

Data61 Trustworthy Systems
<https://ts.data61.csiro.au/projects/TS/>

seL4 Reference Manual

Version 12.0.0-dev

Trustworthy Systems Team, Data61
<https://sel4.systems/contact/>

25 May 2021

© 2021 General Dynamics C4 Systems.

ALL RIGHTS RESERVED.

Acknowledgements

The primary authors of this document are Matthew Grosvenor and Adam Walker, with contributions from Adrian Danis, Andrew Boyton, Anna Lyons, David Greenaway, Etienne Le Sueur, Gernot Heiser, Gerwin Klein, Godfrey van der Linden, Kevin Elphinstone, Matthew Fernandez, Matthias Daum, Michael von Tessin, Peter Chubb, Simon Winwood, Thomas Sewell, Timothy Bourke and Toby Murray. All authors and contributors can be contacted at firstname.lastname@data61.csiro.au.

Contents

1	Introduction	1
2	Kernel Services and Objects	2
2.1	Capability-based Access Control	3
2.2	System Calls	3
2.3	Kernel Objects	5
2.4	Kernel Memory Allocation	7
2.4.1	Reusing Memory	8
2.4.2	Summary of Object Sizes	8
3	Capability Spaces	10
3.1	Capability and CSpace Management	11
3.1.1	CSpace Creation	11
3.1.2	CNode Methods	11
3.1.3	Capabilities to Newly-Retyped Objects	12
3.1.4	Capability Rights	12
3.1.5	Capability Derivation Tree	13
3.2	Deletion and Revocation	14
3.3	CSpace Addressing	15
3.3.1	Capability Address Lookup	15
3.3.2	Addressing Capabilities	16
3.4	Lookup Failure Description	18
3.4.1	Invalid Root	18
3.4.2	Missing Capability	18
3.4.3	Depth Mismatch	19
3.4.4	Guard Mismatch	19

4	Message Passing (IPC)	20
4.1	Message Registers	20
4.2	Endpoints	21
4.2.1	Endpoint Badges	22
4.2.2	Capability Transfer	22
4.2.3	Errors	23
4.2.4	Calling and Replying	23
5	Notifications	25
5.1	Notification Objects	25
5.2	Signalling, Polling and Waiting	25
5.3	Binding Notifications	26
6	Threads and Execution	27
6.1	Threads	27
6.1.1	Thread control blocks	27
6.1.2	Thread Creation	27
6.1.3	Thread Deactivation	28
6.1.4	Scheduling	28
6.1.5	MCS Scheduling	28
6.1.6	Scheduling Contexts	28
6.1.7	Passive Threads	30
6.1.8	Scheduling Context Creation	30
6.1.9	Scheduling Context Donation	31
6.1.10	Scheduling algorithm	32
6.1.11	Exceptions	32
6.1.11.1	Standard Exceptions	32
6.1.11.2	Timeout Exceptions (MCS Only)	33
6.1.12	Message Layout of the Read-/Write-Registers Methods	33
6.2	Faults	33
6.2.1	Capability Faults	34
6.2.2	Unknown Syscall	34
6.2.3	User Exception	34
6.2.4	Debug Exception: Breakpoints and Watchpoints	35

6.2.5	Debug Exception: Single-stepping	36
6.2.6	Timeout Fault (MCS only)	37
6.2.7	VM Fault	37
6.3	Domains	38
6.4	Virtualisation	38
6.4.1	ARM	39
6.4.2	x86	39
7	Address Spaces and Virtual Memory	40
7.1	Objects	40
7.1.1	Hardware Virtual Memory Objects	40
7.1.1.1	IA-32	41
7.1.1.2	x64	41
7.1.1.3	AArch32	41
7.1.1.4	AArch64	42
7.1.2	RISC-V	42
7.1.2.1	RISC-V 32-bit	42
7.1.2.2	RISC-V 64-bit	42
7.1.3	Page	42
7.1.3.1	AArch32 page sizes	43
7.1.3.2	AArch64 page sizes	43
7.1.3.3	IA-32 page sizes	43
7.1.3.4	X64 page sizes	44
7.1.3.5	RISC-V 32-bit page sizes	44
7.1.3.6	RISC-V 64-bit page sizes	44
7.1.4	ASID Control	44
7.1.5	ASID Pool	44
7.2	Mapping Attributes	45
7.3	Sharing Memory	45
7.4	Page Faults	45
8	Hardware I/O	46
8.1	Interrupt Delivery	46
8.2	x86-Specific I/O	46

8.2.1	Interrupts	46
8.2.2	I/O Ports	47
8.2.3	I/O Space	47
8.3	Arm-Specific I/O	48
8.3.1	Arm SMMU version 2.0	48
8.3.1.1	Creating seL4_ARM_SID capabilities	50
8.3.1.2	Creating seL4_ARM_CB capabilities	50
8.3.1.3	Configuring context banks	51
8.3.1.4	Configuring streams (transactions)	51
8.3.1.5	Copying and Deleting caps	51
8.3.1.6	TLB invalidation	52
8.3.1.7	Fault handling	53
9	System Bootstrapping	55
9.1	Initial Thread's Environment	55
9.2	BootInfo Frame	55
9.3	Boot Command-line Arguments	58
10	seL4 API Reference	60
10.1	Error Codes	60
10.1.1	Invalid Argument	60
10.1.2	Invalid Capability	60
10.1.3	Illegal Operation	60
10.1.4	Range Error	61
10.1.5	Alignment Error	61
10.1.6	Failed Lookup	61
10.1.7	Truncated Message	61
10.1.8	Delete First	61
10.1.9	Revoke First	62
10.1.10	Not Enough Memory	62
10.2	System Calls	62
10.2.1	General System Calls	62
10.2.1.1	Send	62
10.2.1.2	Recv	63

10.2.1.3	Call	63
10.2.1.4	Reply	63
10.2.1.5	Non-Blocking Send	64
10.2.1.6	Reply Recv	64
10.2.1.7	NBRecv	65
10.2.1.8	Yield	65
10.2.1.9	Signal	66
10.2.1.10	Wait	66
10.2.1.11	Poll	67
10.2.2	General System Calls (MCS)	68
10.2.2.1	Send	68
10.2.2.2	Recv	68
10.2.2.3	Call	69
10.2.2.4	Non-Blocking Send	69
10.2.2.5	Reply Recv	70
10.2.2.6	NBRecv	70
10.2.2.7	NBSend Recv	71
10.2.2.8	NBSend Wait	71
10.2.2.9	Yield	72
10.2.2.10	Wait	72
10.2.2.11	NBWait	73
10.2.2.12	Poll	73
10.2.2.13	Signal	74
10.2.3	Debugging System Calls	75
10.2.3.1	Put Char	75
10.2.3.2	Dump scheduler	75
10.2.3.3	Halt	76
10.2.3.4	Snapshot	76
10.2.3.5	Cap Identify	76
10.2.3.6	Name Thread	77
10.2.3.7	Send SGI 0-15	77
10.2.3.8	Run	78
10.2.4	Benchmarking System Calls	79

10.2.4.1	Reset Log	79
10.2.4.2	Finalize Log	80
10.2.4.3	Set Log Buffer	80
10.2.4.4	Null Syscall	81
10.2.4.5	Flush Caches	81
10.2.4.6	Flush L1 Caches	81
10.2.4.7	Get Thread Utilisation	82
10.2.4.8	Reset Thread Utilisation	82
10.2.4.9	Dump All Threads Utilisation	82
10.2.4.10	Reset All Threads Utilisation	83
10.2.5	X86 System Calls	84
10.2.5.1	VMEnter	84
10.3	Architecture-Independent Object Methods	86
10.3.1	seL4_CNode	86
10.3.1.1	Cancel Badged Sends	86
10.3.1.2	Copy	87
10.3.1.3	Delete	88
10.3.1.4	Mint	89
10.3.1.5	Move	90
10.3.1.6	Mutate	91
10.3.1.7	Revoke	92
10.3.1.8	Rotate	93
10.3.1.9	Save Caller	94
10.3.2	seL4_DomainSet	94
10.3.2.1	Set	94
10.3.3	seL4_IRQControl	95
10.3.3.1	Get	95
10.3.4	seL4_IRQHandler	95
10.3.4.1	Acknowledge	95
10.3.4.2	Clear	96
10.3.4.3	Set Notification	96
10.3.5	seL4_SchedContext	97
10.3.5.1	Bind	97

10.3.5.2	Consumed	97
10.3.5.3	Unbind	98
10.3.5.4	UnbindObject	98
10.3.5.5	YieldTo	99
10.3.6	seL4_SchedControl	100
10.3.6.1	ConfigureFlags	100
10.3.7	seL4_TCB	101
10.3.7.1	Bind Notification	101
10.3.7.2	Configure (MCS)	101
10.3.7.3	Configure Single Stepping	102
10.3.7.4	Configure	103
10.3.7.5	Copy Registers	104
10.3.7.6	Get Breakpoint	105
10.3.7.7	Read Registers	106
10.3.7.8	Resume	106
10.3.7.9	Set Breakpoint	107
10.3.7.10	Set CPU Affinity	108
10.3.7.11	Set IPC Buffer	108
10.3.7.12	Set Maximum Controlled Priority	109
10.3.7.13	Set Priority	109
10.3.7.14	Set Sched Params (MCS)	110
10.3.7.15	Set Sched Params	111
10.3.7.16	Set Space	111
10.3.7.17	Set TLS Base	112
10.3.7.18	Set Timeout Endpoint	112
10.3.7.19	Suspend	113
10.3.7.20	Unbind Notification	113
10.3.7.21	Unset Breakpoint	114
10.3.7.22	Write Registers	114
10.3.8	seL4_Untyped	115
10.3.8.1	Retype	115
10.4	x86-Specific Object Methods	116
10.4.1	seL4_IRQControl	116

10.4.1.1	Get I/O APIC	116
10.4.1.2	Get MSI	117
10.4.2	seL4_TCB	118
10.4.2.1	Set EPT Root	118
10.4.3	seL4_X86_ASIDControl	119
10.4.3.1	Make Pool	119
10.4.4	seL4_X86_ASIDPool	120
10.4.4.1	Assign	120
10.4.5	seL4_X86_ETPD	121
10.4.5.1	Map	121
10.4.5.2	Unmap	121
10.4.6	seL4_X86_ETPDPT	122
10.4.6.1	Map	122
10.4.6.2	Unmap	122
10.4.7	seL4_X86_EPTPT	123
10.4.7.1	Map	123
10.4.7.2	Unmap	123
10.4.8	seL4_X86_IOPageTable	124
10.4.8.1	Map	124
10.4.8.2	Unmap	124
10.4.9	seL4_X86_IOPort	125
10.4.9.1	In16	125
10.4.9.2	In32	125
10.4.9.3	In8	125
10.4.9.4	Out16	126
10.4.9.5	Out32	126
10.4.9.6	Out8	127
10.4.10	seL4_X86_IOPortControl	127
10.4.10.1	Issue	127
10.4.11	seL4_X86_Page	128
10.4.11.1	Get Address	128
10.4.11.2	Map EPT	128
10.4.11.3	Map I/O	129

10.4.11.4	Map	130
10.4.11.5	Unmap	130
10.4.12	seL4_X86_PageDirectory	131
10.4.12.1	Get Status Bits	131
10.4.12.2	Map	131
10.4.12.3	Unmap	132
10.4.13	seL4_X86_PageTable	132
10.4.13.1	Map	132
10.4.13.2	Unmap	133
10.4.14	seL4_X86_VCPU	133
10.4.14.1	Disable IO Port	133
10.4.14.2	Enable IO Port	134
10.4.14.3	Read VMCS	134
10.4.14.4	Set TCB	135
10.4.14.5	Write Registers	135
10.4.14.6	Write VMCS	136
10.5	IA32-Specific Object Methods	137
10.6	x86_64-Specific Object Methods	138
10.6.1	seL4_X86_PDPT	138
10.6.1.1	Map	138
10.6.1.2	Unmap	138
10.7	ARM-Specific Object Methods	139
10.7.1	seL4_ARM_ASIDControl	139
10.7.1.1	Make Pool	139
10.7.2	seL4_ARM_ASIDPool	140
10.7.2.1	Asid Pool Assign	140
10.7.3	seL4_ARM_CB	140
10.7.3.1	AssignVspace	140
10.7.3.2	CBClearFault	141
10.7.3.3	CBGetFault	141
10.7.3.4	TLBInvalidate	142
10.7.3.5	UnassignVspace	142
10.7.4	seL4_ARM_CBControl	143

10.7.4.1	GetCB	143
10.7.4.2	TLBInvalidateAll	143
10.7.5	seL4_ARM_IOPageTable	144
10.7.5.1	Map	144
10.7.5.2	Unmap	144
10.7.6	seL4_ARM_Page	145
10.7.6.1	Clean Data	145
10.7.6.2	Clean and Invalidate Data	145
10.7.6.3	Get Address	146
10.7.6.4	Invalidate Data	146
10.7.6.5	Map I/O	147
10.7.6.6	Map	147
10.7.6.7	Unify Instruction	148
10.7.6.8	Unmap	148
10.7.7	seL4_ARM_PageTable	149
10.7.7.1	Map	149
10.7.7.2	Unmap	150
10.7.8	seL4_ARM_SID	150
10.7.8.1	BindCB	150
10.7.8.2	UnbindCB	151
10.7.9	seL4_ARM_SIDControl	151
10.7.9.1	ClearFault	151
10.7.9.2	GetFault	152
10.7.9.3	GetSID	152
10.7.10	seL4_ARM_VCPU	153
10.7.10.1	Acknowledge Virtual PPI IRQ	153
10.7.10.2	Inject IRQ	153
10.7.10.3	Read Registers	154
10.7.10.4	Set TCB	154
10.7.10.5	Write Registers	155
10.7.11	seL4_IRQControl	155
10.7.11.1	GetTrigger	155
10.7.11.2	GetTriggerCore	156

10.8	Aarch32-Specific Object Methods	157
10.8.1	seL4_ARM_PageDirectory	157
10.8.1.1	Clean Data	157
10.8.1.2	Clean and Invalidate Data	157
10.8.1.3	Invalidate Data	158
10.8.1.4	Unify Instruction	158
10.9	Aarch64-Specific Object Methods	159
10.9.1	seL4_ARM_PageDirectory	159
10.9.1.1	Map	159
10.9.1.2	Unmap	159
10.9.2	seL4_ARM_PageUpperDirectory	160
10.9.2.1	Map	160
10.9.2.2	Unmap	160
10.9.3	seL4_ARM_VSpace	161
10.9.3.1	Clean Data	161
10.9.3.2	Clean and Invalidate Data	161
10.9.3.3	Invalidate Data	162
10.9.3.4	Unify Instruction	162
10.10	RISCV-Specific Object Methods	163
10.10.1	General RISCV Object Methods	163
10.10.2	seL4_IRQControl	163
10.10.3	seL4_RISCV_ASIDControl	163
10.10.3.1	MakePool	163
10.10.4	seL4_RISCV_ASIDPool	164
10.10.4.1	Assign	164
10.10.5	seL4_RISCV_Page	164
10.10.5.1	GetAddress	164
10.10.5.2	Map	165
10.10.5.3	Unmap	165
10.10.6	seL4_RISCV_PageTable	166
10.10.6.1	Map	166
10.10.6.2	Unmap	166

List of Tables

2.1	Meaning of <code>size_bits</code> for object types of variable size	9
3.1	seL4 access rights: What a specific right entitles a capability to do . . .	12
3.2	Capability derivation.	13
4.1	Fields of the <code>seL4_IPCBuffer</code> structure. Note that <code>badges</code> and <code>caps</code> use the same area of memory in the structure.	21
6.1	Contents of an IPC message.	34
6.2	Debug fault message layout. The register API-ID is not returned in the fault message from the kernel on single-step faults.	36
6.3	Single-step fault message layout.	37
6.4	Timeout fault outcome on 32-bit architectures.	37
6.5	VM Fault outcome on all architectures.	38
7.1	Virtual memory attributes for ARM page table entries.	45
7.2	Virtual memory attributes for x86 page table entries.	45
9.1	Initial thread's CNode content.	56
9.2	BootInfo struct.	57
9.3	BootInfoHeader struct.	58
9.4	seL4_UntypedDesc struct	58
9.5	IA-32 boot command-line arguments.	59

List of Figures

3.1	Example capability derivation tree.	13
3.2	An example CSpace demonstrating object references at all levels, various guard and radix sizes and internal CNode references.	16
3.3	An arbitrary CSpace layout.	17

Chapter 1

Introduction

The seL4 microkernel is an operating-system kernel designed to be a secure, safe, and reliable foundation for systems in a wide variety of application domains. As a microkernel, it provides a small number of services to applications, such as abstractions to create and manage virtual address spaces, threads, and inter-process communication (IPC). The small number of services provided by seL4 directly translates to a small implementation of approximately 8700 lines of C code. This has allowed the ARMv6 version of the kernel to be formally proven in the Isabelle/HOL theorem prover to adhere to its formal specification [Boy09, CKS08, DEK⁺06, EKE08, KEH⁺09, TKN07, WKS⁺09], which in turn enabled proofs of the kernel's enforcement of integrity [SWG⁺11] and confidentiality [MMB⁺13]. The kernel's small size was also instrumental in performing a complete and sound analysis of worst-case execution time [BSC⁺11, BSH12].

This manual describes the seL4 kernel's API from a user's point of view. The document starts by giving a brief overview of the seL4 microkernel design, followed by a reference of the high-level API exposed by the seL4 kernel to userspace.

While we have tried to ensure that this manual accurately reflects the behaviour of the seL4 kernel, this document is by no means a formal specification of the kernel. When the precise behaviour of the kernel under a particular circumstance needs to be known, users should refer to the seL4 abstract specification, which gives a formal description of the seL4 kernel.

Chapter 2

Kernel Services and Objects

A limited number of service primitives are provided by the microkernel; more complex services may be implemented as applications on top of these primitives. In this way, the functionality of the system can be extended without increasing the code and complexity in privileged mode, while still supporting a potentially wide number of services for varied application domains.

Note that some services are available only when the kernel is configured for MCS¹ support.

The basic services seL4 provides are as follows:

Threads are an abstraction of CPU execution that supports running software;

Scheduling contexts (MCS only) are an abstraction of CPU execution time.

Address spaces are virtual memory spaces that each contain an application. Applications are limited to accessing memory in their address space;

Inter-process communication (IPC) via *endpoints* allows threads to communicate using message passing;

Reply objects (MCS only) are used to store single-use reply capabilities, and are provided by the receiver during message passing.

Notifications provide a non-blocking signalling mechanism similar to binary semaphores;

Device primitives allow device drivers to be implemented as unprivileged applications. The kernel exports hardware device interrupts via IPC messages; and

Capability spaces store capabilities (i.e., access rights) to kernel services along with their book-keeping information.

This chapter gives an overview of these services, describes how kernel objects are accessed by userspace applications, and describes how new objects can be created.

¹“mixed-criticality system”

2.1 Capability-based Access Control

The seL4 microkernel provides a capability-based access-control model. Access control governs all kernel services; in order to perform an operation, an application must *invoke* a capability in its possession that has sufficient access rights for the requested service. With this, the system can be configured to isolate software components from each other, and also to enable authorised, controlled communication between components by selectively granting specific communication capabilities. This enables software-component isolation with a high degree of assurance, as only those operations explicitly authorised by capability possession are permitted.

A capability is an unforgeable token that references a specific kernel object (such as a thread control block) and carries access rights that control what methods may be invoked. Conceptually, a capability resides in an application's *capability space*; an address in this space refers to a *slot* which may or may not contain a capability. An application may refer to a capability—to request a kernel service, for example—using the address of the slot holding that capability. This means, the seL4 capability model is an instance of a *segregated* (or *partitioned*) capability system, where capabilities are managed by the kernel.

Capability spaces are implemented as a directed graph of kernel-managed *capability nodes* (CNodes). A CNode is a table of slots, where each slot may contain further CNode capabilities. An address of a capability in a capability space is the concatenation of the indices of slots within CNodes forming the path to the destination slot; we discuss CNode objects in detail in Chapter 3.

Capabilities can be copied and moved within capability spaces, and also sent via IPC. This allows creation of applications with specific access rights, the delegation of authority to another application, and passing to an application authority to a newly created (or selected) kernel service. Furthermore, capabilities can be *minted* to create a derived capability with a subset of the rights of the original capability (never with more rights). A newly minted capability can be used for partial delegation of authority.

Capabilities can also be revoked to withdraw authority. Revocation recursively removes any capabilities that have been derived from the original capability being revoked. The propagation of capabilities through the system is controlled by a *take-grant*-based model [EKE08, Boy09].

2.2 System Calls

The seL4 kernel provides a message-passing service for communication between threads. This mechanism is also used for communication with kernel-provided services. There is a standard message format, each message containing a number of data words and possibly some capabilities. The structure and encoding of these messages are described in detail in Chapter 4.

Threads send messages by invoking capabilities within their capability space. When an endpoint, notification or reply capability is invoked in this way, the message will be transferred through the kernel to another thread. When other capabilities to kernel

objects are invoked, the message will be interpreted as a method invocation in a manner specific to the type of kernel object. For example, invoking a thread control block (TCB) capability with a correctly formatted message will suspend the target thread.

Fundamentally, we can regard the kernel as providing three system calls: *Send*, *Receive* and *Yield*. However, there are also combinations and variants of the basic *Send* and *Receive* calls. An important variant is the *Call* operation, which consists of a standard *Send* operation atomically followed by a variant of *Receive* which waits for a *Reply*. A *reply* message is always delivered via a special resource instead of using the standard IPC mechanism; see `seL4_Call()` below for details.

Invoking methods on kernel objects other than endpoints and notifications is done with *Send* or *Call*, depending on whether the invoker wants a reply from the kernel (*Call*) or not (*Send*). By using functions provided by the libsel4 API you are guaranteed to always use the more appropriate one. The *Yield* system call is not associated with any kernel object and is the only operation that does not invoke a capability. In the MCS configuration, *Wait* is a variant of *Receive* that does not require a reply object to be provided—on non-MCS configurations, *Wait* is synonymous with *Receive*, because neither call takes a reply object.

The fundamental system calls are:

seL4_Yield() is the only system call that does not require a capability to be used.

It forfeits the remainder of the calling thread’s timeslice and causes invocation of the kernel’s scheduler. If there are no other runnable threads with the same priority as the caller, the calling thread will immediately be scheduled with a fresh timeslice. In the MCS configuration, this behaviour depends on the state of the scheduling context; see Section 6.1.6.

seL4_Send() delivers a message through the named capability. If the invoked capability is an endpoint, and no receiver is ready to receive the message immediately, the sending thread will block until the message can be delivered. No error code or response will be returned by the receiving object.

seL4_Recv() (“receive”) is used by a thread to receive messages through endpoints or notifications. If no sender or notification is pending, the caller will block until a message or notification can be delivered. This system call works only on *Endpoint* or *Notification* capabilities, raising a fault (see section 6.2) when attempted with other capability types.

In the MCS configuration, *Receive* takes a reply capability—a capability to a reply object—as a parameter.

The remaining system calls are variants and combinations of `seL4_Send()` and `seL4_Recv()` efficiently accommodate common use cases in systems programming.

seL4_NBSend() performs a polling send on an endpoint. If the message cannot be delivered immediately, i.e., there is no receiver waiting on the destination *Endpoint*, the message is silently dropped. The sending thread continues execution. As with `seL4_Send()`, no error code or response will be returned.

seL4_NBRecv() is used by a thread to check for signals pending on a notification object or messages pending on an endpoint without blocking. This system call works only on endpoints and notification object capabilities, raising a fault (see section 6.2) when attempted with other capability types.

seL4_Call() combines **seL4_Send()** and **seL4_Recv()** with some important differences. The call blocks the sending thread until its message is delivered and a reply message is received.

When invoking capabilities to kernel services other than endpoints, using **seL4_Call()** allows the kernel to return an error code or other response through the reply message.

When the sent message is delivered to another thread via an **Endpoint**, the kernel does the same operation as **seL4_Send()**. What happens next depends on the kernel configuration. For MCS configurations, the kernel then updates the *reply object* provided by the receiver. A *reply object* is a vessel for tracking reply messages, used to send a reply message and wake up the caller. In non-MCS configurations, the kernel then deposits a special *reply capability* in a dedicated slot in the receiver's TCB. This *reply capability* is a single-use right to send a reply message and wake up the caller, meaning that the kernel invalidates it as soon as it has been invoked. For both variants, the calling thread is blocked until a capability to the reply object is invoked. For more information, see Section 4.2.4.

seL4_Reply() is used to respond to a **seL4_Call()**, by invoking the reply capability generated by the **seL4_Call()** system call and stored in a dedicated slot in the replying thread's TCB. It has exactly the same behaviour as invoking the reply capability with **seL4_Send()** which is described in Section 4.2.4.

seL4_ReplyRecv() combines **seL4_Reply()** and **seL4_Recv()**. It exists mostly for efficiency reasons, namely the common case of replying to a request and waiting for the next can be performed in a single kernel system call instead of two. The transition from the reply to the receive phase is also atomic.

seL4_Wait() works like **seL4_Recv()**; on non-MCS configurations, they are in fact synonymous. In the MCS configuration, **seL4_Wait()** is used when no reply is expected. Unlike **seL4_Recv()**, **seL4_Wait()** takes no reply capability.

seL4_NBWait() (MCS only) is used by a thread to poll for messages through endpoints or notifications. If no sender or notification is pending, the system call returns immediately.

seL4_NBSendWait() (MCS only) combines an **seL4_NBSend()** and **seL4_Wait()** into one atomic system call.

seL4_NBSendRecv() (MCS only) combines an **seL4_NBSend()** and **seL4_Recv()** into one atomic system call.

2.3 Kernel Objects

In this section we give a brief overview of the kernel-implemented object types whose instances (also simply called *objects*) can be invoked by applications. The interface to

these objects forms the interface to the kernel itself. The creation and use of kernel services is achieved by the creation, manipulation, and combination of these kernel objects:

CNodes (see Chapter 3) store capabilities, giving threads permission to invoke methods on particular objects. Each CNode has a fixed number of slots, always a power of two, determined when the CNode is created. Slots can be empty or contain a capability.

Thread Control Blocks (TCBs; see Chapter 6) represent a thread of execution in seL4. Threads are the unit of execution that is scheduled, blocked, unblocked, etc., depending on the application's interaction with other threads.

Scheduling contexts (MCS only) (SchedulingContexts; see Chapter 6) represent CPU time in seL4. Users can create scheduling contexts from untyped objects, however on creation scheduling contexts are *empty* and do not represent any time. Initially, there is a capability to SchedControl for each node, which allows scheduling context to be populated with parameters, which combined with priority control thread's access to CPU time.

Endpoints (see Chapter 4) facilitate message-passing communication between threads. IPC is synchronous: A thread trying to send or receive on an endpoint blocks until the message can be delivered. This means that message delivery only happens if a sender and a receiver rendezvous at the endpoint, and the kernel can deliver the message with a single copy (or without copying for short messages using only registers).

A capability to an endpoint can be restricted to be send-only or receive-only. Additionally, Endpoint capabilities can have the grant right, which allows sending capabilities as part of the message.

Reply objects (MCS only) (see Chapter 4) track scheduling context donation and provide a container for single-use reply capabilities. They are provided by `seL4_Recv()`.

Notification Objects (see Chapter 5) provide a simple signalling mechanism. A Notification is a word-size array of flags, each of which behaves like a binary semaphore. Operations are *signalling* a subset of flags in a single operation, polling to check any flags, and blocking until any are signalled. Notification capabilities can be signal-only or wait-only.

Virtual Address Space Objects (see Chapter 7) are used to construct a virtual address space (or VSpace) for one or more threads. These objects largely directly correspond to those of the hardware, and as such are architecture-dependent. The kernel also includes ASID Pool and ASID Control objects for tracking the status of address spaces.

Interrupt Objects (see Chapter 8) give applications the ability to receive and acknowledge interrupts from hardware devices. Initially, there is a capability to IRQControl, which allows for the creation of IRQHandler capabilities. An IRQHandler capability permits the management of a specific interrupt source associated with

a specific device. It is delegated to a device driver to access an interrupt source. The `IRQHandler` object allows threads to wait for and acknowledge individual interrupts.

Untyped Memory (see Section 2.4) is the foundation of memory allocation in the seL4 kernel. Untyped memory capabilities have a single method which allows the creation of new kernel objects. If the method succeeds, the calling thread gains access to capabilities to the newly-created objects. Additionally, untyped memory objects can be divided into a group of smaller untyped memory objects allowing delegation of part (or all) of the system's memory. We discuss memory management in general in the following sections.

2.4 Kernel Memory Allocation

The seL4 microkernel does not dynamically allocate memory for kernel objects. Instead, objects must be explicitly created from application-controlled memory regions via **Untyped Memory** capabilities. Applications must have explicit authority to memory (through these **Untyped Memory** capabilities) in order to create new objects, and all objects consume a fixed amount of memory once created. These mechanisms can be used to precisely control the specific amount of physical memory available to applications, including being able to enforce isolation of physical memory access between applications or a device. There are no arbitrary resource limits in the kernel apart from those dictated by the hardware², and so many denial-of-service attacks via resource exhaustion are avoided.

At boot time, seL4 pre-allocates the memory required for the kernel itself, including the code, data, and stack sections (seL4 is a single kernel-stack operating system). It then creates an initial user thread (with an appropriate address and capability space). The kernel then hands all remaining memory to the initial thread in the form of capabilities to **Untyped Memory**, and some additional capabilities to kernel objects that were required to bootstrap the initial thread. These **Untyped Memory** regions can then be split into smaller regions or other kernel objects using the `seL4_Untyped_Retype()` method; the created objects are termed *children* of the original untyped memory object.

The user-level application that creates an object using `seL4_Untyped_Retype()` receives full authority over the resulting object. It can then delegate all or part of the authority it possesses over this object to one or more of its clients.

Untyped memory objects represent two different types of memory: general purpose memory, or device memory. *General purpose* memory can be untyped into any other object type and used for any operation on untyped memory provided by the kernel. *Device memory* covers memory regions reserved for devices as determined by the hardware platform, and usage of these objects is restricted by the kernel in the following ways:

- Device untyped objects can only be retyped into frames or other untyped objects; developers cannot, for example, create an endpoint from device memory.

²The treatment of virtual ASIDs imposes a fixed number of address spaces. This limitation is to be removed in future versions of seL4.

- Frame objects retyped from device untyped objects cannot be set as thread IPC buffers, or used in the creation of an ASID pool

The type attribute (whether it represents *general purpose* or *device* memory) of a child untyped object is inherited from its parent untyped object. That is, any child of a device untyped will also be a device untyped. Developers cannot change the type attribute of an untyped.

2.4.1 Reusing Memory

The model described thus far is sufficient for applications to allocate kernel objects, distribute authority among client applications, and obtain various kernel services provided by these objects. This alone is sufficient for a simple static system configuration.

The seL4 kernel also allows Untyped Memory regions to be reused. Reusing a region of memory is allowed only when there are no dangling references (i.e., capabilities) left to the objects inside that memory. The kernel tracks *capability derivations*, i.e., the children generated by the methods `seL4_Untyped_Retype()`, `seL4_CNode_Mint()`, `seL4_CNode_Copy()`, and `seL4_CNode_Mutate()`.

The tree structure so generated is termed the *capability derivation tree* (CDT).³ For example, when a user creates new kernel objects by retyping untyped memory, the newly created capabilities would be inserted into the CDT as children of the untyped memory capability.

For each Untyped Memory region, the kernel keeps a *watermark* recording how much of the region has previously been allocated. Whenever a user requests the kernel to create new objects in an untyped memory region, the kernel will carry out one of two actions: if there are already existing objects allocated in the region, the kernel will allocate the new objects at the current watermark level, and increase the watermark. If all objects previously allocated in the region have been deleted, the kernel will reset the watermark and start allocating new objects from the beginning of the region again.

Finally, the `seL4_CNode_Revoke()` method provided by CNode objects destroys all capabilities derived from the argument capability. Revoking the last capability to a kernel object triggers the *destroy* operation on the now unreferenced object. This simply cleans up any in-kernel dependencies between it, other objects and the kernel.

By calling `seL4_CNode_Revoke()` on the original capability to an untyped memory object, the user removes all of the untyped memory object's children—that is, all capabilities pointing to objects in the untyped memory region. Thus, after this invocation there are no valid references to any object within the untyped region, and the region may be safely retyped and reused.

2.4.2 Summary of Object Sizes

When retyping untyped memory it is useful to know how much memory the object will require. Object sizes are defined in `libsel4`.

³Although the CDT conceptually is a separate data structure, it is implemented as part of the CNode object and so requires no additional kernel meta-data.

Note that CNodes, SchedContexts (MCS only), and Untyped Objects have variable sizes. When retyping untyped memory into CNodes or SchedContexts, or breaking an Untyped Object into smaller Untyped Objects, the `size_bits` argument to `seL4_Untyped_Retype()` is used to specify the size of the resulting objects. For all other object types, the size is fixed, and the `size_bits` argument to `seL4_Untyped_Retype()` is ignored.

Type	Meaning of <code>size_bits</code>	Size in Bytes
CNode	\log_2 number of slots	$2^{\text{size_bits}} \cdot 2^{\text{seL4_SlotBits}}$ <code>seL4_SlotBits</code> is: <i>on 32-bit architectures: 4</i> <i>on 64-bit architectures: 5</i>
SchedContext (MCS only)	\log_2 size in bytes	$2^{\text{size_bits}}$
Untyped	\log_2 size in bytes	$2^{\text{size_bits}}$

Table 2.1: Meaning of `size_bits` for object types of variable size

A single call to `seL4_Untyped_Retype()` can retype a single Untyped Object into multiple objects. The number of objects to create is specified by its `num_objects` argument. All created objects must be of the same type, specified by the `type` argument. In the case of variable-sized objects, each object must also be of the same size. If the size of the memory area needed (calculated by the object size multiplied by `num_objects`) is greater than the remaining unallocated memory of the Untyped Object, an error will result.

Chapter 3

Capability Spaces

Recall from Section 2.1 that seL4 implements a capability-based access control model. Each userspace thread has an associated *capability space* (CSpace) that contains the capabilities that the thread possesses, thereby governing which resources the thread can access.

Recall that capabilities reside within kernel-managed objects known as CNodes. A CNode is a table of slots, each of which may contain a capability. This may include capabilities to further CNodes, forming a directed graph. Conceptually a thread's CSpace is the portion of the directed graph that is reachable starting with the CNode capability that is its CSpace root.

A CSpace address refers to an individual slot (in some CNode in the CSpace), which may or may not contain a capability. Threads refer to capabilities in their CSpaces (e.g. when making system calls) using the address of the slot that holds the capability in question. An address in a CSpace is the concatenation of the indices of the CNode capabilities forming the path to the destination slot; we discuss this further in Section 3.3.

Recall that capabilities can be copied and moved within CSpaces, and also sent in messages (message sending will be described in detail in Section 4.2.2). Furthermore, new capabilities can be *minted* from old ones with a subset of their rights. Recall, from Section 2.4.1, that seL4 maintains a *capability derivation tree* (CDT) in which it tracks the relationship between these copied capabilities and the originals. The revoke method removes all capabilities (in all CSpaces) that were derived from a selected capability. This mechanism can be used by servers to restore sole authority to an object they have made available to clients, or by managers of untyped memory to destroy the objects in that memory so it can be retyped.

seL4 requires the programmer to manage all in-kernel data structures, including CSpaces, from userspace. This means that the userspace programmer is responsible for constructing CSpaces as well as addressing capabilities within them. This chapter first discusses capability and CSpace management, before discussing how capabilities are addressed within CSpaces, i.e. how applications can refer to individual capabilities within their CSpaces when invoking methods.

3.1 Capability and CSpace Management

3.1.1 CSpace Creation

CSpaces are created by creating and manipulating CNode objects. When creating a CNode the user must specify the number of slots that it will have, and this determines the amount of memory that it will use. Each slot requires $2^{\text{seL4_SlotBits}}$ bytes of physical memory and has the capacity to hold exactly one capability. This is 16 bytes on 32-bit architectures and 32 bytes on 64-bit architectures. Like any other object, a CNode must be created by calling `seL4_Untyped_Retype()` on an appropriate amount of untyped memory (see Section 2.4.2). The caller must therefore have a capability to enough untyped memory as well as enough free capability slots available in existing CNodes for the `seL4_Untyped_Retype()` invocation to succeed.

3.1.2 CNode Methods

Capabilities are managed largely through invoking CNode methods.

CNodes support the following methods:

`seL4_CNode_Mint()` creates a new capability in a specified CNode slot from an existing capability. The newly created capability may have fewer rights than the original and a different guard (see Section 3.3.1). `seL4_CNode_Mint()` can also create a badged capability (see Section 4.2.1) from an unbadged one.

`seL4_CNode_Copy()` is similar to `seL4_CNode_Mint()`, but the newly created capability has the same badge and guard as the original.

`seL4_CNode_Move()` moves a capability between two specified capability slots. You cannot move a capability to the slot in which it is currently.

`seL4_CNode_Mutate()` can move a capability similarly to `seL4_CNode_Move()` and also reduce its rights similarly to `seL4_CNode_Mint()`, although without an original copy remaining.

`seL4_CNode_Rotate()` moves two capabilities between three specified capability slots. It is essentially two `seL4_CNode_Move()` invocations: one from the second specified slot to the first, and one from the third to the second. The first and third specified slots may be the same, in which case the capability in it is swapped with the capability in the second slot. The method is atomic; either both or neither capabilities are moved.

`seL4_CNode_Delete()` removes a capability from the specified slot.

`seL4_CNode_Revoke()` is equivalent to calling `seL4_CNode_Delete()` on each derived child of the specified capability. It has no effect on the capability itself, except in very specific circumstances outlined in Section 3.2.

`seL4_CNode_SaveCaller()` moves a kernel-generated reply capability of the current thread from the special TCB slot it was created in, into the designated CSpace slot (non-MCS only).

`seL4_CNode_CancelBadgedSends()` cancels any outstanding sends that use the same badge and object as the specified capability.

3.1.3 Capabilities to Newly-Retyped Objects

When retyping untyped memory into objects with `seL4_Untyped_Retype()`, capabilities to the newly-retyped objects are placed in consecutive slots in a `CNode` specified by its `root`, `node_index`, and `node_depth` arguments. The `node_offset` argument specifies the index into the `CNode` at which the first capability will be placed. The `num_objects` argument specifies the number of capabilities (and, hence, objects) to create. All slots must be empty or an error will result. All resulting capabilities will be placed in the same `CNode`.

3.1.4 Capability Rights

As mentioned previously, some capability types have *access rights* associated with them. Currently, access rights are associated with capabilities for Endpoints (see Chapter 4), Notifications (see Chapter 5), Pages (see Chapter 7) and Replying (see Chapter 4). The access rights associated with a capability determine the methods that can be invoked. `seL4` supports four access rights, which are Read, Write, Grant and GrantReply. Read, Write and Grant are orthogonal to each other. GrantReply is a less powerful form of Grant e.g. if you already have Grant, having GrantReply or not is irrelevant. The meaning of each right is interpreted relative to the various object types, as detailed in Table 3.1.

When an object is first created, the initial capability that refers to it carries the maximum set of access rights. Other, less-powerful capabilities may be manufactured from this original capability, using methods such as `seL4_CNode_Mint()` and `seL4_CNode_Mutate()`. If a greater set of rights than the source capability is specified for the destination capability in either of these invocations, the destination rights are silently downgraded to those of the source.

Type	Read	Write	Grant	GrantReply
Endpoint	Receiving	Sending	Sending any capabilities	Sending reply capabilities
Notification	Waiting	Signaling	N/A	N/A
Page	Mapping the page readable.	Mapping the page writable.	N/A	N/A
Reply	N/A	N/A	Sending any capabilities in reply message	N/A

Table 3.1: `seL4` access rights: What a specific right entitles a capability to do

3.1.5 Capability Derivation Tree

As mentioned in Section 2.4.1, seL4 keeps track of capability derivations in a capability derivation tree.

Various methods, such as `seL4_CNode_Copy()` or `seL4_CNode_Mint()`, may be used to create derived capabilities. Not all capabilities support derivation. In general, only *original* capabilities support derivation invocations, but there are exceptions. Table 3.2 summarises the conditions that must be met for capability derivation to succeed for the various capability types, and how capability-derivation failures are reported in each case. The capability types not listed can be derived once.

Cap Type	Conditions for Derivation	Error Code on Derivation Failure
ReplyCap	Cannot be derived	Dependent on syscall
IRQControl	Cannot be derived	Dependent on syscall
Untyped	Must not have children (Section 3.2)	<code>seL4_RevokeFirst</code>
Page Table	Must be mapped	<code>seL4_IllegalOperation</code>
Page Directory	Must be mapped	<code>seL4_IllegalOperation</code>
IO Page Table (IA-32 only)	Must be mapped	<code>seL4_IllegalOperation</code>

Table 3.2: Capability derivation.

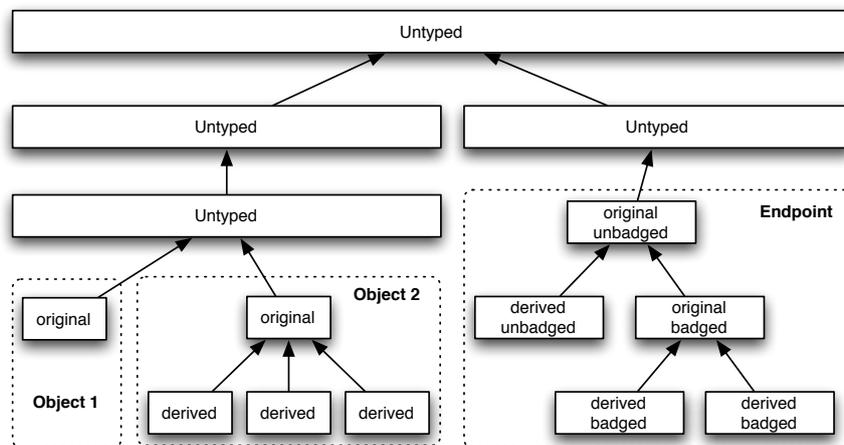


Figure 3.1: Example capability derivation tree.

Figure 3.1 shows an example capability derivation tree that illustrates a standard scenario: the top level is a large untyped capability, the second level splits this capability into two regions covered by their own untyped caps, both are children of the first level. The third level on the left is a copy of the level 2 untyped capability. Untyped capabilities when copied always create children, never siblings. In this scenario, the untyped capability was typed into two separate objects, creating two capabilities on level 4, both are the original capability to the respective object, both are children of

the untyped capability they were created from.

Ordinary original capabilities can have one level of derived capabilities. Further copies of these derived capabilities will create siblings, in this case remaining on level 5. There is an exception to this scheme for `Endpoint` and `Notification` capabilities — they support an additional layer of depth though *badging*. The original `Endpoint` or `Notification` capability will be unbadged. Using the `mint` method, a copy of the capability with a specific *badge* can be created (see Section 4.2.1, Section 5.1). This new, badged capability to the same object is treated as an original capability (the “original badged endpoint capability”) and supports one level of derived children like other capabilities.

3.2 Deletion and Revocation

Capabilities in `seL4` can be deleted and revoked. Both methods primarily affect capabilities, but they can have side effects on objects in the system where the deletion or revocation results in the destruction of the last capability to an object.

As described above, `seL4_CNode_Delete()` will remove a capability from the specified `CNode` slot. Usually, this is all that happens. If, however, it was the last typed capability to an object, this object will now be destroyed by the kernel, cleaning up all remaining in-kernel references and preparing the memory for re-use.

If the object to be destroyed was a capability container, i.e. a `TCB` or `CNode`, the destruction process will delete each capability held in the container, prior to destroying the container. This may result in the destruction of further objects if the contained capabilities are the last capabilities.¹

The `seL4_CNode_Revoke()` method will `seL4_CNode_Delete()` all `CDT` children of the specified capability, but will leave the capability itself intact. If any of the revoked child capabilities were the last capabilities to an object, the appropriate destroy operation is triggered.

Note: `seL4_CNode_Revoke()` may only partially complete in two specific circumstances. The first being where a `CNode` containing the last capability to the `TCB` of the thread performing the revoke (or the last capability to the `TCB` itself) is deleted as a result of the revoke. In this case the thread performing the revoke is destroyed during the revoke and the revoke does not complete. The second circumstance is where the storage containing the capability that is the target of the revoke is deleted as a result of the revoke. In this case, the authority to perform the revoke is removed during

¹The recursion is limited as if the last capability to a `CNode` is found within the container, the found `CNode` is not destroyed. Instead, the found `CNode` is made unreachable by moving the capability pointing to the found `CNode` into the found `cnode` itself, by swapping the capability with the first capability in the found `cnode`, and then trying to delete the swapped capability instead. This breaks the recursion.

The result of this approach is that deleting the last cap to the root `CNode` of a `Cspace` does not recursively delete the entire `Cspace`. Instead, it deletes the root `CNode`, and the branches of the tree become unreachable, potentially including the deleting of some of the unreachable `CNode`’s caps to make space for the self-referring capability. The practical consequence of this approach is that `Cspace` deletion requires user-level to delete the tree leaf first if unreachable `CNodes` are to be avoided. Alternatively, any resulting unreachable `CNodes` can be cleaned up via revoking a covering untyped capability, however this latter approach may be more complex to arrange by construction at user-level.

the operation and the operation stops part way through. Both these scenarios can be and should be avoided at user-level by construction.

Note that for page tables and page directories `seL4_CNode_Revoke()` will not revoke frame capabilities mapped into the address space. They will only be unmapped from the space.

3.3 CSpace Addressing

When performing a system call, a thread specifies to the kernel the capability to be invoked by giving an address in its CSpace. This address refers to the specific slot in the caller's CSpace that contains the capability to be invoked.

CSpaces are designed to permit sparsity, and the process of looking-up a capability address must be efficient. Therefore, CSpaces are implemented as *guarded page tables*.

As explained earlier, a CSpace is a directed graph of CNode objects, and each CNode is a table of slots, where each slot can either be empty, or contain a capability, which may refer to another CNode. Recall from Section 2.3 that the number of slots in a CNode must be a power of two. A CNode is said to have a *radix*, which is the power to which two is raised in its size. That is, if a CNode has 2^k slots, its radix would be k . The kernel stores a capability to the root CNode of each thread's CSpace in the thread's TCB. Conceptually, a CNode capability stores not only a reference to the CNode to which it refers, but also carries a *guard* value, explained in Section 3.3.1.

3.3.1 Capability Address Lookup

Like a virtual memory address, a capability address is simply an integer. Rather than referring to a location of physical memory (as does a virtual memory address), a capability address refers to a capability slot. When looking up a capability address presented by a userspace thread, the kernel first consults the CNode capability in the thread's TCB that defines the root of the thread's CSpace. It then compares that CNode's *guard* value against the most significant bits of the capability address. If the two values are different, lookup fails. Otherwise, the kernel then uses the next most-significant *radix* bits of the capability address as an index into the CNode to which the CNode capability refers. The slot s identified by these next *radix* bits might contain another CNode capability or contain something else (including nothing). If s contains a CNode capability c and there are remaining bits (following the *radix* bits) in the capability address that have yet to be translated, the lookup process repeats, starting from the CNode capability c and using these remaining bits of the capability address. Otherwise, the lookup process terminates successfully; the capability address in question refers to the capability slot s .

Figure 3.2 demonstrates a valid CSpace with the following features:

- a top level CNode object with a 12-bit guard set to 0x000 and 256 slots;
- a top level CNode with direct object references;
- a top level CNode with two second-level CNode references;

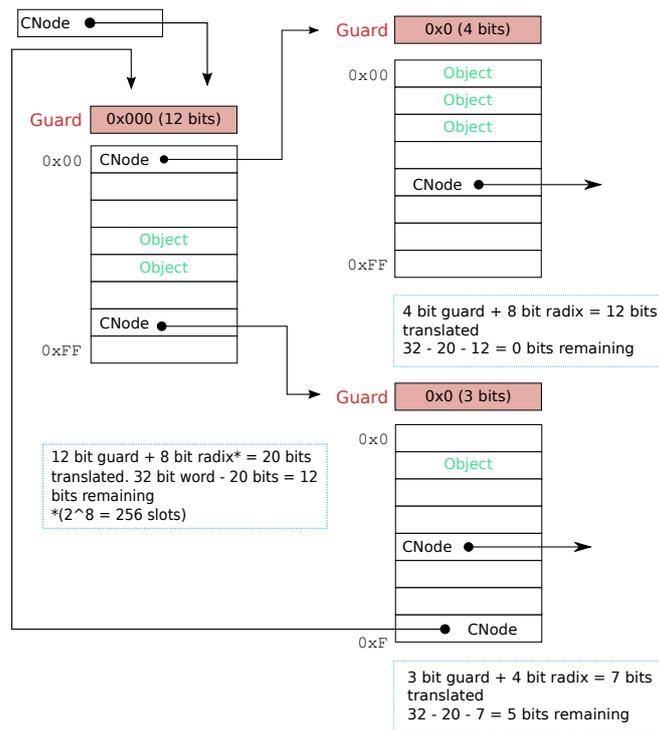


Figure 3.2: An example CSpace demonstrating object references at all levels, various guard and radix sizes and internal CNode references.

- second level CNodes with different guards and slot counts;
- a second level CNode that contains a reference to a top level CNode;
- a second level CNode that contains a reference to another CNode where there are some bits remaining to be translated;
- a second level CNode that contains a reference to another CNode where there are no bits remaining to be translated; and
- object references in the second level CNodes.

It should be noted that Figure 3.2 demonstrates only what is possible, not what is usually practical. Although the CSpace is legal, it would be reasonably difficult to work with due to the small number of slots and the circular references within it.

3.3.2 Addressing Capabilities

A capability address is stored in a CPointer (abbreviated CPTR), which is an unsigned integer variable. Capabilities are addressed in accordance with the translation algorithm described above. Two special cases involve addressing CNode capabilities themselves and addressing a range of capability slots.

Recall that the translation algorithm described above will traverse CNode capabilities while there are address bits remaining to be translated. Therefore, in order to address a

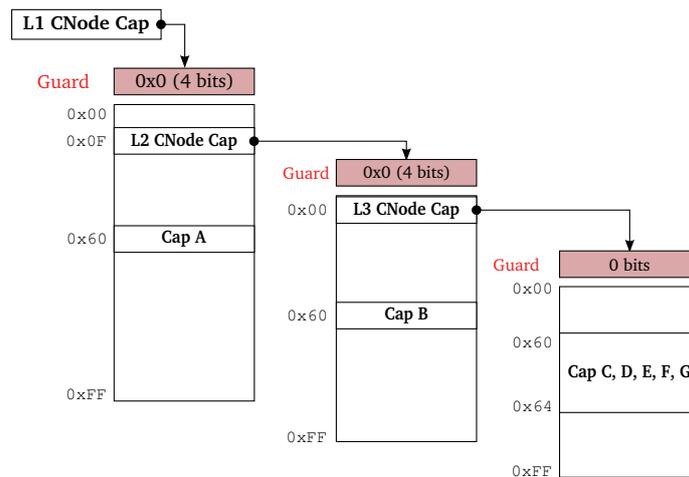


Figure 3.3: An arbitrary CSpace layout.

capability which may be a CNode capability, the user must supply not only a capability address but also specify the maximum number of bits of the capability address that are to be translated, called the *depth limit*. When a CPointer is paired with depth limit *depth*, only its *depth* least significant bits are used in translation.

Certain methods, such as `seL4_Untyped_Retype()`, require the user to provide a range of capability slots. This is done by providing a base capability address, which refers to the first slot in the range, together with a window size parameter, specifying the number of slots (with consecutive addresses, following the base slot) in the range.

Figure 3.3 depicts an example CSpace. In order to illustrate these ideas, we determine the address of each of the 10 capabilities in this CSpace.

Cap A. The first CNode has a 4-bit guard set to 0x0, and an 8-bit radix. Cap A resides in slot 0x60 so, provided that it is not a CNode capability, it may be referred to by any address of the form 0x060nnnnn (where *nnnnn* is any sequence of 5 hexadecimal digits, because the translation process terminates after translating the first 12 bits of the address). For simplicity, we usually set unused address bits to 0, which in this case yields the address 0x06000000.

Cap B. Again, the first CNode has a 4-bit guard set to 0x0, and an 8-bit radix. The second CNode is reached via the L2 CNode Cap. It also has a 4-bit guard of 0x0 and Cap B resides at index 0x60. Hence, Cap B's address is 0x00F06000. Translation of this address terminates after the first 24 bits.

Cap C. This capability is addressed via both CNodes. The third CNode is reached via the L3 CNode Cap, which resides at index 0x00 of the second CNode. The third CNode has no guard and Cap C is at index 0x60. Hence, its address is 0x00F00060. Translation of this address leaves 0 bits untranslated.

Caps C–G. This range of capability slots is addressed by providing a base address (which refers to the slot containing Cap C) of 0x00F00060 and a window size of 5.

L2 CNode Cap. Recall that to address a CNode capability, the user must supply not only a capability address but also specify the depth limit, which is the maximum number of bits to be translated. L2 CNode Cap resides at offset 0x0F of the first CNode, which has a 4-bit guard of 0x0. Hence, it may be referred to by any address of the form 0xnnnnn00F with a depth limit of 12 bits, where *nnnnn* is any sequence of 5 hexadecimal digits.

L3 CNode Cap. This capability resides at index 0x00 of the second CNode, which is reached by the L2 CNode Cap. The second CNode has a 4-bit guard of 0x0. Hence, the capability may be referred to by any address of the form 0xnn00F000 with a depth limit of 24 bits, where *nn* is any sequence of 2 hexadecimal digits.

In summary, to refer to any capability (or slot) in a CSpace, the user must supply its address. When the capability might be a CNode, the user must also supply a depth limit. To specify a range of capability slots, the user supplies a starting address and a window size.

3.4 Lookup Failure Description

When a capability lookup fails, a description of the failure is given to either the calling thread or the thread's exception handler in its IPC buffer. The format of the description is always the same but may occur at varying offsets in the IPC buffer depending on how the error occurred. The description format is explained below. The first word indicates the type of lookup failure and the meaning of later words depend on this.

3.4.1 Invalid Root

A CSpace CPTR root (within which a capability was to be looked up) is invalid. For example, the capability is not a CNode cap.

Data	Meaning
Offset + 0	seL4_InvalidRoot

3.4.2 Missing Capability

A capability required for an invocation is not present or does not have sufficient rights.

Data	Meaning
Offset + 0	seL4_MissingCapability
Offset + seL4_CapFault_BitsLeft	Bits left

3.4.3 Depth Mismatch

When resolving a capability, a CNode was traversed that resolved more bits than was left to decode in the CPTR or a non-CNode capability was encountered while there were still bits remaining to be looked up.

Data	Meaning
Offset + 0	seL4_DepthMismatch
Offset + seL4_CapFault_BitsLeft	Bits of CPTR remaining to decode
Offset + seL4_CapFault_DepthMismatch_BitsFound	Bits that the current CNode being traversed resolved

3.4.4 Guard Mismatch

When resolving a capability, a CNode was traversed with a guard size larger than the number of bits remaining or the CNode's guard did not match the next bits of the CPTR being resolved.

Data	Meaning
Offset + 0	seL4_GuardMismatch
Offset + seL4_CapFault_BitsLeft	Bits of CPTR remaining to decode
Offset + seL4_CapFault_GuardMismatch_GuardFound	The CNode's guard
Offset + seL4_CapFault_GuardMismatch_BitsFound	The CNode's guard size

Chapter 4

Message Passing (IPC)

The seL4 microkernel provides a message-passing IPC mechanism for communication between threads. The same mechanism is also used for communication with kernel-provided services. Messages are sent by invoking a capability to a kernel object. Messages sent to Endpoints are destined for other threads, while messages sent to other objects are processed by the kernel. This chapter describes the common message format, endpoints, and how they can be used for communication between applications.

4.1 Message Registers

Each message contains a number of message words and optionally a number of capabilities. The message words are sent to or received from a thread by placing them in its *message registers*. The message registers are numbered and the first few message registers are implemented using physical CPU registers, while the rest are backed by a fixed region of memory called the *IPC buffer*. The reason for this design is efficiency: very short messages need not use the memory. The IPC buffer is assigned to the calling thread (see Section 6.1 and Section 10.3.7.11).

Every IPC message also has a tag (structure `seL4_MessageInfo_t`). The tag consists of four fields: the label, message length, number of capabilities (the `extraCaps` field) and the `capsUnwrapped` field. The message length and number of capabilities determine either the number of message registers and capabilities that the sending thread wishes to transfer, or the number of message registers and capabilities that were actually transferred. The label is not interpreted by the kernel and is passed unmodified as the first data payload of the message. The label may, for example, be used to specify a requested operation. The `capsUnwrapped` field is used only on the receive side, to indicate the manner in which capabilities were received. It is described in Section 4.2.2.

The kernel assumes that the IPC buffer contains a structure of type `seL4_IPCBuffer` as defined in Table 4.1. The kernel uses as many physical registers as possible to transfer IPC messages. When more arguments are transferred than physical message registers are available, the kernel begins using the IPC buffer's `msg` field to transfer arguments. However, it leaves room in this array for the physical message registers. For example, if an IPC transfer or kernel object invocation required 4 message registers

Type	Name	Description
<code>seL4_MessageInfo_t</code>	<code>tag</code>	Message tag
<code>seL4_Word[]</code>	<code>msg</code>	Message contents
<code>seL4_Word</code>	<code>userData</code>	Base address of the structure, used by supporting user libraries
<code>seL4_CPtr[]</code> (<i>in</i>)	<code>caps</code>	Capabilities to transfer
<code>seL4_CapData_t[]</code> (<i>out</i>)	<code>badges</code>	Badges for endpoint capabilities received
<code>seL4_CPtr</code>	<code>receiveCNode</code>	CPTR to a CNode from which to find the receive slot
<code>seL4_CPtr</code>	<code>receiveIndex</code>	CPTR to the receive slot relative to <code>receiveCNode</code>
<code>seL4_Word</code>	<code>receiveDepth</code>	Number of bits of <code>receiveIndex</code> to use

Table 4.1: Fields of the `seL4_IPCBuffer` structure. Note that `badges` and `caps` use the same area of memory in the structure.

(and there are only 2 physical message registers available on this architecture) then arguments 1 and 2 would be transferred via message registers and arguments 3 and 4 would be in `msg[2]` and `msg[3]`. This allows the user-level object-invocation stubs to copy the arguments passed in physical registers to the space left in the `msg` array if desired. The situation is similar for the `tag` field. There is space for this field in the `seL4_IPCBuffer` structure, which the kernel ignores. User level stubs may wish to copy the message tag from its CPU register to this field, although the user level stubs provided with the kernel do not do this.

4.2 Endpoints

Endpoints allow a small amount of data and capabilities (namely the IPC buffer) to be transferred between two threads. Endpoint objects are invoked directly using the `seL4` system calls described in Section 2.2.

IPC Endpoints uses a rendezvous model and as such is synchronous and blocking. An Endpoint object may queue threads either to send or to receive. If no receiver is ready, threads performing the `seL4_Send()` or `seL4_Call()` system calls will wait in a queue for the first available receiver. Likewise, if no sender is ready, threads performing the `seL4_Recv()` system call or the second half of `seL4_ReplyRecv()` will wait for the first available sender.

Trying to Send or Call without the Write right will fail and return an error. In the case of Send the error is ignored (The kernel isn't allowed to reply). Thus there is no way of knowing that a send has failed because of missing right. On the other hand calling `seL4_Recv()` with a endpoint capability that does not have the Read right will raise a fault, see Section 6.2. This because otherwise the error message would be indistinguishable from a normal message received from another thread via the endpoint.

4.2.1 Endpoint Badges

Endpoint capabilities may be *minted* to create a new endpoint capability with a *badge* attached to it, a data word chosen by the invoker of the *mint* operation. When a message is sent to an endpoint using a badged capability, the badge is transferred to the receiving thread's `badge` register.

An endpoint capability with a zero badge is said to be *unbadged*. Such a capability can be badged with the `seL4_CNode_Mutate()` or `seL4_CNode_Mint()` invocations on the CNode containing the capability. Endpoint capabilities with badges cannot be unbadged, rebadged or used to create child capabilities with different badges.

On 32-bit platforms, only the low 28 bits of the badge are available for use. The kernel will silently ignore any usage of the high 4 bits. On 64-bit platforms, 64 bits are available for badges.

4.2.2 Capability Transfer

Messages may contain capabilities, which will be copied to the receiver, provided that the endpoint capability invoked by the sending thread has Grant rights. An attempt to send capabilities using an endpoint capability without the Grant right will result in transfer of the raw message, without any capability transfer.

Capabilities to be sent in a message are specified in the sending thread's IPC buffer in the `caps` field. Each entry in that array is interpreted as a CPTR in the sending thread's capability space. The number of capabilities to send is specified in the `extraCaps` field of the message tag.

The receiver specifies the slot in which it is willing to receive a capability, with three fields within the IPC buffer: `receiveCNode`, `receiveIndex` and `receiveDepth`. These fields specify the root CNode, capability address and number of bits to resolve, respectively, to find the slot in which to put the capability. Capability addressing is described in Section 3.3.2.

Note that receiving threads may specify only one receive slot, whereas a sending thread may include multiple capabilities in the message. Messages containing more than one capability may be interpreted by kernel objects. They may also be sent to receiving threads in the case where some of the extra capabilities in the message can be *unwrapped*.

If the *n*-th capability in the message refers to the endpoint through which the message is sent, the capability is *unwrapped*: its badge is placed into the *n*-th position of the receiver's badges array, and the kernel sets the *n*-th bit (counting from the least significant) in the `capsUnwrapped` field of the message tag. The capability itself is not transferred, so the receive slot may be used for another capability.

A capability that is not unwrapped is transferred by copying it from the sender's CNode slot to the receiver's CNode slot. The sender retains access to the sent capability.

If a receiver gets a message whose tag has an `extraCaps` of 2 and a `capsUnwrapped` of 2, then the first capability in the message was transferred to the specified receive slot and the second capability was unwrapped, placing its badge in `badges[1]`. There may

have been a third capability in the sender's message which could not be unwrapped.

4.2.3 Errors

Errors in capability transfers can occur at two places: in the send phase or in the receive phase. In the send phase, all capabilities that the caller is attempting to send are looked up to ensure that they exist before the send is initiated in the kernel. If the lookup fails for any reason, `seL4_Send()` and `seL4_Call()` system calls immediately abort and no IPC or capability transfer takes place. The system call will return a lookup failure error as described in Section 10.1.

In the receive phase, seL4 transfers capabilities in the order that they are found in the sending thread's IPC buffer `caps` array and terminates as soon as an error is encountered. Possible error conditions are:

- A source capability cannot be looked up. Although the presence of the source capabilities is checked when the sending thread performs the send system call, this error may still occur. The sending thread may have been blocked on the endpoint for some time before it was paired with a receiving thread. During this time, its CSpace may have changed and the source capability pointers may no longer be valid.
- The destination slot cannot be looked up. Unlike the send system call, the `seL4_Recv()` system call does not check that the destination slot exists and is empty before it initiates the receive. Hence, the `seL4_Recv()` system call will not fail with an error if the destination slot is invalid and will instead transfer badged capabilities until an attempt to save a capability to the destination slot is made.
- The capability being transferred cannot be derived. See Section 3.1.5 for details.

An error will not void the entire transfer, it will just end it prematurely. The capabilities processed before the failure are still transferred and the `extraCaps` field in the receiver's IPC buffer is set to the number of capabilities transferred up to failure. No error message will be returned to the receiving thread in any of the above cases.

4.2.4 Calling and Replying

As explained in Section 2.2, when the user calls `seL4_Call()` on an endpoint capability, some specific actions are taken. First a call will do exactly the same action as a normal `seL4_Send()`. Then after the rendezvous and all the normal IPC procedure happened, instead of returning directly to the caller, `seL4_Call()` will check if either Grant or GrantReply are present on the invoked endpoint capability:

- If this is not the case, the caller thread is suspended as if `seL4_TCB_Suspend()` was called on it. The send part of the call would still have been performed as usual.

- If this is the case. A reply capability is set in a specific slot of the receiver TCB. The Grant right of that reply capability is set by copying the Grant right of the endpoint capability invoked by the receiver in `seL4_Recv()`. Then, the caller thread is blocked waiting for the reply.

A reply capability points directly to the caller thread and once the call has been performed is completely unrelated to the original Endpoint. Even if the latter was destroyed, the reply capability would still exist and point to the caller who would still be waiting for a reply.

The generated reply capability can then be either invoked in place (in the specific TCB slot) with the `seL4_Reply()` or saved to an addressable slot using `seL4_CNode_SaveCaller()` to be invoked later with `seL4_Send()`. The specific slot cannot be directly addressed with any CPtr as it is not part of any CSpace.

A reply capability is invoked in the same way as a normal send on a Endpoint. A reply capability has implicitly the Write right, so the message will always go through. Transferring caps in the reply can only happen if the reply capability has the Grant right and is done in exactly the same way as in a normal IPC transfer as described in Section 4.2.2.

The main difference with a normal endpoint transfer is that the kernel guarantees that invoking a reply capability cannot block: If you own a reply capability, then the thread it points to is waiting for a reply. However a reply capability is a non-owning reference, contrary to all the other capabilities. That means that if the caller thread is destroyed or modified in any way that would render a reply impossible (for example being suspended with `seL4_TCB_Suspend()`), the kernel would immediately destroy the reply capability.

Once the reply capability has been invoked, the caller receives the message as if it has been performing a `seL4_Recv()` and just received the message. In particular, it starts running again.

The `seL4_Call()` operation exists not only for efficiency reasons (combining two operations into a single system call). It differs from `seL4_Send()` immediately followed by `seL4_Recv()` in ways that allow certain system setup to work much more efficiently with much less setup than with a traditional setup. In particular, it is guaranteed that the reply received by the caller comes from the thread that received the call without having to check any kind of badge.

Chapter 5

Notifications

Notifications are a simple, non-blocking signalling mechanism that logically represents a set of binary semaphores.

5.1 Notification Objects

A Notification object contains a single data word, called the *notification word*. Such an object supports two operations: `seL4_Signal()` and `seL4_Wait()`.

Notification capabilities can be badged, using `seL4_CNode_Mutate()` or `seL4_CNode_Mint()`, just like Endpoint capabilities (see Section 4.2.1). As with Endpoint capabilities, badged Notification capabilities cannot be unbadged, rebadged or used to create child capabilities with different badges.

5.2 Signalling, Polling and Waiting

The `seL4_Signal()` method updates the notification word by bit-wise `or`-ing it with the *badge* of the invoked notification capability. It also unblocks the first thread waiting on the notification (if any). As such, `seL4_Signal()` works like concurrently signalling multiple semaphores (those indicated by the bits set in the badge). If the signal sender capability was unbadged or 0-badged, the operation degrades to just waking up the first thread waiting on the notification (also see below).

The `seL4_Wait()` method works similarly to a `select`-style wait on the set of semaphores: If the notification word is zero at the time `seL4_Wait()` is called, the invoker blocks. Else, the call returns immediately, setting the notification word to zero and returning to the invoker the previous notification-word value.

The `seL4_Poll()` is the same as `seL4_Wait()`, except if no signals are pending (the notification word is 0) the call will return immediately without blocking.

If threads are waiting on the Notification object at the time `seL4_Signal()` is invoked, the first queued thread receives the notification. All other threads keep waiting until the next time the notification is signalled.

5.3 Binding Notifications

Notification objects and TCBs can be bound together in a 1-to-1 relationship through the `seL4_TCB_BindNotification()` invocation. When a Notification is bound to a TCB, signals to that notification object will be delivered even if the thread is receiving from an IPC endpoint. To distinguish whether the received message was a notification or an IPC, developers should check the badge value. By reserving a specific badge (or range of badges) for capabilities to the bound notification — distinct from endpoint badges — the message source can be determined.

Once a notification has been bound, the only thread that may perform `seL4_Wait()` on the notification is the bound thread.

Chapter 6

Threads and Execution

6.1 Threads

seL4 provides threads to represent an execution context. On MCS configurations of the kernel, scheduling contexts are used to manage processor time. Without MCS, processor time is also represented by the thread abstraction. A thread is represented in seL4 by its thread control block object (TCB).

With MCS, a scheduling context is represented by a scheduling context object (SCO), and threads cannot run unless they are bound to, or receive a scheduling context.

6.1.1 Thread control blocks

Each TCB has an associated CSpace (see Chapter 3) and VSpace (see Chapter 7) which may be shared with other threads. A TCB may also have an IPC buffer (see Chapter 4), which is used to pass extra arguments during IPC or kernel object invocation that do not fit in the architecture-defined message registers. While it is not compulsory that a thread has an IPC buffer, it will not be able to perform most kernel invocations, as they require cap transfer. Each thread belongs to exactly one security domain (see Section 6.3).

6.1.2 Thread Creation

Like other objects, TCBs are created with the `seL4_Untyped_Retype()` method (see Section 2.4). A newly created thread is initially inactive. It is configured by setting its CSpace and VSpace with the `seL4_TCB_SetSpace()` or `seL4_TCB_Configure()` methods and then calling `seL4_TCB_WriteRegisters()` with an initial stack pointer and instruction pointer. The thread can then be activated either by setting the `resume_target` parameter in the `seL4_TCB_WriteRegisters()` invocation to true or by separately calling the `seL4_TCB_Resume()` method. Both of these methods place the thread in a runnable state.

On the master kernel, this will result in the thread immediately being added to the scheduler. On the MCS kernel, the thread will only begin running if it has a scheduling

context object.

In a SMP configuration of the kernel, the thread will resume on the core corresponding to the affinity of the thread. For master, this is set using `seL4_TCB_SetAffinity()`, while on the MCS kernel the affinity is derived from the scheduling context object.

6.1.3 Thread Deactivation

The `seL4_TCB_Suspend()` method deactivates a thread. Suspended threads can later be resumed. Their suspended state can be retrieved with the `seL4_TCB_ReadRegisters()` and `seL4_TCB_CopyRegisters()` methods. They can also be reconfigured and reused or left suspended indefinitely if not needed. Threads will be automatically suspended when the last capability to their TCB is deleted.

6.1.4 Scheduling

seL4 uses a preemptive, tickless scheduler with 256 priority levels (0 — 255). All threads have a maximum controlled priority (MCP) and a priority, the latter being the effective priority of the thread. When a thread modifies another thread's priority (including itself) it must provide a thread capability from which to use the MCP from. Threads can only set priorities and MCPs to be less than or equal to the provided thread's MCP. The initial task starts with an MCP and priority as the highest priority in the system (`seL4_MaxPrio`). Thread priority and MCP can be set with `seL4_TCB_SetSchedParams()` and `seL4_TCB_SetPriority()`, `seL4_TCB_SetMCPriority()` methods.

Of threads eligible for scheduling, the highest priority thread in a runnable state is chosen.

Thread priority (structure `seL4_PrioProps_t`) consists of two values as follows:

Priority the priority a thread will be scheduled with.

Maximum controlled priority (MCP) the highest priority a thread can set itself or another thread to.

6.1.5 MCS Scheduling

This section only applies to configurations with MCS enabled, where threads must have a scheduling context object available in order to be admitted to the scheduler.

6.1.6 Scheduling Contexts

Access to CPU execution time is controlled through scheduling context objects. Scheduling contexts are configured with a tuple of *budget*(b) and *period* (p), both in microseconds, set by `seL4_SchedControl_Configure()` (see Section 6.1.8). The tuple (b, p) forms an upper bound on the thread's execution – the kernel will not permit a thread to run for more than b out of every p microseconds. However, $\frac{b}{p}$ does not represent a

lower bound on execution, as a thread must have the highest or equal highest priority of all runnable threads to be guaranteed to be scheduled at all, and the kernel does not conduct an admission test. As a result the set of all parameters is not necessarily schedulable. If multiple threads have available budget concurrently they are scheduled first-in first-out, and round-robin scheduling is applied once the budget is expired.

A scheduling context that is eligible to be picked by the scheduler, i.e has budget available, is referred to as *active*. Budget charging and replenishment rules are different for round-robin and sporadic threads. For round-robin threads, the budget is charged each time the current node's scheduling context is changed, until it is depleted and then refilled immediately.

Threads where $b == p$ are treated as round robin threads, where b acts as a timeslice. Otherwise the kernel uses *sporadic servers* to enforce temporal isolation, which enforce the property that $\frac{b}{p}$ cannot be exceeded for all possible p . In theory, sporadic servers provide temporal isolation – preventing threads from exceeding their allocated budget – by using the following algorithm:

- When a thread starts executing at current time T , record T_s
- When a thread stops executing (blocks or is preempted), schedule a replenishment at $T_s + p$ for the amount of time consumed ($T - T_s$) and subtract that from the current replenishment being used.

seL4 implements this algorithm by maintaining an ordered list of sporadic replenishments – **refills** for brevity – in each scheduling context. Each replenishment contains a tuple of the time it is eligible for use (*rTime*) and the amount that replenishment is for (**rAmount**). While a thread is executing, it constantly drains the budget from the **rAmount** at the head of the replenishment list. If the **rTime** is in the future, the thread bound to that scheduling context is placed in a queue of threads waiting for more budget.

Round-robin threads are treated that same as sporadic threads except for one case: how the budget is charged. Round-robin threads have two refills only, both of which are always ready to be used. When a round-robin thread stops executing, budget is moved from the head to the tail replenishment. Once the head budget is consumed, the thread is removed from the scheduling queue for its priority and appended at the tail.

Sporadic threads behave differently depending on the amount of replenishments available, which must be bounded. Developers have two options to configure the size of the replenishment list:

- The maximum number of refills in a single scheduling context is determined by the size of the scheduling context when created by `seL4_Untyped_Retype()`.
- A per scheduling context parameter, **extra_refills** that limits the number of refills for that specific scheduling context. This value is added to the base value of 2 and is limited by the size of the scheduling context.

Threads that have short execution times (e.g. interrupt handlers) and are not frequently preempted should have less refills, while longer running threads with long values of b should have a higher value. Threads bound to a scheduling context with 0 extra refills will behave periodically – tasks that use their head replenishment, or call `yield`, will not be scheduled again until the start of their next period.

Given the number of replenishments is limited, if a node's SC changes and the outgoing SC does not have enough space to store the new replenishment, space is created by removing the current replenishment which can result in preemption if the next replenishment is not yet available. Scheduling contexts with a higher number of configured refills will consume closer to their whole budget, as they can be preempted or switch threads more often without filling their replenishment queue. However, the scheduling overhead will be higher as the replenishment list is subject to fragmentation.

Whenever a thread is executing it consumes the budget from its current scheduling context. The system call `seL4_Yield()` can be used to sacrifice any remaining budget and block until the next replenishment is ready to be used.

Threads can be bound to scheduling contexts using `seL4_TCB_Configure()` or `seL4_SchedContext_Bind()`, both invocations have the same effect although `seL4_TCB_Configure()` allows more thread fields to be set with only one kernel entry. When a thread is bound to a scheduling context, if it is in a runnable state and the scheduling context is active, it will be added to the scheduler.

6.1.7 Passive Threads

Threads can be unbound from a scheduling context with `seL4_SchedContext_UnbindObject()`. This is distinct from suspending a thread, in that threads that are blocked waiting in an endpoint or notification queue will remain in the queue and can still receive messages and signals. However, the unbound thread will not be schedulable again until it receives a scheduling context. Threads without scheduling contexts are referred to as *passive* threads, as they cannot execute without the action of another thread.

6.1.8 Scheduling Context Creation

Like other objects, scheduling contexts are created from untyped memory using `seL4_UntypedRetype()`. On creation, scheduling contexts are empty, representing 0% of CPU execution time. To populate a scheduling context with parameters, one must invoke the appropriate `SchedControl` capability, which provides access to CPU time management on a single node. A scheduling control cap for each node is provided to the initial task at run time. Threads run on the node that their scheduling context is configured for. Scheduling context parameters can then be set and updated using `seL4_SchedControl_ConfigureFlags()`, which allows the budget and period to be specified along with a bitwise OR'd set of the following flags.

`seL4_SchedContext_Sporadic` : constrain the execution time only according to the sporadic server algorithm rather than to a continuous constant bandwidth.

The kernel does not conduct any schedulability tests, as task admission is left to user-level policy and can be conducted online or offline, statically or dynamically or not at all.

6.1.9 Scheduling Context Donation

In addition to explicitly binding and removing scheduling contexts through `seL4_SchedContext_Bind()` and `seL4_SchedContext_UnbindObject()`, scheduling contexts can move between threads over IPC. Scheduling contexts are donated implicitly when the system calls `seL4_Call()` and `seL4_NBSendRecv()` are used to communicate with a passive thread. When an active thread invokes an endpoint with `seL4_Call()` and rendezvous with a passive thread, the active thread's scheduling context is donated to the passive thread. The generated reply cap ensures that the callee is merely borrowing the scheduling context: when the reply cap is consumed by a reply message being sent the scheduling context will be returned to the caller. If the reply cap is revoked, and the callee holds the scheduling context, the scheduling context will be returned to the caller. However, if in a deep call chain and a reply cap in the middle of the call chain is revoked, such that the callee does not possess the scheduling context, the thread will be removed from the call chain and the scheduling context will remain where it is. If the receiver does not provide a reply object to track the donation in (i.e uses `seL4_Wait()` instead of `seL4_Recv()`) scheduling context donation will not occur but the message will be delivered. The passive receiver will be set to inactive as it does not have a scheduling context.

Consider an example where thread A calls thread B which calls thread C. If whilst C holds the scheduling context, B's reply cap to A is revoked, then the scheduling context will remain with C. However, a call chain will remain between A and C, such that if C's reply cap is revoked, or invoked, the scheduling context will return to A.

`seL4_NBSendRecv()` can also result in scheduling context donation. If the non-blocking send phase of the operation results in message delivery to a passive thread, the scheduling context will be donated to that passive thread and the thread making the system call becomes passive on the receiving endpoint in the receive phase. No reply capability generated, so there is no guarantee that the scheduling context will return, which increases book keeping complexity but allows for data-flow like architectures rather than remote-procedure calls. Note that `seL4_Call()` does not guarantee the return of a scheduling context: this is an inherently trusted operation as the server could never reply and return the scheduling context.

Scheduling contexts can also be bound to notification objects using `seL4_SchedContext_Bind()` and unbound using `seL4_SchedContext_UnbindObject()`. If a signal is delivered to a notification object with a passive thread blocked waiting on it, the passive thread will receive the scheduling context that is bound to the notification object. The scheduling context is returned when the thread blocks on the notification object. This feature allows for passive servers to use notification binding (See Section 5.3).

Scheduling contexts can be unbound from all objects (notification objects and TCBS that are bound or have received a scheduling context through donation) using `seL4_SchedContext_Unbind()`.

Passive threads will run on the CPU node that the scheduling context was configured with, and will be migrated on IPC.

6.1.10 Scheduling algorithm

Threads are only eligible for scheduling if they have an active scheduling context. Of threads eligible for scheduling, the highest priority thread in a runnable state is chosen.

Threads of sufficient maximum controlled priority and with possession of the appropriate scheduling context capability can manipulate the scheduler and implement user-level schedulers using IPC.

Scheduling contexts provide access to and an upper bound on execution CPU time, however when a thread executes is determined by thread priority. Consequently, access to CPU is a function of thread MCPs, scheduling contexts and the SchedControl capability. The kernel will enforce that threads do not exceed the budget in their scheduling context for any given period, and that the highest priority thread will always run, however it is up to the system designer to make sure the entire system is schedulable.

6.1.11 Exceptions

Each thread has two associated exception-handler endpoints, a *standard* exception handler and a *timeout* exception handler, where the latter is MCS only. If the thread causes an exception, the kernel creates an IPC message with the relevant details and sends this to the endpoint. This thread can then take the appropriate action. Fault IPC messages are described in Section 6.2. Standard exception-handler endpoints can be set with the `seL4_TCB_SetSpace()` or `seL4_TCB_SetSchedParams()` methods while Timeout exception handlers can be set with `seL4_TCB_SetTimeoutEndpoint()` (MCS only). With these methods, a capability address for the exception handler can be associated with a thread. This address is then used to lookup the handler endpoint, and the capability to the endpoint is installed into the threads' kernel CNode. For threads without an exception handler, a null capability can be used, however the consequences are different per exception handler type. Before raising an exception the handler capability is validated. The kernel does not perform another lookup, but checks that the capability is an endpoint with the correct rights.

The exception endpoint must have Write and either Grant or GrantReply rights. Replying to the exception message restarts the thread. For certain exception types, the contents of the reply message may be used to set the values in the registers of the thread being restarted. See Section 6.2 for details.

6.1.11.1 Standard Exceptions

The standard exception handler is used when a fault is triggered by a thread which cannot be recovered without action by another thread. For example, if a thread raises a fault due to an unmapped virtual memory page, the thread cannot make any more progress until the page is mapped. If a thread experiences a fault that would trigger

the standard exception handler while it is set to a null capability, the kernel will pause the thread and it will not run again. This is because without action by another thread, standard exceptions cannot be recovered from. Consequently threads without standard exception handlers should be trusted not to fault at all.

Standard exception handlers can be passive, in which case they will run on the scheduling context of the faulting thread.

6.1.11.2 Timeout Exceptions (MCS Only)

Timeout faults are raised when a thread attempts to run but has no available budget, and if that thread has a valid timeout exception handler capability. The handling of timeout faults is not compulsory: if a thread does not have a timeout fault handler, a fault will not be raised and the thread will continue running when its budget is replenished. This allows temporally sensitive threads to handle budget overruns while other threads may ignore them.

Timeout faults are registered per thread, which means that while clients may not have a timeout fault handler, servers may, allowing single-threaded, time-sensitive, passive servers to use a timeout exception handler to recover from malicious or untrusted clients whose budget expires while the server is completing the request. Timeout fault handlers can access server reply objects and reply with an error to the client, then reset the server to handle the next client request.

If a reply message is sent to a nested server and a scheduling context without available budget returned, another timeout fault will be generated if the nested server also has a timeout fault handler.

6.1.12 Message Layout of the Read-/Write-Registers Methods

The registers of a thread can be read and written with the `seL4_TCB_ReadRegisters()` and `seL4_TCB_WriteRegisters()` methods. For some registers, the kernel will silently mask certain bits or ranges of bits off, and force them to contain certain values to ensure that they cannot be maliciously set to values that would compromise the running system, or to respect values that the architecture specifications have mandated to be certain values. The register contents are transferred via the IPC buffer.

6.2 Faults

A thread's actions may result in a fault. Faults are delivered to the thread's exception handler so that it can take the appropriate action. The fault type is specified in the message label and is one of: `seL4_Fault_CapFault`, `seL4_Fault_VMFault`, `seL4_Fault_UnknownSyscall`, `seL4_Fault_UserException`, `seL4_Fault_DebugException`, `seL4_Fault_TimeoutFault`, or `seL4_Fault_NullFault` (indicating no fault occurred and this is a normal IPC message).

Faults are delivered in such a way as to imitate a Call from the faulting thread. This means that to send a fault message the fault endpoint must have Write and either

Grant or GrantReply permissions. If this is not the case, a double fault happens (generally the thread is simply suspended).

6.2.1 Capability Faults

Capability faults may occur in two places. Firstly, a capability fault can occur when lookup of a capability referenced by a `seL4_Call()` or `seL4_Send()` system call failed (`seL4_NBSend()` calls on invalid capabilities silently fail). In this case, the capability on which the fault occurred may be the capability being invoked or an extra capability passed in the `caps` field in the IPC buffer.

Secondly, a capability fault can occur when `seL4_Recv()` or `seL4_NBRecv()` is called on a capability that does not exist, is not an endpoint or notification capability or does not have receive permissions.

Replying to the fault IPC will restart the faulting thread. The contents of the IPC message are given in Table 6.1.

Meaning	IPC buffer location
Address at which to restart execution	<code>seL4_CapFault_IP</code>
Capability address	<code>seL4_CapFault_Addr</code>
In receive phase (1 if the fault happened during a receive system call, 0 otherwise)	<code>seL4_CapFault_InRecvPhase</code>
Lookup failure description. As described in Section 3.4	<code>seL4_CapFault_LookupFailureType</code>

Table 6.1: Contents of an IPC message.

6.2.2 Unknown Syscall

This fault occurs when a thread executes a system call with a syscall number that is unknown to seL4. The register set of the faulting thread is passed to the thread's exception handler so that it may, for example, emulate the system call if a thread is being virtualised.

Replying to the fault IPC allows the thread to be restarted and/or the thread's register set to be modified. If the reply has a label of zero, the thread will be restarted. Additionally, if the message length is non-zero, the faulting thread's register set will be updated. In this case, the number of registers updated is controlled with the length field of the message tag.

6.2.3 User Exception

User exceptions are used to deliver architecture-defined exceptions. For example, such an exception could occur if a user thread attempted to divide a number by zero.

Replying to the fault IPC allows the thread to be restarted and/or the thread's register set to be modified. If the reply has a label of zero, the thread will be restarted. Additionally, if the message length is non-zero, the faulting thread's register set will be updated. In this case, the number of registers updated is controlled with the length field of the message tag.

6.2.4 Debug Exception: Breakpoints and Watchpoints

Debug exceptions are used to deliver trace and debug related events to threads. Breakpoints, watchpoints, trace-events and instruction-performance sampling events are examples. These events are supported for userspace threads when the kernel is configured to include them (when `CONFIG_HARDWARE_DEBUG_API` is set). The hardware debugging extensions API is supported on the following subset of the platforms that the kernel has been ported to:

- PC99: IA-32 and x86_64
- Sabrelite (i.MX6)
- Jetson TegraK1
- HiSilicon Hikey
- Raspberry Pi 3
- Odroid-X (Exynos4)
- Xilinx zynq7000

Information on the available hardware debugging resources is presented in the form of the following constants:

seL4_NumHWBreakpoints : Defines the total number of hardware break registers available, of all types available on the hardware platform. On the ARM Cortex A7 for example, there are 6 exclusive instruction breakpoint registers, and 4 exclusive data watchpoint registers, for a total of 10 monitor registers. On this platform therefore, **seL4_NumHWBreakpoints** is defined as 10. The instruction breakpoint registers will always be assigned the lower API-IDs, and the data watchpoints will always be assigned following them.

Additionally, **seL4_NumExclusiveBreakpoints**, **seL4_NumExclusiveWatchpoints** and **seL4_NumDualFunctionMonitors** are defined for each target platform to reflect the number of available hardware breakpoints/watchpoints of a certain type.

seL4_NumExclusiveBreakpoints : Defines the number of hardware registers capable of generating a fault **only** on instruction execution. Currently this will be set only on ARM platforms. The API-ID of the first exclusive breakpoint is given in **seL4_FirstBreakpoint**. If there are no instruction-break exclusive registers, **seL4_NumExclusiveBreakpoints** will be set to 0 and **seL4_FirstBreakpoint** will be set to -1.

seL4_NumExclusiveWatchpoints : Defines the number of hardware registers capable of generating a fault **only** on data access. Currently this will be set only on ARM platforms. The API-ID of the first exclusive watchpoint is given in **seL4_FirstWatchpoint**. If there are no data-break exclusive registers, **seL4_NumExclusiveWatchpoints** will be set to 0 and **seL4_FirstWatchpoint** will be set to -1.

seL4_NumDualFunctionMonitors : Defines the number of hardware registers capable of generating a fault on either type of access – i.e, the register supports both instruction and data breaks. Currently this will be set only on x86 platforms. The API-ID of the first dual-function monitor is given in **seL4_FirstDualFunctionMonitor**. If there are no dual-function break registers, **seL4_NumDualFunctionMonitors** will be set to 0 and **seL4_FirstDualFunctionMonitor** will be set to -1.

Value sent	IPC buffer location
Breakpoint instruction address	IPCBuffer[0]
Exception reason	IPCBuffer[1]
Watchpoint data access address	IPCBuffer[2]
Register API-ID	IPCBuffer[3]

Table 6.2: Debug fault message layout. The register API-ID is not returned in the fault message from the kernel on single-step faults.

6.2.5 Debug Exception: Single-stepping

The kernel provides support for the use of hardware single-stepping of userspace threads when configured to do so (when **CONFIG_HARDWARE_DEBUG_API** is set). To this end it exposes the invocation, **seL4_TCB_ConfigureSingleStepping**.

The caller is expected to select an API-ID that corresponds to an instruction breakpoint, to use when setting up the single-stepping functionality (i.e, API-ID from 0 to **seL4_NumExclusiveBreakpoints** - 1). However, not all hardware platforms require an actual hardware breakpoint register to provide single-stepping functionality. If the caller’s hardware platform requires the use of a hardware breakpoint register, it will use the breakpoint register given to it in **bp_num**, and return **true** in **bp_was_consumed**. If the underlying platform does not need a hardware breakpoint to provide single-stepping, **seL4** will return **false** in **bp_was_consumed** and leave **bp_num** unchanged.

If **bp_was_consumed** is **true**, the caller should not attempt to re-configure **bp_num** for Breakpoint or Watchpoint usage until the caller has disabled single-stepping and released that register, via a subsequent call to **seL4_TCB_ConfigureSingleStepping**, or a fault-reply with **n_instr** being 0. Setting **num_instructions** to 0 **disables single stepping**.

On architectures that require an actual hardware registers to be configured for single-stepping functionality, **seL4** will restrict the number of registers that can be configured as single-steppers, to one at any given time. The register that is currently configured

(if any) for single-stepping will be the implicit `bp_num` argument in a single-step debug fault reply.

The kernel's single-stepping, also supports skipping a certain number of instructions before delivering the single-step fault message. `Num_instructions` should be set to 1 when single-stepping, or any non-zero integer value to skip that many instructions before resuming single-stepping. This skip-count can also be set in the fault-reply to a single-step debug fault.

Value sent	Register set by reply	IPC buffer location
Breakpoint instruction address	<code>num_instructions</code> skip	to <code>IPCBuffer[0]</code>
Exception reason	—	<code>IPCBuffer[1]</code>

Table 6.3: Single-step fault message layout.

6.2.6 Timeout Fault (MCS only)

Timeout faults are raised when a thread consumes all of its budget and has a timeout fault handler that is not a null capability. They allow a timeout exception handler to take some action to restore the thread, and deliver a message containing the scheduling context data word, as well as the amount of time consumed since the last timeout fault occurred on this scheduling context, or since `seL4_SchedContext_YieldTo()` or `seL4_SchedContext_Consumed()` was last called. Timeout exception handlers can reply to a temporal fault with the registers set in the same format as outlined in Section 6.1.12.

Meaning	IPC buffer location
Data word from the scheduling context object that the thread was running on when the fault occurred.	<code>seL4_TimeoutFault_Data</code>
Upper 32-bits of microseconds consumed since last reset	<code>seL4_TimeoutFault_Consumed</code>
Lower 32-bits of microseconds consumed since last reset	<code>seL4_TimeoutFault_Consumed_Low-Bits</code>

Table 6.4: Timeout fault outcome on 32-bit architectures.

6.2.7 VM Fault

The thread caused a page fault. Replying to the fault IPC will restart the thread. The contents of the IPC message are given below.

Meaning	IPC buffer location
Program counter to restart execution at.	<code>seL4_VMFault_IP</code>
Address that caused the fault.	<code>seL4_VMFault_SP</code>
Instruction fault (1 if the fault was caused by an instruction fetch).	<code>seL4_VMFault_PrefetchFault</code>
Fault status register (FSR). Contains information about the cause of the fault. Architecture dependent.	<code>seL4_VMFault_FSR</code>

Table 6.5: VM Fault outcome on all architectures.

6.3 Domains

Domains are used to isolate independent subsystems, so as to limit information flow between them. The kernel switches between domains according to a fixed, time-triggered schedule. The fixed schedule is compiled into the kernel via the constant `CONFIG_NUM_DOMAINS` and the global variable `ksDomSchedule`.

A thread belongs to exactly one domain, and will only run when that domain is active. The `seL4_DomainSet_Set()` method changes the domain of a thread. The caller must possess a Domain cap and the thread's TCB cap. The initial thread starts with a Domain cap (see Section 4.1).

6.4 Virtualisation

Hardware execution virtualisation is supported on specific arm and x86 platforms. The interface is exposed through a series of kernel objects, invocations and syscalls that allow the user to take advantage of hardware virtualisation features.

Hardware virtualisation allows for a thread to perform instructions and operations as if it were running at a higher privilege level. As higher privilege levels typically have access to additional machine registers and other pieces of state a VCPU object is introduced to act as storage for this state. For simplicity we refer to this virtualised higher privileged level as 'guest mode'. VCPUs are bound in a one-to-one relationship with a TCB in order to provide a thread with this ability to run in higher privilege mode. See the section on ARM or x86 for more precise details.

VCPU objects also have additional, architecture specific, invocations for manipulating the additional state or other virtualisation controls provided by the hardware. Binding of a VCPU to a TCB is done by an invocation on the VCPU only, and not the TCB.

The provided objects and invocations are, generally speaking, the thinnest possible shim over the underlying hardware primitives and operations. As a result an in depth familiarity with the underlying architecture specific hardware mechanisms is required to use these objects, and such familiarity is therefore assumed in description.

6.4.1 ARM

When a TCB has a bound VCPU it is allowed to have the mode portion of the `cpsr` register set to values other than `user`. Specifically it may have any value other than `hypervisor`.

TODO: this section needs more detail

6.4.2 x86

A TCB with a bound VCPU has two execution modes; one is the original thread just as if there was no bound VCPU, and the other is the guest mode execution using the VCPU. Switching from regular execution mode into the guest execution mode is done by using the `seL4_VMEnter()` syscall. Executing this syscall causes the thread, whenever it is scheduled thereafter, to execute using the higher privileged mode controlled by the VCPU. Should the guest execution mode generate any kind of fault, or if a message arrives on the TCBs bound notification, the TCB will be switched back to regular mode and the `seL4_VMEnter()` syscall will return with a message indicating the reason for return.

VCPU state and execution is controlled through the `seL4_VCPU_ReadVMCS()` and `seL4_VCPU_WriteVMCS()` invocations. These are very thin wrappers around the hardware `vmread` and `vmwrite` instructions and the kernel merely does enough validation on the parameters to ensure the VCPU is not configured to run in such a way as to violate any kernel properties. For example, it is not possible to disable the use of External Interrupt Exiting, as this would prevent the kernel from receiving timer interrupts and allow the thread to monopolise CPU time.

Memory access of the guest execution mode is controlled by requiring the use of Extended Page Tables (EPT). A series of EPT related paging structure objects (`EPTPML4`, `EPTPDPT`, `EPTPD`, `EPTPT`) exist and are manipulated in exactly the same manner as the objects for the regular virtual address space. Once constructed a TCB can be given an `EPTPML4` as an EPT root with `seL4_TCB_SetEPTRoot()`, which serves as the vspace root when executing in guest mode, with the vspace root set with `seL4_TCB_SetSpace()` or `seL4_TCB_Configure()` continuing to provide translation when the TCB is executing in its normal mode.

Direct access to I/O ports can be given to the privileged execution mode through the `seL4_X86_VCPU_EnableIOPort()` invocation and allows the provided I/O port capability to be linked to the VCPU, and a subset of its I/O port range to be made accessible to the VCPU. Linking means that an I/O port capability can only be used in a single `seL4_X86_VCPU_EnableIOPort()` invocation and a second invocation will undo the previous one. The link also means that if the I/O port capability is deleted for any reason the access will be correspondingly removed from the VCPU.

Chapter 7

Address Spaces and Virtual Memory

A virtual address space in seL4 is called a VSpace. In a similar way to a CSpace (see Chapter 3), a VSpace is composed of objects provided by the microkernel. Unlike CSpaces, these objects for managing virtual memory largely correspond to those of the hardware. Consequently, each architecture defines its own objects for the top-level VSpace and further intermediate paging structures. Common to every architecture is the Page, representing a frame of physical memory. The kernel also includes ASID Pool and ASID Control objects for tracking the status of address spaces.

These VSpace-related objects are sufficient to implement the hardware data structures required to create, manipulate, and destroy virtual memory address spaces. It should be noted that, as usual, the manipulator of a virtual memory space needs the appropriate capabilities to the required objects.

7.1 Objects

7.1.1 Hardware Virtual Memory Objects

Each architecture has a top-level paging structure (level 0) and a number of intermediate levels. The top-level paging structure corresponds directly to the higher-level concept of a VSpace in seL4. For each architecture, the VSpace is realised as a different object, as determined by the architectural details.

In general, each paging structure at each level contains slots where the next level paging structure, or a specifically sized frame of memory, can be mapped. If the previous level is not mapped, a mapping operation will fail. Developers need to manually create and map all paging structures. The size and type of structure at each level, and the number of bits in the virtual address resolved for that level, is hardware defined.

seL4 provides methods for operating on these hardware paging structures including mapping and cache operations. Mapping operations are invoked on the capability being mapped, e.g. to map a level 1 paging structure at a specific virtual address, the capability to the corresponding object is invoked with a map operation, where the

top-level structure is passed as an argument.

In general, the top-level structure has no invocations for mapping, but is used as an argument to several other virtual-memory related object invocations. For some architectures, the top-level page table can be invoked for cache operations. By making these cache related operations invocations on page directory capabilities in addition to the page capabilities themselves, the API allows users more flexible policy options. For example, a process that has delegated a page directory can conduct cache operations on all frames mapped from that capability without access to those capabilities directly.

The rest of this section details the paging structures for each architecture.

7.1.1.1 IA-32

On IA-32, the VSpace is realised as a `PageDirectory`, which covers the entire 4 GiB range in the 32-bit address space, and forms the top-level paging structure. Second level page-tables (`PageTable` objects) each cover a 4 MiB range. Structures at both levels are indexed by 10 bits in the virtual address.

<i>Object</i>	<i>Address Bits</i>	<i>Level</i>	<i>Methods</i>
<code>PageDirectory</code>	22—31	0	Section 10.4.12
<code>PageTable</code>	12—21	1	Section 10.4.13

7.1.1.2 x64

On x86-64, the VSpace is realised as a PML4. Three further levels of paging structure are defined, as shown in the table below. All structures are indexed with 9 bits of the virtual address.

<i>Object</i>	<i>Address Bits</i>	<i>Level</i>	<i>Methods</i>
PML4	39—47	0	None
PDPT	30—38	1	Section 10.6.1
<code>PageDirectory</code>	21—29	2	Section 10.4.12
<code>PageTable</code>	12—20	3	Section 10.4.13

7.1.1.3 AArch32

Like IA-32, ARM AArch32 realise the VSpace as a `PageDirectory`, which covers the entire 4 GiB address range, and a second-level `PageTable`. The second-level structures on AArch32 cover 1 MiB address ranges.

ARM AArch32 processors have a two-level page-table structure. The top-level page directory covers a range of 4 GiB and each page table covers a 1 MiB range.

<i>Object</i>	<i>Address Bits</i>	<i>Level</i>	<i>Methods</i>
<code>PageDirectory</code>	20—31	0	Section 10.8.1
<code>PageTable</code>	12—19	1	Section 10.7.7

7.1.1.4 AArch64

ARM AArch64 processors have a four-level page-table structure, where the VSpace is realised as a `PageGlobalDirectory`. All paging structures are indexed by 9 bits of the virtual address.

<i>Object</i>	<i>Address Bits</i>	<i>Level</i>	<i>Methods</i>
<code>PageGlobalDirectory</code>	39—47	0	??
<code>PageUpperDirectory</code>	30—38	1	Section 10.9.2
<code>PageDirectory</code>	21—29	2	Section 10.9.1
<code>PageTable</code>	12—20	3	Section 10.7.7

7.1.2 RISC-V

RISC-V provides the same paging structure for all levels, `PageTable`. The VSpace is then realised as a `PageTable`.

7.1.2.1 RISC-V 32-bit

32-bit RISC-V `PageTables` are indexed by 10 bits of virtual address.

<i>Object</i>	<i>Address Bits</i>	<i>Level</i>	<i>Methods</i>
<code>PageTable</code>	22—31	0	Section 10.10.6
<code>PageTable</code>	12—21	1	Section 10.10.6

7.1.2.2 RISC-V 64-bit

64-bit RISC-V follows the SV39 model, where `PageTables` are indexed by 9 bits of virtual address. Although RISC-V allows for multiple different numbers of paging levels, currently seL4 only supports exactly three levels of paging structures.

<i>Object</i>	<i>Address Bits</i>	<i>Level</i>	<i>Methods</i>
<code>PageTable</code>	30—38	0	Section 10.10.6
<code>PageTable</code>	21—29	1	Section 10.10.6
<code>PageTable</code>	12—20	2	Section 10.10.6

7.1.3 Page

A `Page` object corresponds to a frame of physical memory that is used to implement virtual memory pages in a virtual address space.

The virtual address for a `Page` mapping must be aligned to the size of the `Page` and must be mapped to a suitable VSpace, and every intermediate paging structure required. To map a page readable, the capability to the page that is being invoked must have read permissions. To map the page writeable, the capability must have write permissions.

The requested mapping permissions are specified with an argument of type `seL4_CapRights` given to the mapping function. If the capability does not have sufficient permissions to authorise the given mapping, then the mapping permissions are silently downgraded. Specific mapping permissions are dependant on the architecture and are documented in the Chapter 10 for each function.

At minimum, each architecture defines `Map`, `Unmap` and `GetAddress` methods for pages. Methods for page objects for each architecture can be found in the Chapter 10, and are indexed per architecture in the table below.

<i>Architectures</i>	<i>Methods</i>
IA32, X64	Section 10.4.11
AArch32, AArch64	Section 10.7.6
RISC-V	Section 10.10.5

Each architecture also defines a range of page sizes. In the next section we show the available page sizes, as well as the *mapping level*, which refers to the level of the paging structure at which this page must be mapped.

7.1.3.1 AArch32 page sizes

<i>Constant</i>	<i>Size</i>	<i>Mapping level</i>
<code>seL4_PageBits</code>	4 KiB	1
<code>seL4_LargePageBits</code>	64 KiB	1
<code>seL4_SectionBits</code>	1 MiB	0
<code>seL4_SuperSectionBits</code>	16 MiB	0

Mappings for sections and super sections consume 16 slots in the page table and page directory respectively.

7.1.3.2 AArch64 page sizes

<i>Constant</i>	<i>Size</i>	<i>Mapping level</i>
<code>seL4_PageBits</code>	4 KiB	3
<code>seL4_LargePageBits</code>	2 MiB	2
<code>seL4_HugePageBits</code>	1 GiB	1

7.1.3.3 IA-32 page sizes

<i>Constant</i>	<i>Size</i>	<i>Mapping level</i>
<code>seL4_PageBits</code>	4 KiB	1
<code>seL4_LargePageBits</code>	4 MiB	0

7.1.3.4 X64 page sizes

<i>Constant</i>	<i>Size</i>	<i>Mapping level</i>
seL4_PageBits	4 KiB	3
seL4_LargePageBits	2 MiB	2
seL4_HugePageBits	1 GiB	1

7.1.3.5 RISC-V 32-bit page sizes

<i>Constant</i>	<i>Size</i>	<i>Mapping level</i>
seL4_PageBits	4 KiB	1
seL4_LargePageBits	4 MiB	0

7.1.3.6 RISC-V 64-bit page sizes

<i>Constant</i>	<i>Size</i>	<i>Mapping level</i>
seL4_PageBits	4 KiB	2
seL4_LargePageBits	2 MiB	1
seL4_HugePageBits	1 GiB	0

7.1.4 ASID Control

For internal kernel book-keeping purposes, there is a fixed maximum number of applications the system can support. In order to manage this limited resource, the microkernel provides an ASID Control capability. The ASID Control capability is used to generate a capability that authorises the use of a subset of available address-space identifiers. This newly created capability is called an ASID Pool. ASID Control only has a single `MakePool` method for each architecture, listed in the table below.

<i>Architectures</i>	<i>Methods</i>
IA32, X64	Section 10.4.3
AArch32, AArch64	Section 10.7.1
RISC-V	Section 10.10.3

7.1.5 ASID Pool

An ASID Pool confers the right to create a subset of the available maximum applications. For a VSpace to be usable by an application, it must be assigned to an ASID. This is done using a capability to an ASID Pool. The ASID Pool object has a single method, `Assign`, for each architecture:

<i>Architectures</i>	<i>Methods</i>
IA32, X64	Section 10.4.4
AArch32, AArch64	Section 10.7.2
RISC-V	Section 10.10.4

7.2 Mapping Attributes

A parameter of type `seL4_ARM_VMAAttributes` or `seL4_x86_VMAAttributes` is used to specify the cache behaviour of the page being mapped; possible values for ARM that can be bitwise OR'd together are shown in Table 7.1 and an enumeration of valid values for IA-32 are shown in Table 7.2. Mapping attributes can be updated on existing mappings using the `Map` invocation with the same virtual address.

Attribute	Meaning
<code>seL4_ARM_PageCacheable</code>	Enable data in this mapping to be cached
<code>seL4_ARM_ParityEnabled</code>	Enable parity checking for this mapping
<code>seL4_ARM_ExecuteNever</code>	Map this memory as non-executable

Table 7.1: Virtual memory attributes for ARM page table entries.

Attribute	Meaning
<code>seL4_x86_WriteBack</code>	Read and writes are cached
<code>seL4_x86_CacheDisabled</code>	Prevent data in this mapping from being cached
<code>seL4_x86_WriteThrough</code>	Enable write through cacheing for this mapping
<code>seL4_x86_WriteCombining</code>	Enable write combining for this mapping

Table 7.2: Virtual memory attributes for x86 page table entries.

7.3 Sharing Memory

seL4 does not allow Page Tables to be shared, but does allow pages to be shared between address spaces. To share a page, the capability to the Page must first be duplicated using the `seL4_CNode_Copy()` method and the new copy must be used in the `seL4_ARM_Page_Map()` or `seL4_x86_Page_Map()` method that maps the page into the second address space. Attempting to map the same capability twice will result in an error.

7.4 Page Faults

Page faults are reported to the exception handler of the executed thread. See Section 6.2.7.

Chapter 8

Hardware I/O

8.1 Interrupt Delivery

Interrupts are delivered as notifications. A thread may configure the kernel to signal a particular Notification object each time a certain interrupt triggers. Threads may then wait for interrupts to occur by calling `seL4_Wait()` or `seL4_Poll()` on that Notification.

IRQHandler capabilities represent the ability of a thread to configure a certain interrupt. They have three methods:

`seL4_IRQHandler_SetNotification()` specifies the Notification the kernel should `signal()` when an interrupt occurs. A driver may then call `seL4_Wait()` or `seL4_Poll()` on this notification to wait for interrupts to arrive.

`seL4_IRQHandler_Ack()` informs the kernel that the userspace driver has finished processing the interrupt and the microkernel can send further pending or new interrupts to the application.

`seL4_IRQHandler_Clear()` de-registers the Notification from the IRQHandler object.

When the system first starts, no IRQHandler capabilities are present. Instead, the initial thread's CSpace contains a single IRQControl capability. This capability may be used to produce a single IRQHandler capability for each interrupt available in the system. Typically, the initial thread of a system will determine which IRQs are required by other components in the system, produce an IRQHandler capability for each interrupt, and then delegate the resulting capabilities as appropriate. Methods on IRQControl can be used for creating IRQHandler capabilities for interrupt sources.

8.2 x86-Specific I/O

8.2.1 Interrupts

In addition to managing IRQHandler capabilities, x86 platforms require the delivery location in the CPU vectors to be configured. Regardless of where an interrupt comes

from (IOAPIC, MSI, etc) it must be assigned a unique vector for delivery, ranging from `VECTOR_MIN` to `VECTOR_MAX`. The rights to allocate a vector are effectively given through the `IRQControl` capability and can be considered as the kernel outsourcing the allocation of this namespace to user level.

`seL4_IRQControl_GetIOAPIC()` creates an `IRQHandler` capability for an IOAPIC interrupt

`seL4_IRQControl_GetMSI()` creates an `IRQHandler` capability for an MSI interrupt

8.2.2 I/O Ports

On x86 platforms, seL4 provides access to I/O ports to user-level threads. Access to I/O ports is controlled by IO Port capabilities. Each IO Port capability identifies a range of ports that can be accessed with it. Reading from I/O ports is accomplished with the `seL4_X86_IOPort_In8()`, `seL4_X86_IOPort_In16()`, and `seL4_X86_IOPort_In32()` methods, which allow for reading of 8-, 16- and 32-bit quantities. Similarly, writing to I/O ports is accomplished with the `seL4_X86_IOPort_Out8()`, `seL4_X86_IOPort_Out16()`, and `seL4_X86_IOPort_Out32()` methods. Each of these methods takes as arguments an IO Port capability and an unsigned integer `port`, which indicates the I/O port to read from or write to, respectively. In each case, `port` must be within the range of I/O ports identified by the given IO Port capability in order for the method to succeed.

The I/O port methods return error codes upon failure. A `seL4_IllegalOperation` code is returned if port access is attempted outside the range allowed by the IO Port capability. Since invocations that read from I/O ports are required to return two values – the value read and the error code – a structure containing two members, `result` and `error`, is returned from these API calls.

At system initialisation, the initial thread's CSpace contains the `IOPortControl` capability, which can be used to `seL4_X86_IOPort_Issue()` IO Port capabilities to sub ranges of I/O ports. Any range that is issued may not have overlap with any existing issued IO Port capability.

8.2.3 I/O Space

I/O devices capable of DMA present a security risk because the CPU's MMU is bypassed when the device accesses memory. In seL4, device drivers run in user space to keep them out of the trusted computing base. A malicious or buggy device driver may, however, program the device to access or corrupt memory that is not part of its address space, thus subverting security. To mitigate this threat, seL4 provides support for the IOMMU on Intel x86-based platforms. An IOMMU allows memory to be remapped from the device's point of view. It acts as an MMU for the device, restricting the regions of system memory that it can access. More information can be obtained from Intel's IOMMU documentation [Int11].

Two new objects are provided by the kernel to abstract the IOMMU:

IOspace This object represents the address space associated with a hardware device on the PCI bus. It represents the right to modify a device’s memory mappings.

IOPageTable This object represents a node in the multilevel page-table structure used by IOMMU hardware to translate hardware memory accesses.

Page capabilities are used to represent the actual frames that are mapped into the I/O address space. A Page can be mapped into either a VSpace or an IOspace but never into both at the same time.

IOspace and VSpace fault handling differ significantly. VSpace page faults are redirected to the thread’s exception handler (see Section 6.2), which can take the appropriate action and restart the thread at the faulting instruction. There is no concept of an exception handler for an IOspace. Instead, faulting transactions are simply aborted; the device driver must correct the cause of the fault and retry the DMA transaction.

An initial master IOspace capability is provided in the initial thread’s CSpace. An IOspace capability for a specific device is created by using the `seL4_CNode_Mint()` method, passing the PCI identifier of the device as the low 16 bits of the `badge` argument, and a Domain ID as the high 16 bits of the `badge` argument. PCI identifiers are explained fully in the PCI specification [SA99], but are briefly described here. A PCI identifier is a 16-bit quantity. The first 8 bits identify the bus that the device is on. The next 5 bits are the device identifier: the number of the device on the bus. The last 3 bits are the function number. A single device may consist of several independent functions, each of which may be addressed by the PCI identifier. Domain IDs are explained fully in the Intel IOMMU documentation [Int11]. There is presently no way to query seL4 for how many Domain IDs are supported by the IOMMU and the `seL4_CNode_Mint()` method will fail if an unsupported value is chosen.

The IOMMU page-table structure has three levels. Page tables are mapped into an IOspace using the `seL4_X86_IOPageTable_Map()` method. This method takes the IOPageTable to map, the IOspace to map into and the address to map at. Three levels of page tables must be mapped before a frame can be mapped successfully. A frame is mapped with the `seL4_X86_Page_MapIO()` method whose parameters are analogous to the corresponding method that maps Pages into VSpaces (see Chapter 7), namely `seL4_X86_Page_Map()`.

Unmapping is accomplished with the usual `unmap` (see Chapter 7) API call, `seL4_X86_Page_Unmap()`.

More information about seL4’s IOMMU abstractions can be found in [Pal09].

8.3 Arm-Specific I/O

8.3.1 Arm SMMU version 2.0

seL4 provides an API for programming the Arm System MMU (SMMU) version 2.0, which allows system software to manage access rights and address translation for devices that can initiate direct memory accesses (DMA).

An Arm SMMU v2.0 implementation allows device memory transactions to be associated with an identifier (StreamID) that is used to direct the transaction through a SMMU translation context bank (CB). A translation context bank can perform address translation, memory protection and memory attribute transformation. The standard specifies different types of address translations that correspond to stages in the ArmV8 virtual memory system architecture such as either non-secure EL0, EL1 first and second stage translations, Hyp mode translations or secure mode translations. It is possible to associate different StreamIDs with the same context bank and it is possible to share address translation tables between a context bank and software MMU address space if the stage and type of translation is the same.

Faults that occur when a memory transaction conflicts with a StreamID or CB configuration happen asynchronously with respect to a processor element's execution. When this occurs an interrupt is used to allow a PE to handle the SMMU fault. Faults are reported through registers in the SMMU that can be queried in an interrupt handler.

TLB maintenance operations are required to keep SMMU translation caches consistent when there are changes to any valid page table mapping entries.

An SMMU implementation usually has a maximum number of StreamIDs that it supports. The specification allows StreamIDs to be up to 16bits wide. There are also a fixed number of context banks, up to a maximum of 128. Context banks can be generic or support only a single address translation stage. This information is reported by ID registers in each implementation.

The seL4 API allows system software to manage an SMMU by assigning StreamIDs to context banks, bind context banks to page translation structures, implement SMMU fault handling and also perform explicit TLB maintenance. This allows system software to ensure that a device is only able to access and modify memory contents that it has been explicitly given access to and allow devices to be presented with a virtualized address space for performing DMA.

All the StreamIDs and context banks are accessible via capabilities. Control capabilities are used to create capabilities referring to each StreamID and context bank in a system. The kernel tracks the allocation of StreamIDs and context banks with two static CNodes, one for each resource type. These CNodes track which VSpace a context bank has bound to it, and which context bank a StreamID is bound to.

The capabilities allow access control policies to be implemented by a user thread. When StreamID, context bank capabilities are revoked, the kernel will disable the context banks or StreamID mappings.

TLB maintenance is handled by the kernel via tracking which context banks are associated with a particular VSpace. Any TLB maintenance operations that the kernel performs on VSpace invocations are also applied to associated context banks.

SMMU fault handling is delegated to user level via invocations that allow fault statuses to be queried and cleared for each context bank and for the SMMU globally. SMMU fault interrupts can be handled the same as other platform level interrupts.

The kernel implementation only uses translation stages matching what translation the kernel is performing for VSpace objects. When seL4 is operating in EL1, the SMMU only uses stage 1 translation (ASID), that is "stage 1 with stage 2 bypass" in the

context bank attribute configuration. When hypervisor mode is enabled, and seL4 is operating in EL2, the SMMU only does stage 2 translations.

Four capabilities types provide access to SMMU resources:

seL4_ARM_SID A capability granting access to a single transaction stream, which can be used to bind and unbind a stream to a single context bank.

seL4_ARM_CB A capability representing a single specific context bank. It can be used to bind and unbind a VSpace to assign what page tables the context bank should use for translation, assign StreamIDs and process context bank faults.

seL4_ARM_SIDControl A control capability which can be used to create seL4_ARM_SID capabilities to specific transaction streams. The seL4_ARM_SIDControl cap is used for managing rights on StreamID configurations. This capability is provided in the initial thread's CSpace.

seL4_ARM_CBControl A control capability that can be used to derive seL4_ARM_CB capabilities. The seL4_ARM_CBControl cap is used for managing rights on context bank configurations. This capability is provided in the initial thread's CSpace.

8.3.1.1 Creating seL4_ARM_SID capabilities

The Arm SMMU 2.0 specification doesn't specify how StreamIDs need to correspond to different devices. Each platform can define its own policy for how StreamIDs are allocated. A seL4_ARM_SIDControl capability can be used to create a capability to any valid StreamID for the SMMU and delegate access to other tasks in the system.

seL4_ARM_SIDControl_GetSID() uses the seL4_ARM_SIDControl capability to create a new seL4_ARM_SID capability that represents a single StreamID. This new capability is placed in the provided slot. It is expected that whatever thread controls an seL4_ARM_SIDControl capability knows about how StreamIDs are allocated in a system.

The Arm SMMU 2.0 specification describes many ways of associating StreamIDs with context banks. Currently only direct mapping of a StreamID to a context bank is supported.

8.3.1.2 Creating seL4_ARM_CB capabilities

Each context bank allows the SMMU to maintain an active translation context with its own registers for holding context specific information. An SMMU has a fixed number of context banks available for use and these are allocated using the seL4_ARM_CBControl capability.

seL4_ARM_CBControl_GetCB() uses the seL4_ARM_CBControl capability to create a new seL4_ARM_CB capability that represents a single context bank. This new capability is placed in the provided slot. It is expected that whatever thread controls a seL4_ARM_CBControl capability has knowledge of the properties of each context bank that each index refers to.

8.3.1.3 Configuring context banks

By providing a `seL4_ARM_CB` cap, a user-level thread can configure the VSpace used by the bank with the following API:

`seL4_ARM_CB_AssignVspace()` configures the context bank to use the provided VSpace root for translations.

`seL4_ARM_CB_UnassignVspace()` removes the configured VSpace and conducting a TLB invalidation.

The SMMU-v2 uses the same paging structure as the MMU (AArch_64 and AArch_32 formats). Therefore, there is no need to provide a new set of page structure caps nor a separate set of map and unmap functions. To manage the assignment, the kernel has an internal CNode, called `smmuStateCBNode`, that stores copies of the `VSpace_cap` created by executing the above API. The copy of the `VSpace_cap` contains its assigned ContextBank number. Therefore the kernel can conduct context bank invalidation if the `VSpace_cap` is revoked.

8.3.1.4 Configuring streams (transactions)

A user-level thread can bind a context bank with an `seL4_ARM_SID` capability with:

`seL4_ARM_SID_BindCB()` configures the stream to use given context bank for translation. To simplify the process, the binding also enables the stream ID. `seL4_ARM_SID_BindCB` generates a copy of the `seL4_ARM_CB` cap in kernel's internal CNode. This allows the stream ID to be disabled if the `seL4_ARM_CB` cap is revoked.

`seL4_ARM_SID_UnbindCB()` removes the `seL4_ARM_CB` cap from the kernel's internal CNode and disables the stream ID. The kernel provides this API for the conveniences of sharing a stream ID among multiple VSpaces.

If there are any exceptions after the stream ID is enabled, the user-level software should use the fault handling mechanisms to resolve them.

8.3.1.5 Copying and Deleting caps

The kernel allows copying both `ARM_SID` cap and `seL4_ARM_CB` cap. This allows capabilities to be delegated to different threads. The kernel does not allow copying neither the `seL4_ARM_SIDControl` nor the `seL4_ARM_CBControl` capabilities.

Deleting a `seL4_ARM_CB` cap that contains a valid `capBindSID` field will:

- invalidate the `streamID` to `ContextBank` assignment in hardware.

Deleting the last `seL4_ARM_CB` cap will:

- perform an `seL4_ARM_CB_UnassignVspace()`, removing any configured VSpace,
- conduct a TLB invalidation.

Similarly, deleting a `VSpace_cap` that contains assigned context bank number will:

- invalidate the context bank
- conduct a TLB invalidation

Deleting the last `ARM_SID` cap will:

- Perform an `seL4_ARM_SID_UnbindCB()`, (deleting the copy of the assigned `seL4_ARM_CB` cap)
- Disable the stream ID.

8.3.1.6 TLB invalidation

The kernel is expected to perform all required SMMU TLB maintenance operations as part of the API implementation. In addition, the kernel provides two system calls for explicitly performing invalidations:

`seL4_ARM_CBControl_TLBInvalidateAll()` invalidates all TLB entries in all context banks.

`seL4_ARM_CB_TLBInvalidate()` invalidates all TLB entries in a context bank.

The kernel does not impose any restrictions on how a VSpace is used by user-level applications, hence a VSpace can be shared by normal threads and drivers. Sharing a VSpace between threads and drivers also means sharing all mappings in that VSpace between MMUs in CPU cores and SMMU used by device transactions. Moreover, multiple context banks in SMMU can share a VSpace. Therefore, maintaining the coherency between the TLB in MMU and the TLB in SMMU's context banks is important.

The kernel keeps a record of Vspace's usage in context banks in SMMU by maintaining: the number of context banks using a given ASID, and the ASID that a given context bank is using. There are a few reasons behind this design.

- First, the ASID is efficient for representing a VSpace. In seL4, each VSpace has an ASID which is assigned before the VSpace is ready to be used and will never change until the VSpace is deleted. Recording how many context banks are using a VSpace's ASID is equivalent to recording the VSpace's usage in context banks.
- Second, all TLB invalidation operation requires knowledge of the ASID. There are two types of TLB invalidation operations: invalidating a page table entry using its ASID (triggered by updating a page table entry, e.g. unmapping a page), and invalidating all mappings of an ASID (triggered by deleting a VSpace).

- Third, the kernel can easily find a context banks' ASID on all occasions, which is useful to either conduct TLB invalidation requests or unassign VSpace from a context bank.

By knowing how many context banks are using an ASID, the kernel can easily check in every TLB invalidation operation and invoke TLB invalidation in SMMU if the value is not zero. In SMMU's TLB invalidation operation, the kernel searches the context banks using the ASID, and conducts TLB invalidation in those context banks.

Ideally, the SMMU shares the same ASID or VMID name space with the rest of the system. This allows the SMMU to maintain TLB coherency by listening for TLB broadcasting messages. This means the context banks should be configured with the correct ASID or VMID when the StreamID is enabled. This is not a problem for stage 1 translation, as there are a large number of ASID bits and an ASID can be assigned to a vspace root with existing APIs. However, the VMID used in stage 2 only has 8 bits, and the kernel allocates them on demand and can reclaim a VSpace's hardware ASID to reuse if there are more VSpaces than available ASIDS. While it is possible to do this when the VSpace is only used in an MMU, it is not possible with multiple active context banks. Due to this, the context bank in SMMU cannot be configured with the correct VMID. Currently, the SMMU driver uses private VMID space, and uses the context bank number as the corresponding VMID number.

8.3.1.7 Fault handling

The number of IRQs used for reporting transaction faults is hardware dependent. There are two kinds of faults: global faults (general configuration and transaction faults), or context bank faults. For transaction faults, the SMMU reports faulty stream IDs. The global faults reports:

- Invalid context fault.
- Unidentified stream fault.
- Stream match conflict fault.
- Unimplemented context bank fault.
- Unimplemented context interrupt fault.
- Configuration access fault.
- External fault.

Each context bank contains registers to report faults on address translation, for example, faulty addresses, or permission errors. The SMMU driver identifies the cause of a fault by first reading the global fault registers (one state register and three fault syndrome registers), then by reading corresponding context bank fault registers. Note, the SMMU reports the faulty transaction (stream) ID, which can be used to identify its context bank ID.

- System assumption: Both the SMMU's IRQ handler and the owner of the `seL4_ARM_SIDControl` cap (controlling stream ID distributions) are trusted.
- SMMU interrupts are handled as same as other IRQs, i.e, the kernel does not treat the SMMU IRQs special, reporting the interrupt via IRQ notifications.
- The kernel provides a API for reading the global fault registers: `seL4_ARM_SIDControl_GetFault()`. Because the IRQ notification can only deliver information via the badge, the owner of the `StreamControl_cap` can retrieve more information via this API.
- If the fault is related to a transaction, the owner of the `seL4_ARM_SIDControl` cap will notify the holder of the corresponding stream ID cap, which should also have a copy of the context bank cap bound to this transaction.
- The kernel provides an API for reading the context bank fault registers: `seL4_ARM_CB_CBGetFault()`, used by a context bank cap holder (the `seL4_ARM_CB` cap holder).
- Once the fault handling is done, the server can call `seL4_ARM_CB_CBClearFault()` to clear the fault status on a context bank, and `seL4_ARM_SIDControl_ClearFault()` to clear the fault status on SMMU.

Chapter 9

System Bootstrapping

9.1 Initial Thread's Environment

The seL4 kernel creates a minimal boot environment for the initial thread, which is started at priority `seL4_MaxPrio` and maximum control priority `seL4_MaxPrio`. This environment consists of the initial thread's TCB, CSpace and VSpace, consisting of frames that contain the userland image (code/data of the initial thread) and the IPC buffer.

On the MCS kernel, the initial thread is configured with a round-robin scheduling context with `CONFIG_BOOT_THREAD_TIME_SLICE` milliseconds timeslice. Without MCS, all threads including the initial thread are scheduled round-robin with `CONFIG_TIMER_TICK_MS * CONFIG_TIME_SLICE` timeslices.

The initial thread's CSpace consists of exactly one CNode which contains capabilities to the initial thread's own resources as well as to all available global resources. The CNode size can be configured at compile time (default is 2^{12} slots), but the guard is always chosen so that the CNode resolves exactly the number of bits in the architecture (32 bits or 64 bits). This means, the first slot of the CNode has CPTR 0x0, the second slot has CPTR 0x1 etc.

The first 15 slots (or 14 slots if not MCS) contain specific capabilities as listed in Table 9.1.

9.2 BootInfo Frame

CNode slots with CPTR `seL4_NumInitialCaps` (defined in the seL4 userland library) and above are filled dynamically during bootstrapping. Their exact contents depend on the userland image size, platform configuration (devices) etc. In order to tell the initial thread which capabilities are stored where in its CNode, the kernel provides a *BootInfo Frame* which is mapped into the initial thread's address space. The mapped address is chosen by the kernel and given to the initial thread via a CPU register.

The BootInfo Frame contains the C struct described in Table 9.2. It is defined in the seL4 userland library. Besides talking about capabilities, it also informs the initial

Table 9.1: Initial thread's CNode content.

Enum Constant	Capability
<code>seL4_CapNull</code>	null
<code>seL4_CapInitThreadTCB</code>	initial thread's TCB
<code>seL4_CapInitThreadCNode</code>	initial thread's CNode
<code>seL4_CapInitThreadVSpace</code>	initial thread's VSpace
<code>seL4_CapIRQControl</code>	global IRQ controller (see Section 8.1)
<code>seL4_CapASIDControl</code>	global ASID controller (see Chapter 7)
<code>seL4_CapInitThreadASIDPool</code>	initial thread's ASID pool (see Chapter 7)
<code>seL4_CapIOPort</code>	global I/O port cap, null cap if unsupported (see Section 8.2.2)
<code>seL4_CapIOSpace</code>	global I/O space cap, null cap if unsupported (see Section 8.2.3)
<code>seL4_CapBootInfoFrame</code>	BootInfo frame (see Section 9.2)
<code>seL4_CapInitThreadIPCBuffer</code>	initial thread's IPC buffer (see Section 4.1)
<code>seL4_CapDomain</code>	domain cap (see Section 6.3)
<code>seL4_CapSMMUSIDControl</code>	global ARM SMMU SID controller, null cap if unsupported (see Section 8.3.1)
<code>seL4_CapSMMUCBControl</code>	global ARM SMMU CB controller, null cap if unsupported (see Section 8.3.1)
<code>seL4_CapInitThreadSC</code>	initial thread's scheduling context (MCS only)

thread about the current platform's configuration.

The type `seL4_SlotRegion` is a C struct which contains `start` and `end` slot CPTRs. It denotes a region of slots in the initial thread's CNode, starting with CPTR `start` and with `end` being the CPTR of the first slot after the region ends, i.e. `end - 1` points to the last slot of the region.

Depending on the architecture and platform there might be additional pieces of boot information. If `extraLen` is greater than zero then 4K after the start of bootinfo is a region of `extraLen` bytes containing additional bootinfo structures. Each chunk starts with a `seL4_BootInfoHeader`, described in Table 9.3, that describes what the chunk is and how long it is, where the length includes the header. The length can be used to skip over chunks that you do not understand. The only generally defined chunk type is `SEL4_BOOTINFO_HEADER_PADDING` and describes an empty chunk that has no data, any other types are platform or architecture specific. The `extraBIPages` slot region gives the frames capabilities for the pages that make up the additional boot info region.

The capabilities in `userImageFrames` are ordered such that the first capability references the first frame of the userland image and so on. The capabilities in `userImagePaging` are ordered in descending order of paging structure size. Within a given paging structure size, capabilities are ordered by the virtual address at which the corresponding objects are mapped into the initial thread's address space.

It is up to userland to infer the virtual address of frames referenced by the capabilities in `userImageFrames` and the virtual address and types of paging structures referenced

Table 9.2: BootInfo struct.

Field Type	Field Name	Description
seL4_Word	extraLen	length of additional bootinfo information in bytes
seL4_NodeId	nodeID	node ID
seL4_Word	numNodes	number of nodes
seL4_Word	numIOPTLevels	number of I/O page-table levels (-1 if CONFIG_IOMMU unset)
seL4_IPCBuffer*	ipcBuffer	pointer to the initial thread's IPC buffer
seL4_SlotRegion	empty	empty slots (null caps)
seL4_SlotRegion	sharedFrames	reserved
seL4_SlotRegion	userImageFrames	frames containing the userland image
seL4_SlotRegion	userImagePaging	userland-image paging structure caps
seL4_SlotRegion	ioSpaceCaps	I/O space capabilities for ARM SMMU
seL4_SlotRegion	extraBIPages	frames backing additional bootinfo information
seL4_Uint8	initThreadCNodeSizeBits	CNode size (2^n slots)
seL4_Word	initThreadDomain	domain of the initial thread (see Section 6.3)
seL4_SlotRegion	schedcontrol	seL4_SchedControl capabilities, one for each node (MCS only).
seL4_SlotRegion	untyped	untyped-memory capabilities
seL4_UntypedDesc []	untypedList	array of information about each untyped

by the capabilities in `userImagePaging`. Userland typically has a way of finding out to which virtual addresses its code and data is mapped (e.g. in GCC, with the standard linker script, the symbols `__executable_start` and `_end` are available). Additionally, the initial thread can assume that its address space is virtually contiguous, and is made up of the smallest frames available on the architecture. It's also assumed that the initial thread knows which paging structures are available on the architecture it's running on. This, along with knowledge of how capabilities in `userImageFrames` and `userImagePaging` are ordered, is sufficient information for userland to infer the virtual address of each frame capability, and the virtual address and type of each paging structure capability.

Untyped memory is given in no particular order. The array entry `untypedList[i]` stores the untyped-memory information of the *i*-th untyped cap of the slot region `untyped`. Therefore, the array length is at least `untyped.end - untyped.start`. The actual length is hardcoded in the kernel and irrelevant to the reader of the array. The

Table 9.3: BootInfoHeader struct.

Field Type	Field Name	Description
seL4_Word	id	Identifier indicating the contents of the chunk
seL4_Uint8	len	Length in bytes of the chunk

untyped memory information is stored in a `seL4_UntypedDesc` struct, described in Table 9.4, and details the address, size and kind of the memory backing the untyped. This allows userland to infer physical memory addresses of retyped frames and use them to initiate DMA transfers when no IOMMU is available. The kernel makes no guarantees about certain sizes of untyped memory being available.

Table 9.4: seL4_UntypedDesc struct

Field Type	Field Name	Description
seL4_Word	paddr	physical base address of the untyped object
seL4_Uint8	sizeBits	size (2^n bytes) of the untyped object
seL4_Uint8	isDevice	is this untyped a device or not (see Section 2.4)
seL4_Uint8[]	padding	manual padding so final struct is a multiple of the word size

If the platform has an seL4-supported IOMMU, `numIOPTLevels` contains the number of IOMMU-page-table levels. This information is needed by userland when constructing an IOMMU address space (`IOspace`). If there is no IOMMU support, `numIOPTLevels` is 0.

On ARM if the platform has any available SMMU units the capabilities for them will be described by the `ioSpaceCaps` slot region. The mapping of a capability from this region to a specific SMMU is platform specific.

9.3 Boot Command-line Arguments

On IA-32, seL4 accepts boot command-line arguments which are passed to the kernel via a multiboot-compliant bootloader (e.g. GRUB, syslinux). Multiple arguments are separated from each other by whitespace. Two forms of arguments are accepted: (1) key-value arguments of the form “key=value” and (2) single keys of the form “key”. The value field of the key-value form may be a string, a decimal integer, a hexadecimal integer beginning with “0x”, or an integer list where list elements are separated by commas. Keys and values can’t have any whitespace in them and there can be no whitespace before or after an “=” or a comma either. Arguments are listed in Table 9.5 along with their default values (if left unspecified).

Table 9.5: IA-32 boot command-line arguments.

Key	Value	Default
<code>console_port</code>	I/O-port base of the serial port that the kernel prints to (if compiled in debug mode)	0x3f8
<code>debug_port</code>	I/O-port base of the serial port that is used for kernel debugging (if compiled in debug mode)	0x3f8
<code>disable_iommu</code>	none	The IOMMU is enabled by default on VT-d-capable platforms

Chapter 10

seL4 API Reference

10.1 Error Codes

Invoking a capability with invalid parameters will result in an error. seL4 system calls return an error code in the message tag and a short error description in the message registers to aid the programmer in determining the cause of errors.

10.1.1 Invalid Argument

A non-capability argument is invalid.

Field	Meaning
Label	<code>seL4_InvalidArgument</code>
IPCBuffer[0]	Invalid argument number

10.1.2 Invalid Capability

A capability argument is invalid.

Field	Meaning
Label	<code>seL4_InvalidCapability</code>
IPCBuffer[0]	Invalid capability argument number

10.1.3 Illegal Operation

The requested operation is not permitted.

Field	Meaning
Label	<code>seL4_IllegalOperation</code>

10.1.4 Range Error

An argument is out of the allowed range.

Field	Meaning
Label	<code>seL4_RangeError</code>
<code>IPCBuffer[0]</code>	Minimum allowed value
<code>IPCBuffer[1]</code>	Maximum allowed value

10.1.5 Alignment Error

A supplied argument does not meet the alignment requirements.

Field	Meaning
Label	<code>seL4_AlignmentError</code>

10.1.6 Failed Lookup

A capability could not be looked up.

Field	Meaning
Label	<code>seL4_FailedLookup</code>
<code>IPCBuffer[0]</code>	1 if the lookup failed for a source capability, 0 otherwise
<code>IPCBuffer[1]</code>	Type of lookup failure
<code>IPCBuffer[2..]</code>	Lookup failure description as described in Section 3.4

10.1.7 Truncated Message

Too few message words or capabilities were sent in the message.

Field	Meaning
Label	<code>seL4_TruncatedMessage</code>

10.1.8 Delete First

A destination slot specified in the syscall arguments is occupied.

Field	Meaning
Label	<code>seL4_DeleteFirst</code>

10.1.9 Revoke First

The object currently has other objects derived from it and the requested invocation cannot be performed until either these objects are deleted or the revoke invocation is performed on the capability.

Field	Meaning
Label	<code>seL4_RevokeFirst</code>

10.1.10 Not Enough Memory

The Untyped Memory object does not have enough unallocated space to complete the `seL4_Untyped_Retype()` request.

Field	Meaning
Label	<code>seL4_NotEnoughMemory</code>
<code>IPCBuffer[0]</code>	Amount of memory available in bytes

10.2 System Calls

10.2.1 General System Calls

This section provides the system call API for non-MCS kernel configurations.

10.2.1.1 Send

```
LIBSEL4_INLINE_FUNC void seL4_Send
```

Send to a capability.

Type	Name	Description
<code>seL4_CPtr</code>	<code>dest</code>	The capability to be invoked.
<code>seL4_MessageInfo_t</code>	<code>msgInfo</code>	The messageinfo structure for the IPC.

Return value: This method does not return anything.

Description: See Section 2.2

10.2.1.2 Recv

LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_Recv

Block until a message is received on an endpoint.

Type	Name	Description
sel4_CPtr	src	The capability to be invoked.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: A sel4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.1.3 Call

LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_Call

Call a capability.

Type	Name	Description
sel4_CPtr	dest	The capability to be invoked.
sel4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.

Return value: A sel4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.1.4 Reply

LIBSEL4_INLINE_FUNC void sel4_Reply

Perform a send to a one-off reply capability stored when the thread was last called.

Type	Name	Description
sel4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.

Return value: This method does not return anything.

Description: See Section 2.2

10.2.1.5 Non-Blocking Send

LIBSEL4_INLINE_FUNC void sel4_NBSend

Perform a non-blocking send to a capability.

Type	Name	Description
sel4_CPtr	dest	The capability to be invoked.
sel4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.

Return value: This method does not return anything.

Description: See Section 2.2

10.2.1.6 Reply Recv

LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_ReplyRecv

Perform a reply followed by a receive in one system call.

Type	Name	Description
sel4_CPtr	dest	The capability to be invoked.
sel4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: A sel4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.1.7 NBRcv

LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_NBRcv

Receive a message from an endpoint but do not block in the case that no messages are pending.

Type	Name	Description
sel4_CPtr	src	The capability to be invoked.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: A sel4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.1.8 Yield

LIBSEL4_INLINE_FUNC void sel4_Yield

Donate the remaining timeslice to a thread of the same priority.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: See Section 2.2

10.2.1.9 Signal

```
LIBSEL4_INLINE_FUNC void sel4_Signal
```

Signal a notification.

Type	Name	Description
sel4_CPtr	dest	The capability to be invoked.

Return value: This method does not return anything.

Description: This is not a proper system call known by the kernel. Rather, it is a convenience wrapper which calls `sel4_Send()`. It is useful for signalling a notification.

See the description of `sel4_Send()` in Section 2.2.

10.2.1.10 Wait

```
LIBSEL4_INLINE_FUNC void sel4_Wait
```

Perform a receive on a notification object.

Type	Name	Description
sel4_CPtr	src	The capability to be invoked.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: This method does not return anything.

Description: This is not a proper system call known by the kernel. Rather, it is a convenience wrapper which calls `sel4_Recv()`.

See the description of `sel4_Recv()` in Section 2.2.

10.2.1.11 Poll

```
LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_Poll
```

Perform a non-blocking recv on a notification object.

Type	Name	Description
<code>sel4_CPtr</code>	<code>src</code>	The capability to be invoked.
<code>sel4_Word *</code>	<code>sender</code>	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: A `sel4_MessageInfo_t` structure as described in Section 4.1

Description: This is not a proper system call known by the kernel. Rather, it is a convenience wrapper which calls `sel4_NBRecv()`. It is useful for doing a non-blocking wait on a notification.

See the description of `sel4_NBRecv()` in Section 2.2.

10.2.2 General System Calls (MCS)

This section provides the system call API for MCS kernel configurations.

10.2.2.1 Send

```
LIBSEL4_INLINE_FUNC void sel4_Send
```

Send to a capability.

Type	Name	Description
sel4_CPtr	dest	The capability to be invoked.
sel4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.

Return value: This method does not return anything.

Description: See Section 2.2

10.2.2.2 Recv

```
LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_Recv
```

Block until a message is received on an endpoint.

Type	Name	Description
sel4_CPtr	src	The capability to be invoked.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.
sel4_CPtr	reply	The capability to the reply object to use on a call (only used on MCS).

Return value: A `sel4_MessageInfo_t` structure as described in Section 4.1

Description: See Section 2.2

10.2.2.3 Call

LIBSEL4_INLINE_FUNC `seL4_MessageInfo_t seL4_Call`

Call a capability.

Type	Name	Description
<code>seL4_CPtr</code>	<code>dest</code>	The capability to be invoked.
<code>seL4_MessageInfo_t</code>	<code>msgInfo</code>	The messageinfo structure for the IPC.

Return value: A `seL4_MessageInfo_t` structure as described in Section 4.1

Description: See Section 2.2

10.2.2.4 Non-Blocking Send

LIBSEL4_INLINE_FUNC `void seL4_NBSend`

Perform a non-blocking send to a capability.

Type	Name	Description
<code>seL4_CPtr</code>	<code>dest</code>	The capability to be invoked.
<code>seL4_MessageInfo_t</code>	<code>msgInfo</code>	The messageinfo structure for the IPC.

Return value: This method does not return anything.

Description: See Section 2.2

10.2.2.5 Reply Recv

LIBSEL4_INLINE_FUNC seL4_MessageInfo_t seL4_ReplyRecv

Perform a reply followed by a receive in one system call.

Type	Name	Description
seL4_CPtr	src	The capability to perform the receive on.
seL4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.
seL4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.
seL4_CPtr	reply	The capability to the reply object, which is first invoked and then used for the recv phase to store a new reply capability.

Return value: A seL4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.2.6 NBRecv

LIBSEL4_INLINE_FUNC seL4_MessageInfo_t seL4_NBRecv

Receive a message from an endpoint but do not block in the case that no messages are pending.

Type	Name	Description
seL4_CPtr	src	The capability to receive on.
seL4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.
seL4_CPtr	reply	The capability to the reply object to use on a call.

Return value: A seL4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.2.7 NBSend Recv

LIBSEL4_INLINE_FUNC seL4_MessageInfo_t seL4_NBSendRecv

Non-blocking send on one capability, and a blocking receive on another in a single system call.

Type	Name	Description
seL4_CPtr	dest	The capability to be invoked.
seL4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.
seL4_CPtr	src	The capability to receive on.
seL4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.
seL4_CPtr	reply	The capability to the reply object, which is first invoked and then used for the recv phase to store a new reply capability.

Return value: A seL4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.2.8 NBSend Wait

LIBSEL4_INLINE_FUNC seL4_MessageInfo_t seL4_NBSendWait

Non-blocking invoke of a capability and wait on another in one system call.

Type	Name	Description
seL4_CPtr	dest	The capability to be invoked.
seL4_MessageInfo_t	msgInfo	The messageinfo structure for the IPC.
seL4_CPtr	src	The capability to receive on.
seL4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: A seL4_MessageInfo_t structure as described in Section 4.1

Description: See Section 2.2

10.2.2.9 Yield

```
LIBSEL4_INLINE_FUNC void sel4_Yield
```

Yield the remaining timeslice. Periodic threads will not be scheduled again until their next sporadic replenishment.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: See Section 2.2

10.2.2.10 Wait

```
LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_Wait
```

Perform a wait on an endpoint or notification object.

Type	Name	Description
sel4_CPtr	src	The capability to be invoked.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: TODO

Description: Block on a notification or endpoint waiting for a message. No reply object is required for a Wait. Wait should not be paired with Call, as it does not provide a reply object. If Wait is paired with a Call the waiter will block after receiving the message.

See the description of `sel4_Wait()` in Section 2.2.

10.2.2.11 NBWait

LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_NBWait

Perform a polling wait on an endpoint or notification object.

Type	Name	Description
sel4_CPtr	src	The capability to be invoked.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: TODO*Description:* Poll a notification or endpoint waiting for a message. No reply object is required for a Wait. Wait should not be paired with Call.See the description of `sel4_NBWait()` in Section 2.2.**10.2.2.12 Poll**

LIBSEL4_INLINE_FUNC sel4_MessageInfo_t sel4_Poll

Perform a non-blocking recv on a notification object.

Type	Name	Description
sel4_CPtr	src	The capability to be invoked.
sel4_Word *	sender	The address to write sender information to. The sender information is the badge of the endpoint capability that was invoked by the sender, or the notification word of the notification object that was signalled. This parameter is ignored if NULL.

Return value: A `sel4_MessageInfo_t` structure as described in Section 4.1*Description:* This is not a proper system call known by the kernel. Rather, it is a convenience wrapper which calls `sel4_NBWait()`. It is useful for doing a non-blocking wait on a notification.See the description of `sel4_NBWait()` in Section 2.2.

10.2.2.13 Signal

LIBSEL4_INLINE_FUNC void sel4_Signal

Signal a notification.

Type	Name	Description
sel4_CPtr	dest	The capability to be invoked.

Return value: This method does not return anything.

Description: This is not a proper system call known by the kernel. Rather, it is a convenience wrapper which calls `sel4_Send()`. It is useful for signalling a notification.

See the description of `sel4_Send()` in Section 2.2.

10.2.3 Debugging System Calls

This section documents debugging system calls available when the kernel is build with the `DEBUG_BUILD` configuration. For any system calls that rely on a kernel serial driver, `PRINTING` must also be enabled.

10.2.3.1 Put Char

```
LIBSEL4_INLINE_FUNC void sel4_DebugPutChar
```

Output a single char through the kernel.

Type	Name	Description
char	c	The character to output.

Return value: This method does not return anything.

Description: Use the kernel serial driver to output a single character. This is useful for debugging when a user level serial driver is not available.

10.2.3.2 Dump scheduler

```
LIBSEL4_INLINE_FUNC void sel4_DebugDumpScheduler
```

Output the contents of the kernel scheduler.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: Dump the state of the all TCB objects to kernel serial output. This system call will output a table containing:

- Address: the address of the TCB object for that thread,
- Name: the name of the thread (if set),
- IP: the contents of the instruction pointer the thread is at,
- Priority: the priority of that thread,
- State : the state of the thread.

10.2.3.3 Halt

LIBSEL4_INLINE_FUNC void sel4_DebugHalt

Halt the system.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: This debugging system call will cause the kernel immediately cease responding to system calls. The kernel will switch permanently to the idle thread with interrupts disabled. Depending on the platform, the kernel may switch the hardware into a low-power state.

10.2.3.4 Snapshot

LIBSEL4_INLINE_FUNC void sel4_DebugSnapshot

Output a capDL dump of the current kernel state.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: This debugging system call will output all of the capabilities in the current kernel using capDL.

10.2.3.5 Cap Identify

LIBSEL4_INLINE_FUNC seL4_Uint32 sel4_DebugCapIdentify

Identify the type of a capability in the current cspace.

Type	Name	Description
seL4_CPtr	cap	A capability slot in the current cspace.

Return value: The type of capability passed in.

Description: This debugging system call returns the type of capability in a capability slot in the current cspace. The type returned is not a libsel4 type, but refers to an internal seL4 type. This can be looked up in a built kernel by looking for the (generated) enum `cap_tag`, type `cap_tag_t`.

10.2.3.6 Name Thread

LIBSEL4_INLINE_FUNC void sel4_DebugNameThread

Name a thread.

Type	Name	Description
sel4_CPtr	tcb	A capability to the tcb object for the thread to name.
const char *	name	The name for the thread.

Return value: This method does not return anything.

Description: Name a thread. This name will then be output by the kernel in all debugging output. Note that the max name length that can be passed to this function is limited by the number of chars that will fit in an IPC message (`sel4_MsgMaxLength` multiplied by the amount of chars that fit in a word). However the name is also truncated in order to fit into a TCB object. For some platforms you may need to increase `sel4_TCBBits` by 1 in a debug build in order to fit a long enough name.

10.2.3.7 Send SGI 0-15

LIBSEL4_INLINE_FUNC void sel4_DebugSendIPI

Sends arbitrary SGI.

Type	Name	Description
sel4_Uint8	target	The target core ID.
unsigned	irq	The SGI number (0-15).

Return value: This method does not return anything.

Description: Send an arbitrary SGI (core-specific interrupt 0-15) to the specified target core.

10.2.3.8 Run

LIBSEL4_INLINE_FUNC void seL4_DebugRun

Run a user level function in kernel mode.

Type	Name	Description
void(*) (void *)	userfn	The address in userspace of the function to run.
void *	userarg	A single argument to pass to the function.

Return value: This method does not return anything.

Description: This extremely dangerous function is for running benchmarking and debugging code that needs to be executed in kernel mode from userlevel. It should never be used in a release kernel. This works because the kernel can access all user mappings of device memory, and does not switch page directories on kernel entry.

Unlike the other system calls in this section, `seL4_DebugRun` does not depend on the `DEBUG_BUILD` configuration option, but its own config variable `DANGEROUS_CODE_INJECTION`.

10.2.4 Benchmarking System Calls

This section documents system calls available when the kernel is configured with benchmarking enabled. There are several different benchmarking modes which can be configured when building the kernel:

1. `BENCHMARK_TRACEPOINTS`: Enable using tracepoints in the kernel and timing code.
2. `BENCHMARK_TRACK_KERNEL_ENTRIES`: Keep track of information on kernel entries.
3. `BENCHMARK_TRACK_UTILISATION`: Allow users to get CPU timing info for the system, threads and/or idle thread.

10.2.4.1 Reset Log

```
LIBSEL4_INLINE_FUNC sel4_Error sel4_BenchmarkResetLog
```

Reset benchmark logging.

Type	Name	Description
void		

Return value: A `sel4_Error` error if the user-level log buffer has not been set by the user (`BENCHMARK_TRACEPOINTS/BENCHMARK_TRACK_KERNEL_ENTRIES`).

Description: The behaviour of this system call depends on benchmarking mode in action while invoking this system call:

1. `BENCHMARK_TRACEPOINTS`: resets the log index to 0,
2. `BENCHMARK_TRACK_KERNEL_ENTRIES`: as above,
3. `BENCHMARK_TRACK_UTILISATION`: resets benchmark and current thread start time (to the time of invoking this syscall), resets idle thread utilisation to 0, and starts tracking utilisation.

10.2.4.2 Finalize Log

LIBSEL4_INLINE_FUNC seL4_Word seL4_BenchmarkFinalizeLog

Stop benchmark logging.

Type	Name	Description
void		

Return value: The index of the final entry in the log buffer (if BENCHMARK_TRACEPOINTS/BENCHMARK_TRACK_KERNEL_ENTRIES are enabled).

Description: The behaviour of this system call depends on benchmarking mode in action while invoking this system call:

1. BENCHMARK_TRACEPOINTS: Sets the final log buffer index to the current index,
2. BENCHMARK_TRACK_KERNEL_ENTRIES: as above,
3. BENCHMARK_TRACK_UTILISATION: sets benchmark end time to current time, stops tracking utilisation.

10.2.4.3 Set Log Buffer

LIBSEL4_INLINE_FUNC seL4_Error seL4_BenchmarkSetLogBuffer

Set log buffer.

Type	Name	Description
seL4_Word	frame_cptr	A capability pointer to a user allocated frame of seL4_LargePage size.

Return value: A seL4_IllegalOperation error if frame_cptr is not valid and couldn't set the buffer.

Description: Provide a large frame object for the kernel to use as a log-buffer. The object must not be device memory, and must be seL4_LargePageBits in size.

10.2.4.4 Null Syscall

LIBSEL4_INLINE_FUNC void sel4_BenchmarkNullSyscall

Null system call that enters and exits the kernel immediately, for timing kernel traps in microbenchmarks.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: Used to time kernel traps (in and out).

10.2.4.5 Flush Caches

LIBSEL4_INLINE_FUNC void sel4_BenchmarkFlushCaches

Flush hardware caches.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: Flush all possible hardware caches for this platform.

10.2.4.6 Flush L1 Caches

LIBSEL4_INLINE_FUNC void sel4_BenchmarkFlushL1Caches

Flush L1 caches.

Type	Name	Description
sel4_Word	cache_type	L1 Cache Type to be flushed

Return value: This method does not return anything.

Description: Flush L1 caches for this platform (currently only support for ARM). Allow to specify the cache type to be flushed (i.e. instruction cache only, data cache only and both instruction cache and data cache).

10.2.4.7 Get Thread Utilisation

LIBSEL4_INLINE_FUNC void sel4_BenchmarkGetThreadUtilisation

Get utilisation timing information.

Type	Name	Description
sel4_Word	tcb_cptr	TCB cap pointer to a thread to get CPU utilisation for.

Return value: This method does not return anything.

Description: Get timing information for the system, requested thread and idle thread. Such information is written into the caller's IPC buffer; see the definition of `benchmark_track_util_ipc_index` enum for more details on the data/format returned on the IPC buffer.

10.2.4.8 Reset Thread Utilisation

LIBSEL4_INLINE_FUNC void sel4_BenchmarkResetThreadUtilisation

Reset utilisation timing for a specific thread.

Type	Name	Description
sel4_Word	tcb_cptr	TCB cap pointer to a thread to get CPU utilisation for.

Return value: This method does not return anything.

Description: Reset the kernel's timing information data (start time and utilisation) for a specific thread.

10.2.4.9 Dump All Threads Utilisation

LIBSEL4_INLINE_FUNC void sel4_BenchmarkDumpAllThreadsUtilisation

Print the current accumulated cycle count for every thread on the current node.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: Uses kernel's printf to print number of cycles on each line in the following format: `thread_name,thread_cycles`

10.2.4.10 Reset All Threads Utilisation

`LIBSEL4_INLINE_FUNC void sel4_BenchmarkResetAllThreadsUtilisation`

Reset the accumulated cycle count for every thread on the current node.

Type	Name	Description
void		

Return value: This method does not return anything.

Description: Reset the cycle count for each thread to 0.

10.2.5 X86 System Calls

10.2.5.1 VMEnter

```
LIBSEL4_INLINE_FUNC sel4_Word sel4_VMEnter
```

Change current thread to execute from its bound VCPU.

Type	Name	Description
sel4_Word *	sender	The address to write sender information to. If the syscall returns due to receiving a notification on the bound notification then the sender information is the badge of the notification capability that was invoked. This parameter is ignored if NULL.

Return value: SEL4_VMENTER_RESULT_NOTIF if a notification was received or SEL4_VMENTER_RESULT_FAULT if the guest mode execution faulted for any reason

Description: Changes the execution mode of the current thread from normal TCB execution, to guest execution using its bound VCPU. For details on VCPUs and execution modes see Section 6.4.

Invoking `sel4_VMEnter` is similar to replying to a fault in that updates to the registers can be given in the message, but unlike a fault no message info (see Section 4.1) is sent as the registers are not optional and the number that must be sent is fixed. The mapping of hardware register to message register is

- SEL4_VMENTER_CALL_EIP_MR Address to start executing instructions at in the guest mode
- SEL4_VMENTER_CALL_CONTROL_PPC_MR New value for the Primary Processor Based VM Execution Controls
- SEL4_VMENTER_CALL_CONTROL_ENTRY_MR New value for the VM Entry Controls

On return these same three message registers will be filled with the values at the point that the privileged mode ceased executing. If this function returns with SEL4_VMENTER_RESULT_FAULT then the following additional message registers will be filled out

- SEL4_VMENTER_FAULT_REASON_MR
- SEL4_VMENTER_FAULT_QUALIFICATION_MR
- SEL4_VMENTER_FAULT_INSTRUCTION_LEN_MR
- SEL4_VMENTER_FAULT_GUEST_PHYSICAL_MR
- SEL4_VMENTER_FAULT_RFLAGS_MR
- SEL4_VMENTER_FAULT_GUEST_INT_MR

- SEL4_VMENTER_FAULT_CR3_MR
- SEL4_VMENTER_FAULT_EAX
- SEL4_VMENTER_FAULT_EBX
- SEL4_VMENTER_FAULT_ECX
- SEL4_VMENTER_FAULT_EDX
- SEL4_VMENTER_FAULT_ESI
- SEL4_VMENTER_FAULT_EDI
- SEL4_VMENTER_FAULT_EBP

10.3 Architecture-Independent Object Methods

10.3.1 seL4_CNode

10.3.1.1 Cancel Badged Sends

```
static inline int seL4_CNode_CancelBadgedSends
```

The cancel badged sends method is intend to allow for the reuse of badges by an authority. When used with a badged endpoint capability it will cancel any outstanding send operations for that endpoint and badge. This operation has no effect on un-badged or other objects.

Type	Name	Description
seL4_CNode	<code>_service</code>	CPTR to the CNode at the root of the CSpace where the capability will be found. Must be at a depth equivalent to the wordsize.
seL4_Word	<code>index</code>	CPTR to the capability. Resolved from the root of the <code>_service</code> parameter.
seL4_Uint8	<code>depth</code>	Number of bits of index to resolve to find the capability being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.2 Copy

```
static inline int seL4_CNode_Copy
```

Copy a capability, setting its access rights whilst doing so

Type	Name	Description
seL4_CNode	<code>_service</code>	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	<code>dest_index</code>	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_UInt8	<code>dest_depth</code>	Number of bits of <code>dest_index</code> to resolve to find the destination slot.
seL4_CNode	<code>src_root</code>	CPTR to the CNode that forms the root of the source CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	<code>src_index</code>	CPTR to the source slot. Resolved from the root of the source CSpace.
seL4_UInt8	<code>src_depth</code>	Number of bits of <code>src_index</code> to resolve to find the source slot.
seL4_CapRights_t	<code>rights</code>	The rights inherited by the new capability. Possible values for this type are given in Section 3.1.4 .

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.3 Delete

```
static inline int sel4_CNode_Delete
```

Delete a capability

Type	Name	Description
sel4_CNode	_service	CPTR to the CNode at the root of the CSpace where the capability will be found. Must be at a depth equivalent to the wordsize.
sel4_Word	index	CPTR to the capability. Resolved from the root of the _service parameter.
sel4_Uint8	depth	Number of bits of index to resolve to find the capability being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.4 Mint

```
static inline int seL4_CNode_Mint
```

Copy a capability, setting its access rights and badge whilst doing so

Type	Name	Description
seL4_CNode	_service	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	dest_index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_UInt8	dest_depth	Number of bits of dest_index to resolve to find the destination slot.
seL4_CNode	src_root	CPTR to the CNode that forms the root of the source CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	src_index	CPTR to the source slot. Resolved from the root of the source CSpace.
seL4_UInt8	src_depth	Number of bits of src_index to resolve to find the source slot.
seL4_CapRights_t	rights	The rights inherited by the new capability. Possible values for this type are given in Section 3.1.4 .
seL4_Word	badge	Badge or guard to be applied to the new capability. For badges on 32-bit platforms, the high 4 bits are ignored.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.5 Move

```
static inline int sel4_CNode_Move
```

Move a capability

Type	Name	Description
sel4_CNode	_service	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
sel4_Word	dest_index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
sel4_Uint8	dest_depth	Number of bits of dest_index to resolve to find the destination slot.
sel4_CNode	src_root	CPTR to the CNode that forms the root of the source CSpace. Must be at a depth equivalent to the wordsize.
sel4_Word	src_index	CPTR to the source slot. Resolved from the root of the source CSpace.
sel4_Uint8	src_depth	Number of bits of src_index to resolve to find the source slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.6 Mutate

```
static inline int seL4_CNode_Mutate
```

Move a capability, setting its badge in the process

Type	Name	Description
seL4_CNode	_service	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	dest_index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_UInt8	dest_depth	Number of bits of dest_index to resolve to find the destination slot.
seL4_CNode	src_root	CPTR to the CNode that forms the root of the source CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	src_index	CPTR to the source slot. Resolved from the root of the source CSpace.
seL4_UInt8	src_depth	Number of bits of src_index to resolve to find the source slot.
seL4_Word	badge	Badge or guard to be applied to the new capability. For badges on 32-bit platforms, the high 4 bits are ignored.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.7 Revoke

```
static inline int sel4_CNode_Revoke
```

Delete all child capabilities of a capability

Type	Name	Description
sel4_CNode	_service	CPTR to the CNode at the root of the CSpace where the capability will be found. Must be at a depth equivalent to the wordsize.
sel4_Word	index	CPTR to the capability. Resolved from the root of the _service parameter.
sel4_Uint8	depth	Number of bits of index to resolve to find the capability being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.8 Rotate

```
static inline int seL4_CNode_Rotate
```

Given 3 capability slots - a destination, pivot and source - move the capability in the pivot slot to the destination slot and the capability in the source slot to the pivot slot

Type	Name	Description
seL4_CNode	_service	CPTR to the CNode at the root of the CSpace where the destination slot will be found. Must be at a depth equivalent to the wordsize.
seL4_Word	dest_index	CPTR to the destination slot. Resolved relative to _service. Must be empty unless it refers to the same slot as the source slot.
seL4_UInt8	dest_depth	Depth to resolve dest_index to.
seL4_Word	dest_badge	The new capdata for the capability that ends up in the destination slot.
seL4_CNode	pivot_root	CPTR to the CNode at the root of the CSpace where the pivot slot will be found. Must be at a depth equivalent to the wordsize.
seL4_Word	pivot_index	CPTR to the pivot slot. Resolved relative to pivot_root. The resolved slot must not refer to the source or destination slots.
seL4_UInt8	pivot_depth	Depth to resolve pivot_index to.
seL4_Word	pivot_badge	The new capdata for the capability that ends up in the pivot slot.
seL4_CNode	src_root	CPTR to the CNode at the root of the CSpace where the source slot will be found. Must be at a depth equivalent to the wordsize.
seL4_Word	src_index	CPTR to the source slot. Resolved relative to src_root.
seL4_UInt8	src_depth	Depth to resolve src_index to.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.1.9 Save Caller

```
static inline int sel4_CNode_SaveCaller
```

Save the reply capability from the last time the thread was called in the given CSpace so that it can be invoked later

Type	Name	Description
sel4_CNode	_service	CPTR to the CNode at the root of the CSpace where the capability is to be saved. Must be at a depth equivalent to the wordsize.
sel4_Word	index	CPTR to the slot in which to save the capability. Resolved from the root of the _service parameter.
sel4_Uint8	depth	Number of bits of index to resolve to find the slot being targeted.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 3.1.2.

10.3.2 sel4_DomainSet

10.3.2.1 Set

```
static inline int sel4_DomainSet_Set
```

Change the domain of a thread.

Type	Name	Description
sel4_DomainSet	_service	Capability allowing domain configuration.
sel4_Uint8	domain	The thread's new domain.
sel4_TCB	thread	Capability to the TCB which is being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.3.

10.3.3 seL4_IRQControl

10.3.3.1 Get

```
static inline int seL4_IRQControl_Get
```

Create an IRQ handler capability

Type	Name	Description
seL4_IRQControl	_service	An IRQControl capability. This gives you the authority to make this call.
seL4_Word	irq	The IRQ that you want this capability to handle.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	depth	Number of bits of dest_index to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1.

10.3.4 seL4_IRQHandler

10.3.4.1 Acknowledge

```
static inline int seL4_IRQHandler_Ack
```

Acknowledge the receipt of an interrupt and re-enable it

Type	Name	Description
seL4_IRQHandler	_service	The IRQ handler capability.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1.

10.3.4.2 Clear

```
static inline int sel4_IRQHandler_Clear
```

Clear the handler capability from the IRQ slot

Type	Name	Description
sel4_IRQHandler	_service	The IRQ handler capability.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1.

10.3.4.3 Set Notification

```
static inline int sel4_IRQHandler_SetNotification
```

Set the notification which the kernel will signal on interrupts controlled by the supplied IRQ handler capability

Type	Name	Description
sel4_IRQHandler	_service	The IRQ handler capability.
sel4_CPtr	notification	The notification which the IRQs will signal.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1.

10.3.5 seL4_SchedContext

10.3.5.1 Bind

```
static inline int seL4_SchedContext_Bind
```

Bind an object to a scheduling context. The object can be a notification object or a thread.

If the object is a thread and the thread is in a runnable state and the scheduling context has available budget, this will start the thread running.

If the object is a notification, when passive threads wait on the notification object and a signal arrives, the passive thread will receive the scheduling context and possess it until it waits on the notification object again.

This operation will fail if the scheduling context is already bound to a thread or notification object.

Type	Name	Description
seL4_SchedContext	_service	TODO
seL4_CPtr	cap	Capability to a TCB or a notification object

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1

10.3.5.2 Consumed

```
static inline seL4_SchedContext_Consumed_t seL4_SchedContext_Consumed
```

Return the amount of time used by this scheduling context since this function was last called or a timeout exception triggered.

Type	Name	Description
seL4_SchedContext	_service	TODO

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1

10.3.5.3 Unbind

```
static inline int seL4_SchedContext_Unbind
```

Unbind any objects (threads or notification objects) from a scheduling context. This will render the bound thread passive, see Section 6.1.5.

Type	Name	Description
seL4_SchedContext	_service	TODO

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1

10.3.5.4 UnbindObject

```
static inline int seL4_SchedContext_UnbindObject
```

Unbind an object from a scheduling context. The object can be either a thread or a notification.

If the thread being unbound is the thread that is bound to this scheduling context, this will render the thread passive. However if the thread being unbound received the scheduling context via scheduling context donation over IPC, the scheduling context will be returned to the thread that it was originally bound to.

If the object is a notification and it is bound to the scheduling context, unbind it.

Type	Name	Description
seL4_SchedContext	_service	TODO
seL4_CPtr	cap	Capability to a notification that is bound to the scheduling context or capability to a tcb that is bound to this scheduling context or has received it through scheduling context donation.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.7

10.3.5.5 YieldTo

```
static inline seL4_SchedContext_YieldTo_t seL4_SchedContext_YieldTo
```

If a thread is currently runnable and running on this scheduling context and the scheduling context has available budget, place it at the head of the scheduling queue. If the caller is at an equal priority to the thread this will result in the thread being scheduled. If the caller is at a higher priority the thread will not run until the threads priority is the highest priority in the system. The caller must have a maximum control priority greater than or equal to the threads priority.

Type	Name	Description
seL4_SchedContext	_service	TODO

Return value: TODO

Description: TODO

10.3.6 seL4_SchedControl

10.3.6.1 ConfigureFlags

```
static inline int seL4_SchedControl_ConfigureFlags
```

Set the parameters of a scheduling context by invoking the scheduling control capability. If the scheduling context is bound to a currently running thread, the parameters will take effect immediately: that is the current budget will be increased or reduced by the difference between the new and previous budget and the replenishment time will be updated according to any difference in the period. This can result in active threads being post-poned or released depending on the nature of the parameter change and the state of the thread. Additionally, if the scheduling context was previously empty (no budget) but bound to a runnable thread, this can result in a thread running for the first time since it now has access to CPU time. This call will return seL4 Invalid Argument if the parameters are too small (smaller than the kernel WCET for this platform) or too large (will overflow the timer).

Type	Name	Description
seL4_SchedControl	_service	TODO
seL4_SchedContext	schedcontext	Capability to the scheduling context which is being operated on.
seL4_Time	budget	Timeslice in microseconds, when the budget expires the thread will be pre-empted.
seL4_Time	period	Period in microseconds, if equal to budget, this thread will be treated as a round-robin thread. Otherwise, sporadic servers will be used to assure the scheduling context does not exceed the budget over the specified period.
seL4_Word	extra_refills	Number of extra sporadic replenishments this scheduling context should use. Ignored for round-robin threads.
seL4_Word	badge	Identifier for this scheduling context. Delivered to timeout exception handler. Can be used to determine which scheduling context triggered the timeout.
seL4_Word	flags	Bitwise OR'd set of seL4_SchedContextFlag.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1

10.3.7 seL4_TCB

10.3.7.1 Bind Notification

```
static inline int seL4_TCB_BindNotification
```

Binds a notification object to a TCB

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_CPtr	notification	Notification to bind.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 5.3

10.3.7.2 Configure (MCS)

```
static inline int seL4_TCB_Configure
```

Set the parameters of a TCB

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_CNode	cspace_root	The new CSpace root.
seL4_Word	cspace_root_data	Optionally set the guard and guard size of the new root CNode. If set to zero, this parameter has no effect.
seL4_CPtr	vspace_root	The new VSpace root.
seL4_Word	vspace_root_data	Has no effect on x86 or ARM processors.
seL4_Word	buffer	Location of the thread's IPC buffer. Must be 512-byte aligned. The IPC buffer may not cross a page boundary.
seL4_CPtr	bufferFrame	Capability to a page containing the thread's IPC buffer.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1

10.3.7.3 Configure Single Stepping

```
static inline seL4_TCB_ConfigureSingleStepping_t seL4_TCB_ConfigureSingleStepping
```

Set or modify single stepping options for the target TCB. Subsequent calls to this function overwrite previous configuration. Depending on your processor architecture, this may or may not require the consumption of a hardware register.

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_Uint16	bp_num	The API-ID of a target breakpoint. This ID will be a positive integer, with values ranging from 0 to seL4_NumHWBreakpoints - 1.
seL4_Word	num_instructions	Number of instructions to step over before delivering a fault to the target thread's fault endpoint. Setting this to 0 disables single-stepping.

Return value: A `seL4_TCB_ConfigureSingleStepping_t`: Struct that contains `seL4_Error error`, an seL4 API error value, `seL4_Bool bp_was_consumed`, a boolean which indicates whether or not the `bp_num` breakpoint ID that was passed to the function, was consumed in the setup of the single-stepping functionality: if this is `true`, the caller should not attempt to re-use `bp_num` until it has disabled the single-stepping functionality via a subsequent call to `seL4_TCB_ConfigureSingleStepping` with an `num_instructions` argument of 0.

Description: See Sections 6.2.5 and 6.2.4

10.3.7.4 Configure

```
static inline int seL4_TCB_Configure
```

Set the parameters of a TCB

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_Word	fault_ep	CPTR to the endpoint which receives IPCs when this thread faults. This capability is in the CSpace of the thread being configured.
seL4_CNode	cspace_root	The new CSpace root.
seL4_Word	cspace_root_data	Optionally set the guard and guard size of the new root CNode. If set to zero, this parameter has no effect.
seL4_CPtr	vspace_root	The new VSpace root.
seL4_Word	vspace_root_data	Has no effect on x86 or ARM processors.
seL4_Word	buffer	Location of the thread's IPC buffer. Must be 512-byte aligned. The IPC buffer may not cross a page boundary.
seL4_CPtr	bufferFrame	Capability to a page containing the thread's IPC buffer.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1

10.3.7.5 Copy Registers

```
static inline int sel4_TCB_CopyRegisters
```

Copy the registers from one thread to another

Type	Name	Description
sel4_TCB	_service	Capability to the TCB which is being operated on. This is the destination TCB.
sel4_TCB	source	Cap to the source TCB.
sel4_Bool	suspend_source	The invocation should also suspend the source thread.
sel4_Bool	resume_target	The invocation should also resume the destination thread.
sel4_Bool	transfer_frame	Frame registers should be transferred.
sel4_Bool	transfer_integer	Integer registers should be transferred.
sel4_Uint8	arch_flags	Architecture dependent flags. These have no meaning on either x86 or ARM.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: In the context of this function, frame registers are those that are read, modified or preserved by a system call and integer registers are those that are not. Refer to the sel4 userland library source for specifics. Section 6.1.3

10.3.7.6 Get Breakpoint

```
static inline seL4_TCB_GetBreakpoint_t seL4_TCB_GetBreakpoint
```

Read a breakpoint or watchpoint's current configuration.

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_Uint16	bp_num	The API-ID of a target breakpoint. This ID will be a positive integer, with values ranging from 0 to seL4_NumHWBreakpoints - 1.

Return value: A `seL4_TCB_GetBreakpoint_t`: Struct that contains `seL4_Error error`, an seL4 API error value, `seL4_Word vaddr`, the virtual address at which the breakpoint will currently be triggered; `seL4_Word type`, the type of operation which will currently trigger the breakpoint, whether instruction execution, or data access; `seL4_Word size`, integer value for the span-size of the breakpoint. Usually a power of two (1, 2, 4, etc.); `seL4_Word rw`, the access direction that will currently trigger the breakpoint, whether read, write, or both and `seL4_Bool is_enabled`, which indicates whether or not the breakpoint will currently be triggered if the match conditions are met.

Description: See Section 6.2.4

10.3.7.7 Read Registers

```
static inline int sel4_TCB_ReadRegisters
```

Read a thread's registers into the first `count` fields of a given `sel4_UserContext`

Type	Name	Description
<code>sel4_TCB</code>	<code>_service</code>	Capability to the TCB which is being operated on.
<code>sel4_Bool</code>	<code>suspend_source</code>	The invocation should also suspend the source thread.
<code>sel4_Uint8</code>	<code>arch_flags</code>	Architecture dependent flags. These have no meaning on either x86 or ARM.
<code>sel4_Word</code>	<code>count</code>	The number of registers to read.
<code>sel4_UserContext *</code>	<code>regs</code>	The structure to read the registers into.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.12

10.3.7.8 Resume

```
static inline int sel4_TCB_Resume
```

Resume a thread

Type	Name	Description
<code>sel4_TCB</code>	<code>_service</code>	Capability to the TCB which is being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.3

10.3.7.9 Set Breakpoint

```
static inline int seL4_TCB_SetBreakpoint
```

Set or modify a thread's breakpoints or watchpoints. Calls to this function overwrite previous configurations for the target breakpoint. Do not use this with `seL4_SingleStep`: the API will reject the call and return an error. Instead, use `seL4_TCB_ConfigureSingleStepping` to configure single-stepping.

Type	Name	Description
<code>seL4_TCB</code>	<code>_service</code>	Capability to the TCB which is being operated on.
<code>seL4_Uint16</code>	<code>bp_num</code>	The API-ID of a target breakpoint. This ID will be a positive integer, with values ranging from 0 to <code>seL4_NumHWBreakpoints - 1</code> .
<code>seL4_Word</code>	<code>vaddr</code>	A virtual address which forms part of the match conditions for the triggering of the breakpoint.
<code>seL4_Word</code>	<code>type</code>	One of: <code>seL4_InstructionBreakpoint</code> , which specifies that the breakpoint should occur on instruction execution at the specified <code>vaddr</code> or <code>seL4_DataBreakpoint</code> , which states that the breakpoint should occur on data access at the specified <code>vaddr</code> .
<code>seL4_Word</code>	<code>size</code>	A positive integer indicating the trigger-span of the watchpoint. Must be zero when 'type' is <code>seL4_InstructionBreakpoint</code> .
<code>seL4_Word</code>	<code>rw</code>	One of <code>seL4_BreakOnRead</code> , meaning the breakpoint will only be triggered on read-access; <code>seL4_BreakOnWrite</code> meaning the breakpoint will only be triggered on write-access, and <code>seL4_BreakOnReadWrite</code> meaning the breakpoint will be triggered on any access.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.2.4

10.3.7.10 Set CPU Affinity

```
static inline int sel4_TCB_SetAffinity
```

Change a thread's current CPU in multicore machine

Type	Name	Description
sel4_TCB	_service	Capability to the TCB which is being operated on.
sel4_Word	affinity	The thread's new CPU to run.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.2

10.3.7.11 Set IPC Buffer

```
static inline int sel4_TCB_SetIPCBuffer
```

Set a thread's IPC buffer

Type	Name	Description
sel4_TCB	_service	Capability to the TCB which is being operated on.
sel4_Word	buffer	Location of the thread's IPC buffer. Must be 512-byte aligned. The IPC buffer may not cross a page boundary.
sel4_CPtr	bufferFrame	Capability to a page containing the thread's IPC buffer.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Sections 6.1 and 4.1

10.3.7.12 Set Maximum Controlled Priority

```
static inline int seL4_TCB_SetMCPriority
```

Change a thread's maximum controlled priority

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_TCB	authority	Capability to the TCB to use the MCP from when setting the MCP.
seL4_Word	mcp	The thread's new maximum controlled priority.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.4

10.3.7.13 Set Priority

```
static inline int seL4_TCB_SetPriority
```

Change a thread's priority

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_TCB	authority	Capability to the TCB to use the MCP from when setting the priority.
seL4_Word	priority	The thread's new priority.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.4

10.3.7.14 Set Sched Params (MCS)

```
static inline int sel4_TCB_SetSchedParams
```

Change a thread's priority, maximum controlled priority, scheduling context and fault handler.

Type	Name	Description
sel4_TCB	_service	Capability to the TCB which is being operated on.
sel4_TCB	authority	Capability to the TCB to use the MCP from when setting the priority and MCP.
sel4_Word	mcp	The thread's new maximum controlled priority.
sel4_Word	priority	The thread's new priority.
sel4_CPtr	sched_context	Capability to the scheduling context that the TCB should run on. If the scheduling context is already bound to a notification or TCB that is not this TCB this operation will fail. Similarly, if this TCB is already bound to a scheduling context that is not this scheduling context, this will also fail.
sel4_CPtr	fault_ep	CPTR to the endpoint which receives IPCs when this thread faults.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.4

10.3.7.15 Set Sched Params

```
static inline int seL4_TCB_SetSchedParams
```

Change a thread's priority and maximum controlled priority.

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_TCB	authority	Capability to the TCB to use the MCP from when setting the priority and MCP.
seL4_Word	mcp	The thread's new maximum controlled priority.
seL4_Word	priority	The thread's new priority.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.4

10.3.7.16 Set Space

```
static inline int seL4_TCB_SetSpace
```

Set the fault endpoint, CSpace and VSpace of a thread

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.
seL4_Word	fault_ep	CPTR to the endpoint which receives IPCs when this thread faults. On MCS this cap gets copied into the TCB.
seL4_CNode	cspace_root	The new CSpace root.
seL4_Word	cspace_root_data	Optionally set the guard and guard size of the new root CNode. If set to zero, this parameter has no effect.
seL4_CPtr	vspace_root	The new VSpace root.
seL4_Word	vspace_root_data	Has no effect on x86 or ARM processors.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1

See Section 6.1

10.3.7.17 Set TLS Base

```
static inline int sel4_TCB_SetTLSBase
```

Set the TLS base of the target TCB. This intended for use on architectures where the register used for TLS can only be written to in priviledged mode

Type	Name	Description
sel4_TCB	_service	Capability to the TCB which is being operated on.
sel4_Word	tls_base	The TLS base to set

Return value: TODO

Description: TODO

10.3.7.18 Set Timeout Endpoint

```
static inline int sel4_TCB_SetTimeoutEndpoint
```

Set a thread's timeout endpoint.

Type	Name	Description
sel4_TCB	_service	Capability to the TCB which is being operated on.
sel4_CPtr	timeout_fault_ep	CPtr to the endpoint which receives IPCs when this thread triggers timeout faults. Can be null.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Timeout exception messages will be delivered to this endpoint if it is not a null capability.

10.3.7.19 Suspend

```
static inline int seL4_TCB_Suspend
```

Suspend a thread

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.3

10.3.7.20 Unbind Notification

```
static inline int seL4_TCB_UnbindNotification
```

Unbinds any notification object from a TCB

Type	Name	Description
seL4_TCB	_service	Capability to the TCB which is being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 5.3

10.3.7.21 Unset Breakpoint

```
static inline int seL4_TCB_UnsetBreakpoint
```

Disables a hardware breakpoint or watchpoint. The caller should assume that the underlying configuration of the hardware registers has also been cleared. Do not use this to clear single-stepping; the API will reject the call and return an error. Instead, use `seL4_TCB_ConfigureSingleStepping` to disable single-stepping.

Type	Name	Description
<code>seL4_TCB</code>	<code>_service</code>	Capability to the TCB which is being operated on.
<code>seL4_Uint16</code>	<code>bp_num</code>	The API-ID of a target breakpoint. This ID will be a positive integer, with values ranging from 0 to <code>seL4_NumHWBreakpoints - 1</code> .

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.2.4

10.3.7.22 Write Registers

```
static inline int seL4_TCB_WriteRegisters
```

Set a thread's registers to the first `count` fields of a given `seL4_UserContext`

Type	Name	Description
<code>seL4_TCB</code>	<code>_service</code>	Capability to the TCB which is being operated on.
<code>seL4_Bool</code>	<code>resume_target</code>	The invocation should also resume the destination thread.
<code>seL4_Uint8</code>	<code>arch_flags</code>	Architecture dependent flags. These have no meaning on either x86 or ARM.
<code>seL4_Word</code>	<code>count</code>	The number of registers to be set.
<code>seL4_UserContext *</code>	<code>regs</code>	Data structure containing the new register values.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.1.12

10.3.8 seL4_Untyped

10.3.8.1 Retype

```
static inline int seL4_Untyped_Retype
```

Retype an untyped object

Type	Name	Description
seL4_Untyped	<code>_service</code>	CPTR to an untyped object.
seL4_Word	<code>type</code>	The seL4 object type that we are retyping to.
seL4_Word	<code>size_bits</code>	Used to determine the size of variable-sized objects.
seL4_CNode	<code>root</code>	CPTR to the CNode at the root of the destination CSpace.
seL4_Word	<code>node_index</code>	CPTR to the destination CNode. Resolved relative to the root parameter.
seL4_Word	<code>node_depth</code>	Number of bits of <code>node_index</code> to translate when addressing the destination CNode.
seL4_Word	<code>node_offset</code>	Number of slots into the node at which capabilities start being placed.
seL4_Word	<code>num_objects</code>	Number of capabilities to create.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Given a capability, `_service`, to an untyped object, creates `num_objects` of the requested type. Creates `num_objects` capabilities to the new objects starting at `node_offset` in the CNode specified by `root`, `node_index`, and `node_depth`.

For variable-sized kernel objects, the `size_bits` argument is used to determine the size of objects to create. The relationship between `size_bits` and object size depends on the type of object being created. See Section 2.4.2 for more information about object sizes. See Section 2.4 for more information about how untyped memory is retyped. See Section 3.1.3 for more information about the placement of capabilities to created objects.

10.4 x86-Specific Object Methods

10.4.1 seL4_IRQControl

10.4.1.1 Get I/O APIC

```
static inline int seL4_IRQControl_GetIOAPIC
```

Create an IRQ handler capability for an interrupt from an IOAPIC.

Type	Name	Description
seL4_IRQControl	_service	An IRQControl capability. This gives you the authority to make this call.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_UInt8	depth	Number of bits of index to resolve to find the destination slot.
seL4_Word	ioapic	Zero based index of IOAPIC to get interrupt from, ordered the same as in ACPI tables
seL4_Word	pin	IOAPIC pin that generates the interrupt.
seL4_Word	level	Indicates whether the IOAPIC should be programmed to treat this interrupt as level triggered.
seL4_Word	polarity	Indicates whether the IOAPIC should be programmed to treat this interrupt as high or low triggered
seL4_Word	vector	CPU vector to deliver the interrupt to.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1 and Section 8.2.1.

10.4.1.2 Get MSI

```
static inline int seL4_IRQControl_GetMSI
```

Create an IRQ handler capability for an interrupt from an MSI.

Type	Name	Description
seL4_IRQControl	_service	An IRQControl capability. This gives you the authority to make this call.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	depth	Number of bits of index to resolve to find the destination slot.
seL4_Word	pci_bus	PCI bus ID of the device that will generate the interrupt.
seL4_Word	pci_dev	PCI device ID of the device that will generate the interrupt.
seL4_Word	pci_func	PCI function ID of the device that will generate the interrupt.
seL4_Word	handle	Value of the handle programmed into the data portion of the MSI.
seL4_Word	vector	CPU vector to deliver the interrupt to.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1 and Section 8.2.1.

10.4.2 seL4_TCB

10.4.2.1 Set EPT Root

```
static inline int seL4_TCB_SetEPTRoot
```

Set the EPT root of a thread

Type	Name	Description
seL4_TCB	_service	TODO
seL4_X86_EPTPML4	eptpml4	CPTR to an EPT PML4 object to act as the guest mode vspace root

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 6.4.

10.4.3 seL4_X86_ASIDControl

10.4.3.1 Make Pool

```
static inline int seL4_X86_ASIDControl_MakePool
```

Create an X86 ASID pool.

Type	Name	Description
seL4_X86_ASIDControl	<code>_service</code>	The master ASIDControl capability.
seL4_Untyped	<code>untyped</code>	Capability to an untyped memory object that will become the pool. Must be 4K bytes.
seL4_CNode	<code>root</code>	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	<code>index</code>	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	<code>depth</code>	Number of bits of index to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Together with a capability to `Untyped Memory`, which is passed as an argument, create an `ASID Pool`. The untyped capability must represent a 4K memory object. This will create an ASID pool with enough space for 1024 VSpaces.

10.4.4 seL4_X86_ASIDPool

10.4.4.1 Assign

```
static inline int seL4_X86_ASIDPool_Assign
```

Assign an ASID pool.

Type	Name	Description
seL4_X86_ASIDPool	_service	The ASID pool which is being assigned to. Must not be full. Each ASID pool can contain 1024 entries.
seL4_CPtr	vspace	The page directory that is being assigned to an ASID pool. Must not already be assigned to an ASID pool.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Assigns an ASID to the VSpace associated with the Page Directory passed in as an argument.

10.4.5 seL4_X86_EPTPD

10.4.5.1 Map

```
static inline int seL4_X86_EPTPD_Map
```

Map an EPT page directory.

Type	Name	Description
seL4_X86_EPTPD	_service	Capability to the EPT PD being operated on.
seL4_X86_EPTPML4	eptpm14	Capability to the EPT root which will contain the mapping
seL4_Word	gpa	Guest physical address to map the page into.
seL4_X86_VMAttributes	attr	VM attributes for the mapping. Possible values for this type are given in Chapter 7

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.5.2 Unmap

```
static inline int seL4_X86_EPTPD_Unmap
```

Unmap an EPT page directory.

Type	Name	Description
seL4_X86_EPTPD	_service	Capability to the EPT PD being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.6 seL4_X86_EPTPDPT

10.4.6.1 Map

```
static inline int seL4_X86_EPTPDPT_Map
```

Map an EPT page directory page table.

Type	Name	Description
seL4_X86_EPTPDPT	_service	Capability to the EPT PDPT being operated on.
seL4_X86_EPTPML4	eptpm14	Capability to the EPT root which will contain the mapping
seL4_Word	gpa	Guest physical address to map the page into.
seL4_X86_VMAttributes	attr	VM attributes for the mapping. Possible values for this type are given in Chapter 7

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.6.2 Unmap

```
static inline int seL4_X86_EPTPDPT_Unmap
```

Unmap an EPT page directory page table.

Type	Name	Description
seL4_X86_EPTPDPT	_service	Capability to the EPT PDPT being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.7 seL4_X86_EPTPT

10.4.7.1 Map

```
static inline int seL4_X86_EPTPT_Map
```

Map an EPT page table.

Type	Name	Description
seL4_X86_EPTPT	_service	Capability to the EPT PT being operated on.
seL4_X86_EPTPML4	eptpm14	Capability to the EPT root which will contain the mapping
seL4_Word	gpa	Guest physical address to map the page into.
seL4_X86_VMAttributes	attr	VM attributes for the mapping. Possible values for this type are given in Chapter 7

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.7.2 Unmap

```
static inline int seL4_X86_EPTPT_Unmap
```

Unmap an EPT page table.

Type	Name	Description
seL4_X86_EPTPT	_service	Capability to the EPT PT being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.8 seL4_X86_IOPageTable

10.4.8.1 Map

```
static inline int seL4_X86_IOPageTable_Map
```

Map an IO page table into an IOSpace.

Type	Name	Description
seL4_X86_IOPageTable	_service	Capability to the I/O page table being operated on.
seL4_X86_IOSpace	iospace	The IOSpace to map the page table into.
seL4_Word	ioaddr	The address to map the page table at.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.2.3

10.4.8.2 Unmap

```
static inline int seL4_X86_IOPageTable_Unmap
```

Unmap an IO page table from an IOSpace.

Type	Name	Description
seL4_X86_IOPageTable	_service	Capability to the I/O page table being operated on. The page table to unmap.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.2.3

10.4.9 seL4_X86_IOPort

10.4.9.1 In16

```
static inline seL4_X86_IOPort_In16_t seL4_X86_IOPort_In16
```

Read 16 bits from an IO port.

Type	Name	Description
seL4_X86_IOPort	_service	An I/O Port capability.
seL4_Uint16	port	The port to read from.

Return value: A `seL4_X86_IOPort_In16_t` structure as described in Section 8.2.2.

Description: See Section 8.2.2

10.4.9.2 In32

```
static inline seL4_X86_IOPort_In32_t seL4_X86_IOPort_In32
```

Read 32 bits from an IO port.

Type	Name	Description
seL4_X86_IOPort	_service	An I/O Port capability.
seL4_Uint16	port	The port to read from.

Return value: A `seL4_X86_IOPort_In32_t` structure as described in Section 8.2.2.

Description: See Section 8.2.2

10.4.9.3 In8

```
static inline seL4_X86_IOPort_In8_t seL4_X86_IOPort_In8
```

Read 8 bits from an IO port.

Type	Name	Description
seL4_X86_IOPort	_service	An I/O Port capability.
seL4_Uint16	port	The port to read from.

Return value: A `seL4_X86_IOPort_In8_t` structure as described in Section 8.2.2.

Description: See Section 8.2.2

10.4.9.4 Out16

```
static inline int sel4_X86_IOPort_Out16
```

Write 16 bits to an IO port.

Type	Name	Description
sel4_X86_IOPort	_service	An I/O Port capability.
sel4_Word	port	The port to write to.
sel4_Word	data	Data to write to the IO port.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.2.2

10.4.9.5 Out32

```
static inline int sel4_X86_IOPort_Out32
```

Write 32 bits to an IO port.

Type	Name	Description
sel4_X86_IOPort	_service	An I/O Port capability.
sel4_Word	port	The port to write to.
sel4_Word	data	Data to write to the IO port.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.2.2

10.4.9.6 Out8

```
static inline int seL4_X86_IOPort_Out8
```

Write 8 bits to an IO port.

Type	Name	Description
seL4_X86_IOPort	_service	An I/O Port capability.
seL4_Word	port	The port to write to.
seL4_Word	data	Data to write to the IO port.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.2.2

10.4.10 seL4_X86_IOPortControl

10.4.10.1 Issue

```
static inline int seL4_X86_IOPortControl_Issue
```

Issue an IO port sub range.

Type	Name	Description
seL4_X86_IOPortControl	_service	Control capability for I/O ports.
seL4_Word	first_port	First port of the range of the issued capability.
seL4_Word	last_port	Last port of the range of the issued capability.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	depth	Number of bits of dest_index to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.2.2

10.4.11 seL4_X86_Page

10.4.11.1 Get Address

```
static inline seL4_X86_Page_GetAddress_t seL4_X86_Page_GetAddress
```

Get the physical address of the underlying frame.

Type	Name	Description
seL4_X86_Page	_service	Capability to the page being operated on.

Return value: A `seL4_IA32_Page_GetAddress_t` struct that contains a `seL4_Word` `paddr`, which holds the physical address of the page, and `int` `error`. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.11.2 Map EPT

```
static inline int seL4_X86_Page_MapEPT
```

TODO

Type	Name	Description
seL4_X86_Page	_service	Capability to the page being operated on.
seL4_X86_EPTPML4	vspace	TODO
seL4_Word	vaddr	TODO
seL4_CapRights_t	rights	TODO
seL4_X86_VMAAttributes	attr	TODO

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.4.11.3 Map I/O

```
static inline int seL4_X86_Page_MapIO
```

Map a page into an IOSpace.

Type	Name	Description
seL4_X86_Page	_service	Capability to the page being operated on.
seL4_X86_IOSpace	iospace	The IOSpace that the frame is being mapped into
seL4_CapRights_t	rights	Rights for the mapping. Possible values for this type are given in Section 3.1.4
seL4_Word	ioaddr	The address that the frame is being mapped at.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.11.4 Map

```
static inline int seL4_X86_Page_Map
```

Map a page into an address space or update the mapping attributes.

Type	Name	Description
seL4_X86_Page	_service	Capability to the page being operated on.
seL4_CPtr	vspace	Capability to the VSpace which will contain the mapping
seL4_Word	vaddr	Virtual address to map the page into.
seL4_CapRights_t	rights	Rights for the mapping. Possible values for this type are given in Section 3.1.4
seL4_X86_VMAAttributes	attr	VM attributes for the mapping. Possible values for this type are given in Chapter 7

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Takes a VSpace capability, as an argument and installs a reference to the given Page in the lowest-level unmapped paging structure corresponding to the given address, or updates the mapping attributes if the page is already mapped at this address. If the required paging structures are not present this operation will fail, returning a seL4_FailedLookup error.

10.4.11.5 Unmap

```
static inline int seL4_X86_Page_Unmap
```

Unmap a page.

Type	Name	Description
seL4_X86_Page	_service	Capability to the page being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Removes an existing mapping.

10.4.12 seL4_X86_PageDirectory

10.4.12.1 Get Status Bits

```
static inline seL4_X86_PageDirectory_GetStatusBits_t seL4_X86_PageDirectory_GetStatusBits
```

Retrieve the accessed and dirty bits of a page mapped into an address space.

Type	Name	Description
seL4_X86_PageDirectory	_service	Capability to the page directory being operated on. Capability to the address space to query.
seL4_Word	vaddr	Virtual address of the page to query

Return value: A `seL4_X86_PageDirectory_GetStatusBits_t` structure.

Description: See Chapter 7

10.4.12.2 Map

```
static inline int seL4_X86_PageDirectory_Map
```

Map a page directory.

Type	Name	Description
seL4_X86_PageDirectory	_service	Capability to the page directory being operated on.
seL4_CPtr	vspace	Capability to the VSpace which will contain the mapping
seL4_Word	vaddr	Virtual address to map the page into.
seL4_X86_VMAttributes	attr	VM attributes for the mapping. Possible values for this type are given in Chapter 7

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.12.3 Unmap

```
static inline int sel4_X86_PageDirectory_Unmap
```

Unmap a page directory.

Type	Name	Description
sel4_X86_PageDirectory	_service	Capability to the page directory being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

10.4.13 sel4_X86_PageTable

10.4.13.1 Map

```
static inline int sel4_X86_PageTable_Map
```

Map a page table into an address space.

Type	Name	Description
sel4_X86_PageTable	_service	Capability to the page table being operated on.
sel4_CPtr	vspace	Capability to the VSpace which will contain the mapping
sel4_Word	vaddr	Virtual address to map the page into.
sel4_X86_VMAttributes	attr	VM attributes for the mapping. Possible values for this type are given in Chapter 7

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Takes a `PageDirectory` capability as an argument, and installs a reference to the invoked `PageTable` in a specified slot in the `PageDirectory`.

10.4.13.2 Unmap

```
static inline int seL4_X86_PageTable_Unmap
```

Unmap a page table from its address space and zero it out.

Type	Name	Description
seL4_X86_PageTable	_service	Capability to the page table being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Removes the reference to the invoked `PageTable` from its containing `PageDirectory`. See Chapter 7

10.4.14 seL4_X86_VCPU

10.4.14.1 Disable IO Port

```
static inline int seL4_X86_VCPU_DisableIOPort
```

Disable I/O port range in privileged execution

Type	Name	Description
seL4_X86_VCPU	_service	VCPU object to operate on
seL4_Word	low	Start of the I/O port range to disable
seL4_Word	high	Last I/O port in the range to disable

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Disable a range of I/O ports for direct access by the execution mode in the VCPU.

10.4.14.2 Enable IO Port

```
static inline int seL4_X86_VCPU_EnableIOPort
```

Enable I/O port range in guest execution

Type	Name	Description
seL4_X86_VCPU	_service	VCPU object to operate on
seL4_X86_IOPort	ioPort	I/O port capability whose authority is being delegating
seL4_Word	low	Start of the I/O port range to enable
seL4_Word	high	Last I/O port in the range to enable

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Enables a range of I/O ports for direct access by the execution mode in the VCPU. The requested port range must be a sub range of the provided I/O port capability.

This also establishes a link between the provided I/O port capability and the VCPU, see Section 6.4 for details.

10.4.14.3 Read VMCS

```
static inline seL4_X86_VCPU_ReadVMCS_t seL4_X86_VCPU_ReadVMCS
```

Read VMCS field from the hardware

Type	Name	Description
seL4_X86_VCPU	_service	VCPU object to operate on
seL4_Word	field	Field to give to vmread instruction

Return value: A `seL4_X86_VCPU_ReadVMCS_t` struct that contains a `seL4_Word` value, which holds the return result of the `vmread` instruction, and `int error`. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Thin wrapper around the `vmread` instruction that is performed on the VMCS region that is part of the VCPU object. After validating that a legal field is requested the value of 'vmread' is returned directly in the result.

10.4.14.4 Set TCB

```
static inline int seL4_X86_VCPU_SetTCB
```

Bind TCB to VCPU

Type	Name	Description
seL4_X86_VCPU	_service	VCPU object to operate on
seL4_TCB	tcb	CPTR of the TCB to bind to

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Configures the one-to-one binding of a VCPU and TCB, overwriting any previous binding in both. See Section 6.4.

10.4.14.5 Write Registers

```
static inline int seL4_X86_VCPU_WriteRegisters
```

Set guest mode registers to the fields of a given `seL4_VCPUContext`

Type	Name	Description
seL4_X86_VCPU	_service	VCPU object to operate on
seL4_VCPUContext *	regs	Data structure containing the new register values.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Sets the guest mode registers, which is any registers not already part of the VMCS.

10.4.14.6 Write VMCS

```
static inline seL4_X86_VCPU_WriteVMCS_t seL4_X86_VCPU_WriteVMCS
```

Write VMCS field to the hardware

Type	Name	Description
seL4_X86_VCPU	_service	VCPU object to operate on
seL4_Word	field	Field to give to <code>vmwrite</code> instruction
seL4_Word	value	Value to write using <code>vmwrite</code> instruction

Return value: A `seL4_X86_VCPU_WriteVMCS_t` struct that contains a `seL4_Word` `written`, which holds the final value written with the `vmwrite` instruction, and `int error`. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Thin wrapper around the ‘`vmwrite`’ instruction that is performed on the VMCS region that is part of the VCPU object. As well as validating that a legal field is requested, the value may be modified to ensure any bits that are fixed in the hardware are correct, and that any features required for kernel correctness are not disabled (see Section 6.4).

The final value written to the hardware is returned and can be compared to the input parameter to determine what bits the kernel changed.

10.5 IA32-Specific Object Methods

No methods.

10.6 x86_64-Specific Object Methods

10.6.1 seL4_X86_PDPT

10.6.1.1 Map

```
static inline int seL4_X86_PDPT_Map
```

TODO

Type	Name	Description
seL4_X86_PDPT	_service	TODO
seL4_X64_PML4	pml4	TODO
seL4_Word	vaddr	TODO
seL4_X86_VMAttributes	attr	TODO

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.6.1.2 Unmap

```
static inline int seL4_X86_PDPT_Unmap
```

TODO

Type	Name	Description
seL4_X86_PDPT	_service	TODO

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.7 ARM-Specific Object Methods

10.7.1 seL4_ARM_ASIDControl

10.7.1.1 Make Pool

```
static inline int seL4_ARM_ASIDControl_MakePool
```

Create an ASID Pool.

Type	Name	Description
seL4_ARM_ASIDControl	_service	The master ASIDControl capability being operated on.
seL4_Untyped	untyped	Capability to an untyped memory object that will become the pool. Must be 4K bytes.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	depth	Number of bits of index to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Together with a capability to **Untyped Memory**, which is passed as an argument, create an **ASID Pool**. The untyped capability must represent a 4K memory object. This will create an ASID pool with enough space for 1024 VSpaces.

10.7.2 seL4_ARM_ASIDPool

10.7.2.1 Asid Pool Assign

```
static inline int seL4_ARM_ASIDPool_Assign
```

Assign an ASID Pool.

Type	Name	Description
seL4_ARM_ASIDPool	_service	The ASID pool which is being assigned to. Must not be full. Each ASID pool can contain 1024 entries.
seL4_CPtr	vspace	The VSpace that is being assigned to an ASID pool. Must not already be assigned to an ASID pool.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Assigns an ASID to the VSpace passed in as an argument.

10.7.3 seL4_ARM_CB

10.7.3.1 AssignVspace

```
static inline int seL4_ARM_CB_AssignVspace
```

Assigning a vspace to a context bank.

Type	Name	Description
seL4_ARM_CB	_service	A CB capability. This gives you the authority to make this call.
seL4_CPtr	vspace	The VSpace that is being assigned to a context bank. Must already has an assigned ASID.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.3.

10.7.3.2 CBClearFault

```
static inline int seL4_ARM_CB_CBClearFault
```

Clear the fault status of the context bank.

Type	Name	Description
seL4_ARM_CB	_service	A CB capability. This gives you the authority to make this call.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.7.

10.7.3.3 CBGetFault

```
static inline seL4_ARM_CB_CBGetFault_t seL4_ARM_CB_CBGetFault
```

Get the fault status of the context bank.

Type	Name	Description
seL4_ARM_CB	_service	A CB capability. This gives you the authority to make this call.

Return value: A `seL4_ARM_SMMU_CB_GetFault_t` struct that contains a `seL4_Word status`, which holds the fault status of the context bank, `seL4_Word address`, which holds the faulty address, and `int error`. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.7.

10.7.3.4 TLBInvalidate

```
static inline int seL4_ARM_CB_TLBInvalidate
```

Invalidating TLB entries used by the current ASID in this context bank.

Type	Name	Description
seL4_ARM_CB	_service	A CB capability. This gives you the authority to make this call.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.6.

10.7.3.5 UnassignVspace

```
static inline int seL4_ARM_CB_UnassignVspace
```

Unassigning a vspace to a context bank.

Type	Name	Description
seL4_ARM_CB	_service	A CB capability. This gives you the authority to make this call.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.3.

10.7.4 seL4_ARM_CBControl

10.7.4.1 GetCB

```
static inline int seL4_ARM_CBControl_GetCB
```

Create a CB capability.

Type	Name	Description
seL4_ARM_CBControl	_service	A CBControl capability. This gives you the authority to make this call.
seL4_Word	cb	The CB that you want this capability to manage.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	depth	Number of bits of dest_index to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.2.

10.7.4.2 TLBInvalidateAll

```
static inline int seL4_ARM_CBControl_TLBInvalidateAll
```

Invalidate all TLB entries.

Type	Name	Description
seL4_ARM_CBControl	_service	A CBControl capability. This gives you the authority to make this call.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.6.

10.7.5 seL4_ARM_IOPageTable

10.7.5.1 Map

```
static inline int seL4_ARM_IOPageTable_Map
```

TODO

Type	Name	Description
seL4_ARM_IOPageTable	_service	Capability to the I/O page table being operated on.
seL4_ARM_IOSpace	iospace	TODO
seL4_Word	ioaddr	TODO

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.7.5.2 Unmap

```
static inline int seL4_ARM_IOPageTable_Unmap
```

TODO

Type	Name	Description
seL4_ARM_IOPageTable	_service	Capability to the I/O page table being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.7.6 seL4_ARM_Page

10.7.6.1 Clean Data

```
static inline int seL4_ARM_Page_Clean_Data
```

Cleans the data cache out to RAM. The start and end are relative to the page being serviced.

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.
seL4_Word	start_offset	The offset, relative to the start of the page inclusive.
seL4_Word	end_offset	The offset, relative to the start of the page exclusive.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.7.6.2 Clean and Invalidate Data

```
static inline int seL4_ARM_Page_CleanInvalidate_Data
```

Clean and invalidates the cache range within the given page. The range will be flushed out to RAM. The start and end are relative to the page being serviced.

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.
seL4_Word	start_offset	The offset, relative to the start of the page inclusive.
seL4_Word	end_offset	The offset, relative to the start of the page exclusive.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.7.6.3 Get Address

```
static inline seL4_ARM_Page_GetAddress_t seL4_ARM_Page_GetAddress
```

Get the physical address of the underlying frame.

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.

Return value: A `seL4_ARM_Page_GetAddress_t` struct that contains a `seL4_Word paddr`, which holds the physical address of the page, and `int error`. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.7.6.4 Invalidate Data

```
static inline int seL4_ARM_Page_Invalidate_Data
```

Invalidates the cache range within the given page. The start and end are relative to the page being serviced and should be aligned to a cache line boundary where possible. An additional clean is performed on the outer cache lines if the start and end are not aligned, to clean out the bytes between the requested and the cache line boundary.

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.
seL4_Word	start_offset	The offset, relative to the start of the page inclusive.
seL4_Word	end_offset	The offset, relative to the start of the page exclusive.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.7.6.5 Map I/O

```
static inline int seL4_ARM_Page_MapIO
```

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.
seL4_ARM_IOSpace	iospace	TODO
seL4_CapRights_t	rights	TODO
seL4_Word	ioaddr	TODO

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.7.6.6 Map

```
static inline int seL4_ARM_Page_Map
```

Map a page into an address space or update the mapping attributes.

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.
seL4_CPtr	vspace	Capability to the VSpace which will contain the mapping. Must be assigned to an ASID pool.
seL4_Word	vaddr	Virtual address to map the page into.
seL4_CapRights_t	rights	Rights for the mapping. Possible values for this type are given in Section 3.1.4 .
seL4_ARM_VMAttributes	attr	VM Attributes for the mapping. Possible values for this type are given in Chapter 7 .

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Takes a VSpace capability, as an argument and installs a reference to the given Page in the lowest-level unmapped paging structure corresponding to the given address, or updates the mapping attributes if the page is already mapped at this address. If the required paging structures are not present this operation will fail, returning a seL4_FailedLookup error. The page must not already be mapped through this capability in a different VSpace or at a different address; the page may be mapped in multiple VSpaces by copying the capability.

10.7.6.7 Unify Instruction

```
static inline int seL4_ARM_Page_Unify_Instruction
```

Unify Instruction Cache. Cleans data lines to point of unification, invalidate corresponding instruction lines to point of unification, then invalidates branch predictors. The start and end are relative to the page being serviced.

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.
seL4_Word	start_offset	The offset, relative to the start of the page inclusive.
seL4_Word	end_offset	The offset, relative to the start of the page exclusive.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.7.6.8 Unmap

```
static inline int seL4_ARM_Page_Unmap
```

Unmap a page.

Type	Name	Description
seL4_ARM_Page	_service	Capability to the page being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Removes an existing mapping.

10.7.7 seL4_ARM_PageTable

10.7.7.1 Map

```
static inline int seL4_ARM_PageTable_Map
```

Map a page table into an address space.

Type	Name	Description
seL4_ARM_PageTable	_service	Capability to the page table being operated on.
seL4_CPtr	vspace	Capability to the VSpace which will contain the mapping. Must be assigned to an ASID pool.
seL4_Word	vaddr	Virtual address to map the page into.
seL4_ARM_VMAttributes	attr	VM Attributes for the mapping. Possible values for this type are given in Chapter 7 .

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Takes a VSpace capability as an argument, and installs a reference to the invoked PageTable in the VSpace according to the provided virtual address. If higher-level paging structures do not exist for the virtual address this operation will fail, returning a seL4_FailedLookup error. If the page table is already mapped this operation will fail, returning a seL4_InvalidCapability error. If another object is already mapped at this level for the virtual address this operation will fail, returning a seL4_DeleteFirst error.

10.7.7.2 Unmap

```
static inline int seL4_ARM_PageTable_Unmap
```

Unmap a page table from its Page Directory and zero it out.

Type	Name	Description
seL4_ARM_PageTable	_service	Capability to the page table being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Removes the reference to the invoked Page Table from its containing Page Directory.

10.7.8 seL4_ARM_SID

10.7.8.1 BindCB

```
static inline int seL4_ARM_SID_BindCB
```

Binding a context bank to a stream ID.

Type	Name	Description
seL4_ARM_SID	_service	A SID capability. This gives you the authority to make this call.
seL4_CPtr	cb	The CB that is being binded to a stream ID. Must already has an assigned vspace.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.4.

10.7.8.2 UnbindCB

```
static inline int seL4_ARM_SID_UnbindCB
```

Unbinding a context bank from a stream ID.

Type	Name	Description
seL4_ARM_SID	_service	A SID capability. This gives you the authority to make this call.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.4.

10.7.9 seL4_ARM_SIDControl

10.7.9.1 ClearFault

```
static inline int seL4_ARM_SIDControl_ClearFault
```

Clear the fault status of the SMMU.

Type	Name	Description
seL4_ARM_SIDControl	_service	A SIDControl capability. This gives you the authority to make this call.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.7.

10.7.9.2 GetFault

```
static inline seL4_ARM_SIDControl_GetFault_t seL4_ARM_SIDControl_GetFault
```

Get the fault status of the SMMU.

Type	Name	Description
seL4_ARM_SIDControl	_service	A SIDControl capability. This gives you the authority to make this call.

Return value: A `seL4_ARM_SMMU_GetFault_t` struct that contains a `seL4_Word status`, which holds the global fault status of the SMMU, `seL4_Word syndrome_0`, which holds the global fault syndrome 0 of the SMMU, `seL4_Word syndrome_1`, which holds the global fault syndrome 1 of the SMMU, and `int error`. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.7.

10.7.9.3 GetSID

```
static inline int seL4_ARM_SIDControl_GetSID
```

Create a SID capability.

Type	Name	Description
seL4_ARM_SIDControl	_service	A SIDControl capability. This gives you the authority to make this call.
seL4_Word	sid	The SID that you want this capability to manage.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	depth	Number of bits of <code>dest_index</code> to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.3.1.1.

10.7.10 seL4_ARM_VCPU

10.7.10.1 Acknowledge Virtual PPI IRQ

```
static inline int seL4_ARM_VCPU_AckVPPI
```

Acknowledge a PPI IRQ previously forwarded from a VPPIEvent fault

Type	Name	Description
seL4_ARM_VCPU	_service	Capability to the VCPU being operated on.
seL4_Word	irq	irq to ack.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Acknowledge and unmask the PPI interrupt so that further interrupts can be forwarded through VPPIEvent faults.

10.7.10.2 Inject IRQ

```
static inline int seL4_ARM_VCPU_InjectIRQ
```

Inject an IRQ to a virtual CPU

Type	Name	Description
seL4_ARM_VCPU	_service	Capability to the VCPU being operated on.
seL4_Uint16	virq	Virtual IRQ ID
seL4_Uint8	priority	Priority of the IRQ to be injected
seL4_Uint8	group	IRQ group
seL4_Uint8	index	IRQ index

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.7.10.3 Read Registers

```
static inline seL4_ARM_VCPU_ReadRegs_t seL4_ARM_VCPU_ReadRegs
```

Read a virtual CPU register

Type	Name	Description
seL4_ARM_VCPU	_service	Capability to the VCPU being operated on.
seL4_Word	field	Register to read from a VCPU

Return value: TODO

Description: TODO

10.7.10.4 Set TCB

```
static inline int seL4_ARM_VCPU_SetTCB
```

Bind a TCB to a virtual CPU

Type	Name	Description
seL4_ARM_VCPU	_service	Capability to the VCPU being operated on.
seL4_TCB	tcb	Capability to TCB to bind to a virtual CPU

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: There is a 1:1 relationship between a virtual CPU and a TCB. If either (or both) of them is associated with another one, they will be dissociated, and then associated to the ones called in this system calls.

10.7.10.5 Write Registers

```
static inline int seL4_ARM_VCPU_WriteRegs
```

Write a virtual CPU register

Type	Name	Description
seL4_ARM_VCPU	_service	Capability to the VCPU being operated on.
seL4_Word	field	Register ID to write to a VCPU
seL4_Word	value	Value to be written to the VCPU register

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.7.11 seL4_IRQControl

10.7.11.1 GetTrigger

```
static inline int seL4_IRQControl_GetTrigger
```

Create an IRQ handler capability and specify the trigger method (edge or level).

Type	Name	Description
seL4_IRQControl	_service	An IRQControl capability. This gives you the authority to make this call.
seL4_Word	irq	The IRQ that you want this capability to handle.
seL4_Word	trigger	Indicates whether this IRQ is edge (1) or level (0) triggered.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
seL4_Word	index	CPTR to the destination slot. Resolved from the root of the destination CSpace.
seL4_Uint8	depth	Number of bits of dest_index to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1.

See Section 8.1.

10.7.11.2 GetTriggerCore

```
static inline int sel4_IRQControl_GetTriggerCore
```

Create an IRQ handler capability and specify the trigger method (edge or level) and the target core.

Type	Name	Description
sel4_IRQControl	<code>_service</code>	An IRQControl capability. This gives you the authority to make this call.
sel4_Word	<code>irq</code>	The IRQ that you want this capability to handle.
sel4_Word	<code>trigger</code>	Indicates whether this IRQ is edge (1) or level (0) triggered.
sel4_CNode	<code>root</code>	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth equivalent to the wordsize.
sel4_Word	<code>index</code>	CPTR to the destination slot. Resolved from the root of the destination CSpace.
sel4_Uint8	<code>depth</code>	Number of bits of <code>dest_index</code> to resolve to find the destination slot.
sel4_Word	<code>target</code>	Indicates the target core ID to which this irq will be sent.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Section 8.1.

10.8 Aarch32-Specific Object Methods

10.8.1 seL4_ARM_PageDirectory

10.8.1.1 Clean Data

```
static inline int seL4_ARM_PageDirectory_Clean_Data
```

Clean cached pages within a page directory

Type	Name	Description
seL4_ARM_PageDirectory	_service	Capability to the page directory being operated on.
seL4_Word	start	Start address
seL4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.8.1.2 Clean and Invalidate Data

```
static inline int seL4_ARM_PageDirectory_CleanInvalidate_Data
```

Clean and invalidate cached pages within a page directory

Type	Name	Description
seL4_ARM_PageDirectory	_service	Capability to the page directory being operated on.
seL4_Word	start	Start address
seL4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.8.1.3 Invalidate Data

```
static inline int seL4_ARM_PageDirectory_Invalidate_Data
```

Invalidate cached pages within a page directory

Type	Name	Description
seL4_ARM_PageDirectory	_service	Capability to the page directory being operated on.
seL4_Word	start	Start address
seL4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.8.1.4 Unify Instruction

```
static inline int seL4_ARM_PageDirectory_Unify_Instruction
```

Clean and invalidate cached instruction pages to point of unification

Type	Name	Description
seL4_ARM_PageDirectory	_service	Capability to the page directory being operated on.
seL4_Word	start	Start address
seL4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.9 Aarch64-Specific Object Methods

10.9.1 seL4_ARM_PageDirectory

10.9.1.1 Map

```
static inline int seL4_ARM_PageDirectory_Map
```

Map a page directory

Type	Name	Description
seL4_ARM_PageDirectory	_service	Capability to the page directory being operated on.
seL4_CPtr	vspace	Top level translation table. Must be assigned to an ASID pool.
seL4_Word	vaddr	Virtual address
seL4_ARM_VMAttributes	attr	VM Attributes for the mapping. Possible values for this type are given in Chapter 7 .

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Map a page directory (level 2) to an upper page directory (level 1). If an upper page directory does not exist for the virtual address this operation will fail, returning a seL4_FailedLookup error. If the page directory is already mapped this operation will fail, returning a seL4_InvalidCapability error. If another object is already mapped at this level for the virtual address this operation will fail, returning a seL4_DeleteFirst error.

10.9.1.2 Unmap

```
static inline int seL4_ARM_PageDirectory_Unmap
```

Unmap a page directory

Type	Name	Description
seL4_ARM_PageDirectory	_service	Capability to the page directory being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Unmap a page directory (level 2) from an upper page directory (level 1)

10.9.2 seL4_ARM_PageUpperDirectory

10.9.2.1 Map

```
static inline int seL4_ARM_PageUpperDirectory_Map
```

Map an upper page directory

Type	Name	Description
seL4_ARM_PageUpperDirectory	_service	Capability to the upper page directory being operated on.
seL4_CPtr	vspace	Top level translation table. Must be assigned to an ASID pool.
seL4_Word	vaddr	Virtual address
seL4_ARM_VMAttributes	attr	VM Attributes for the mapping. Possible values for this type are given in Chapter 7 .

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Map an upper page directory (level 1) to a top level translation table (level 0). If the upper page directory is already mapped this operation will fail, returning a seL4_InvalidCapability error. If another object is already mapped at this level for the virtual address this operation will fail, returning a seL4_DeleteFirst error.

10.9.2.2 Unmap

```
static inline int seL4_ARM_PageUpperDirectory_Unmap
```

TODO

Type	Name	Description
seL4_ARM_PageUpperDirectory	_service	Capability to the upper page directory being operated on.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: TODO

10.9.3 seL4_ARM_VSpace

10.9.3.1 Clean Data

```
static inline int seL4_ARM_VSpace_Clean_Data
```

Clean cached pages within a top level translation table

Type	Name	Description
seL4_ARM_VSpace	_service	Capability to the top level translation table being operated on.
seL4_Word	start	Start address
seL4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.9.3.2 Clean and Invalidate Data

```
static inline int seL4_ARM_VSpace_CleanInvalidate_Data
```

Clean and invalidate cached pages within a top level translation table

Type	Name	Description
seL4_ARM_VSpace	_service	Capability to the top level translation table being operated on.
seL4_Word	start	Start address
seL4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.9.3.3 Invalidate Data

```
static inline int sel4_ARM_VSpace_Invalidate_Data
```

Invalidate cached pages within a top level translation table

Type	Name	Description
sel4_ARM_VSpace	_service	Capability to the top level translation table being operated on.
sel4_Word	start	Start address
sel4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.9.3.4 Unify Instruction

```
static inline int sel4_ARM_VSpace_Unify_Instruction
```

Clean and invalidate cached instruction pages to point of unification

Type	Name	Description
sel4_ARM_VSpace	_service	Capability to the top level translation table being operated on.
sel4_Word	start	Start address
sel4_Word	end	End address

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.10 RISC-V-Specific Object Methods

10.10.1 General RISC-V Object Methods

10.10.2 seL4_IRQControl

No methods.

10.10.3 seL4_RISC-V__ASIDControl

10.10.3.1 MakePool

```
static inline int seL4_RISC_V_ASIDControl_MakePool
```

Create an ASID Pool.

Type	Name	Description
seL4_RISC_V_ASIDControl	_service	The master ASIDControl capability to invoke.
seL4_Untyped	untyped	Capability to an untyped memory object that will become the pool. Must be 4K bytes.
seL4_CNode	root	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth of 32.
seL4_Word	index	CPTR to the CNode that forms the root of the destination CSpace. Must be at a depth of 32.
seL4_Uint8	depth	Number of bits of index to resolve to find the destination slot.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Together with a capability to **Untyped Memory**, which is passed as an argument, create an **ASID Pool**. The untyped capability must represent a 4K memory object. This will create an ASID pool with enough space for 1024 VSpaces.

10.10.4 seL4_RISCV_ASIDPool

10.10.4.1 Assign

```
static inline int seL4_RISCV_ASIDPool_Assign
```

Assign an ASID Pool.

Type	Name	Description
seL4_RISCV_ASIDPool	_service	The ASID Pool capability to invoke, which must be to an ASID pool that is not full.
seL4_CPtr	vspace	The top-level PageTable that is being assigned to an ASID pool. Must not already be assigned to an ASID pool.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Assigns an ASID to the VSpace passed in as an argument.

10.10.5 seL4_RISCV_Page

10.10.5.1 GetAddress

```
static inline seL4_RISCV_Page_GetAddress_t seL4_RISCV_Page_GetAddress
```

Get the physical address of a page.

Type	Name	Description
seL4_RISCV_Page	_service	Capability to the page to invoke.

Return value: A `seL4_RISCV_Page_GetAddress_t` struct that contains a `seL4_Word paddr`, which holds the physical address of the page, and `int error`. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7.

10.10.5.2 Map

```
static inline int seL4_RISCV_Page_Map
```

Map a page into a page table.

Type	Name	Description
seL4_RISCV_Page	_service	Capability to the page to invoke.
seL4_RISCV_PageTable	vspace	VSpace to map the page into.
seL4_Word	vaddr	Virtual address at which to map the page.
seL4_CapRights_t	rights	Rights for the mapping. Possible values for this type are given in Section 3.1.4.
seL4_RISCV_VMAttributes	attr	VM Attributes for the mapping. Possible values for this type are given in Chapter 7.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Takes a VSpace, or top-level Page Table, capability as an argument and installs a reference to the given Page in the page table slot corresponding to the given address. If a page is already mapped at the same virtual address, update the mapping attributes. If the required paging structures are not present this operation will fail, returning a seL4_FailedLookup error.

10.10.5.3 Unmap

```
static inline int seL4_RISCV_Page_Unmap
```

Unmap a page.

Type	Name	Description
seL4_RISCV_Page	_service	Capability to the page to invoke.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Removes an existing mapping.

10.10.6 seL4_RISCV_PageTable

10.10.6.1 Map

```
static inline int seL4_RISCV_PageTable_Map
```

Map a page table at a specific virtual address.

Type	Name	Description
seL4_RISCV_PageTable	_service	Capability to the page table to invoke.
seL4_RISCV_PageTable	vspace	VSpace to map the lower-level page table into.
seL4_Word	vaddr	Virtual address at which to map the page table.
seL4_RISCV_VMAttributes	attr	VM Attributes for the mapping. Possible values for this type are given in Chapter 7.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: Starting from the VSpace, map the page table object at any unpopulated level for the provided virtual address. If all paging structures and mappings are present for this virtual address, return an seL4_DeleteFirst error.

10.10.6.2 Unmap

```
static inline int seL4_RISCV_PageTable_Unmap
```

Unmap a page table.

Type	Name	Description
seL4_RISCV_PageTable	_service	Capability to the page table to invoke.

Return value: A return value of 0 indicates success. A non-zero value indicates that an error occurred. See Section 10.1 for a description of the message register and tag contents upon error.

Description: See Chapter 7

Bibliography

- [Boy09] Andrew Boyton. A verified shared capability model. In Gerwin Klein, Ralf Huuck, and Bastian Schlich, editors, *Proceedings of the 4th Workshop on Systems Software Verification*, volume 254 of *Electronic Notes in Computer Science*, pages 25–44, Aachen, Germany, October 2009. Elsevier.
- [BSC⁺11] Bernard Blackham, Yao Shi, Sudipta Chattopadhyay, Abhik Roychoudhury, and Gernot Heiser. Timing analysis of a protected operating system kernel. In *IEEE Real-Time Systems Symposium*, pages 339–348, Vienna, Austria, November 2011.
- [BSH12] Bernard Blackham, Yao Shi, and Gernot Heiser. Improving interrupt response time in a verifiable protected microkernel. In *EuroSys*, pages 323–336, Bern, Switzerland, April 2012.
- [CKS08] David Cock, Gerwin Klein, and Thomas Sewell. Secure microkernels, state monads and scalable refinement. In Otmane Ait Mohamed, César Muñoz, and Sofiène Tahar, editors, *Proceedings of the 21st International Conference on Theorem Proving in Higher Order Logics*, volume 5170 of *Lecture Notes in Computer Science*, pages 167–182, Montreal, Canada, August 2008. Springer-Verlag.
- [DEK⁺06] Philip Derrin, Kevin Elphinstone, Gerwin Klein, David Cock, and Manuel M. T. Chakravarty. Running the manual: An approach to high-assurance microkernel development. In *Proceedings of the ACM SIGPLAN Haskell Workshop*, Portland, OR, USA, September 2006.
- [EKE08] Dhammika Elkaduwe, Gerwin Klein, and Kevin Elphinstone. Verified protection model of the seL4 microkernel. In Jim Woodcock and Natarajan Shankar, editors, *Proceedings of Verified Software: Theories, Tools and Experiments 2008*, volume 5295 of *Lecture Notes in Computer Science*, pages 99–114, Toronto, Canada, October 2008. Springer-Verlag.
- [Int11] Intel Corporation. *Intel Virtualization Technology for Directed I/O — Architecture Specification*, February 2011. [http://download.intel.com/technology/computing/vptech/Intel\(r\)_VT_for_Direct_IO.pdf](http://download.intel.com/technology/computing/vptech/Intel(r)_VT_for_Direct_IO.pdf).
- [KEH⁺09] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, Thomas Sewell, Harvey Tuch, and Simon Winwood. seL4: Formal verification of an OS kernel. In *Proceedings of the 22nd ACM*

- Symposium on Operating Systems Principles*, pages 207–220, Big Sky, MT, USA, October 2009. ACM.
- [MMB⁺13] Toby Murray, Daniel Matichuk, Matthew Brassil, Peter Gammie, Timothy Bourke, Sean Seefried, Corey Lewis, Xin Gao, and Gerwin Klein. seL4: from general purpose to a proof of information flow enforcement. In *IEEE Symposium on Security & Privacy*, pages 415–429, San Francisco, CA, May 2013.
- [Pal09] Ameya Palande. Capability-based secure DMA in seL4. Masters thesis, Vrije Universiteit, Amsterdam, January 2009.
- [SA99] Tom Shanley and Don Anderson. *PCI System Architecture*. Mindshare, Inc., 1999.
- [SWG⁺11] Thomas Sewell, Simon Winwood, Peter Gammie, Toby Murray, June Andronick, and Gerwin Klein. seL4 enforces integrity. In Marko van Eekelen, Herman Geuvers, Julien Schmaltz, and Freek Wiedijk, editor, *Interactive Theorem Proving (ITP)*, pages 325–340, Nijmegen, The Netherlands, August 2011.
- [TKN07] Harvey Tuch, Gerwin Klein, and Michael Norrish. Types, bytes, and separation logic. In Martin Hofmann and Matthias Felleisen, editors, *Proceedings of the 34th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 97–108, Nice, France, January 2007. ACM.
- [WKS⁺09] Simon Winwood, Gerwin Klein, Thomas Sewell, June Andronick, David Cock, and Michael Norrish. Mind the gap: A verification framework for low-level C. In Stefan Berghofer, Tobias Nipkow, Christian Urban, and Makarius Wenzel, editors, *Proceedings of the 22nd International Conference on Theorem Proving in Higher Order Logics*, volume 5674 of *Lecture Notes in Computer Science*, pages 500–515, Munich, Germany, August 2009. Springer-Verlag.