrupt notifications arrive on the `meip/seip/ueip` bits depending on the privilege level of the hart context.

In simple systems, the interrupt notifications will be simple wires connected to the processor implementing a hart. In more complex platforms, the notifications might be routed as messages across a system interconnect.

The PLIC hardware only supports multicasting of interrupts, such that all enabled targets will receive interrupt notifications for a given active interrupt.

Multicasting provides rapid response since the fastest responder claims the interrupt, but can be wasteful in high-interrupt-rate scenarios if multiple harts take a trap for an interrupt that only one can successfully claim. Software can modulate the PLIC IE bits as part of each interrupt handler to provide alternate policies, such as interrupt affinity or round-robin unicasting.

Depending on the platform architecture and the method used to transport interrupt notifications, these might take some time to be received at the targets. The PLIC is guaranteed to eventually deliver all state changes in EIP to all targets, provided there is no intervening activity in the PLIC core.

The value in an interrupt notification is only guaranteed to hold an EIP value that was valid at some point in the past. In particular, a second target can respond and claim an interrupt while a notification to the first target is still in flight, such that when the first target tries to claim the interrupt it finds it has no active interrupts in the PLIC core.

7.10 Interrupt Claims

Sometime after a target receives an interrupt notification, it might decide to service the interrupt. The target sends an *interrupt claim* message to the PLIC core, which will usually be implemented as a non-idempotent memory-mapped I/O control register read. On receiving a claim message, the PLIC core will atomically determine the ID of the highest-priority pending interrupt for the target and then clear down the corresponding source's IP bit. The PLIC core will then return the ID to the target. The PLIC core will return an ID of zero, if there were no pending interrupts for the target when the claim was serviced.

After the highest-priority pending interrupt is claimed by a target and the corresponding IP bit is cleared, other lower-priority pending interrupts might then become visible to the target, and so the PLIC EIP bit might not be cleared after a claim. The interrupt handler can check the local `meip/seip/ueip` bits before exiting the handler, to allow more efficient service of other interrupts without first restoring the interrupted context and taking another interrupt trap.

It is always legal for a hart to perform a claim even if the EIP is not set. In particular, a hart could set the threshold value to maximum to disable interrupt notifications and instead poll for active interrupts using periodic claim requests, though a simpler approach to implement polling would be to clear the external interrupt enable in the corresponding *xie* register for privilege mode *x*.

7.11 Interrupt Completion

After a handler has completed service of an interrupt, the associated gateway must be sent an interrupt completion message, usually as a write to a non-idempotent memory-mapped I/O control register. The gateway will only forward additional interrupts to the PLIC core after receiving the completion message.

7.12 Interrupt Flow

Figure 7.2 shows the messages flowing between agents when handling interrupts via the PLIC.

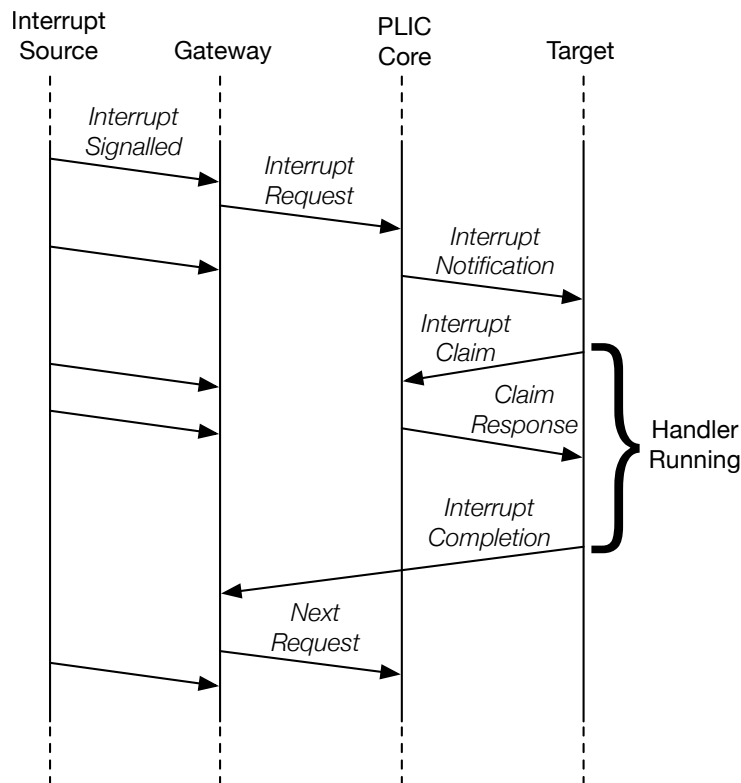


Figure 7.2: Flow of interrupt processing via the PLIC.

The gateway will only forward a single interrupt request at a time to the PLIC, and not forward subsequent interrupts requests until an interrupt completion is received. The PLIC will set the IP bit once it accepts an interrupt request from the gateway, and sometime later forward an interrupt notification to the target. The target might take a while to respond to a new interrupt arriving,

but will then send an interrupt claim request to the PLIC core to obtain the interrupt ID. The PLIC core will atomically return the ID and clear the corresponding IP bit, after which no other target can claim the same interrupt request. Once the handler has processed the interrupt, it sends an interrupt completion message to the gateway to allow a new interrupt request.

7.13 PLIC Core Specification

The operation of the PLIC core can be specified as a non-deterministic finite-state machine with input and output message queues, with the following atomic actions:

- **Write Register:** A message containing a register write request is dequeued. One of the internal registers is written, where an internal register can be a priority, an interrupt-enable (IE), or a threshold.
- **Accept Request:** If the IP bit corresponding to the interrupt source is clear, a message containing an interrupt request from a gateway is dequeued and the IP bit is set.
- **Process Claim:** An interrupt claim message is dequeued. A claim-response message is enqueued to the requester with the ID of the highest-priority active interrupt for that target, and the IP bit corresponding to this interrupt source is cleared.

The value in the EIP bit is determined as a combinational function of the PLIC Core state. Interrupt notifications are sent via an autonomous process that ensures the EIP value is eventually reflected at the target.

Note that the operation of the interrupt gateways is decoupled from the PLIC core. A gateway can handle parsing of interrupt signals and processing interrupt completion messages concurrently with other operations in the PLIC core.

Figure 7.1 is a high-level conceptual view of the PLIC design. The PLIC core can be implemented in many ways provided its behavior can always be understood as following from some sequential ordering of these atomic actions. In particular, the PLIC might process multiple actions in a single clock cycle, or might process each action over many clock cycles.

7.14 Controlling Access to the PLIC

In the expected use case, only machine mode accesses the source priority, source pending, and target interrupt enables to configure the interrupt subsystem. Lower-privilege modes access these features via ABI or SBI calls. The interrupt enables act as a protection mechanism where a target can only signal completion to an interrupt gateway that is currently enabled for that target.

Interrupt handlers that run with lower than machine-mode privilege need only be able to perform a claim read and a completion write, and to set their target threshold value. The memory map for these registers should allow machine mode to protect different targets from each other's accesses, using either physical memory protection or virtual memory page protections.

Chapter 8

Machine Configuration Description

To reduce porting effort for OS boots, we have reverted back to using Device Trees to communicate platform information to the kernel, so this chapter is out of date. Config string was designed for other uses in addition, but for now, we are staying with a standard device tree model.

RISC-V platforms may contain myriad devices, processor cores, and configuration parameters. To support higher-level software, including bootloaders and operating systems, it is recommended that hardware platforms embed a description of their components in read-only memory that is directly accessible after processor reset for use by low-level system software, external debuggers, or manufacturing test procedures. We call this low-level embedded information a configuration description. We define here a standard mechanism to encode and locate the configuration information, and to determine the format of the configuration information.

8.1 Configuration String Search Procedure

The platform must describe how to locate a pointer to find this string, for example, by specifying a fixed physical address at which the pointer resides. To support a wide variety of platforms, configuration formats, and chips with manufacturing-time programming of configuration options, a flexible search procedure is defined to locate the configuration information seeded by the initial pointer specified by the platform.

The configuration string pointer provided by the platform points to an initial memory address at which the search for configuration string begins.

The configuration string cannot begin with a padding byte, where a padding byte is defined to contain either 0x0 or 0xff, but can be preceded by up to 63 padding bytes that are ignored. If 64 padding bytes are encountered, then the search terminates without finding a config string.

The padding bytes represent common values returned by unpopulated memory or bus regions or unprogrammed non-volatile memory. Configuration strings can therefore include pointers to regions that are optionally populated or programmed, and these regions will be ignored if there is nothing present. The padding bytes also support alignment of binary data structures.

Otherwise the first non-padding byte is the beginning of the configuration information. For example, configuration information in Device Tree String format would begin with a “/dts-v1/”.

Configuration information in Flattened Device Tree format would begin with the magic number `0xd00dfeed`. Configuration information in the config string format would begin with `“/cs-v1/”`.

Config string is a new format that is backwards-compatible with device tree string (as far as DTS specs exist) but can include additional configuration information in other memory regions.

Chapter 9

History

9.1 Research Funding at UC Berkeley

Development of the RISC-V architecture and implementations has been partially funded by the following sponsors.

- **Par Lab:** Research supported by Microsoft (Award #024263) and Intel (Award #024894) funding and by matching funding by U.C. Discovery (Award #DIG07-10227). Additional support came from Par Lab affiliates Nokia, NVIDIA, Oracle, and Samsung.
- **Project Isis:** DoE Award DE-SC0003624.
- **ASPIRE Lab:** DARPA PERFECT program, Award HR0011-12-2-0016. DARPA POEM program Award HR0011-11-C-0100. The Center for Future Architectures Research (C-FAR), a STARnet center funded by the Semiconductor Research Corporation. Additional support from ASPIRE industrial sponsor, Intel, and ASPIRE affiliates, Google, Huawei, Nokia, NVIDIA, Oracle, and Samsung.

The content of this paper does not necessarily reflect the position or the policy of the US government and no official endorsement should be inferred.

Bibliography

- [1] Robert P. Goldberg. Survey of virtual machine research. *Computer*, 7(6):34–45, June 1974.
- [2] Juan Navarro, Sitaram Iyer, Peter Druschel, and Alan Cox. Practical, transparent operating system support for superpages. *SIGOPS Oper. Syst. Rev.*, 36(SI):89–104, December 2002.