

北京理工大学

本科生毕业设计（论文）外文翻译

From L3 to seL4 What Have We Learnt in 20 Years of L4
Microkernels?

外文原文题目：

从 L3 到 seL4：微内核二十年的演进带来的经验与启示

中文翻译题目：

ReL4 中基于硬件加速的异步系统调用的设计 与实现

**Design and Implementation of Asynchronous System Call based
on Hardware Acceleration in Rel4**

学 院： 计算机学院

专 业： 计算机科学与技术

班 级： 07112103

学生姓名： 王菁芃

学 号： 1120211759

指导教师： 陆慧梅

原文:

北京理工大学



From L3 to seL4

What Have We Learnt in 20 Years of L4 Microkernels?

Kevin Elphinstone and Gernot Heiser
NICTA and UNSW, Sydney
{kevin.elphinstone, gernot}@nicta.com.au

Abstract

The L4 microkernel has undergone 20 years of use and evolution. It has an active user and developer community, and there are commercial versions which are deployed on a large scale and in safety-critical systems. In this paper we examine the lessons learnt in those 20 years about microkernel design and implementation. We revisit the L4 design papers, and examine the evolution of design and implementation from the original L4 to the latest generation of L4 kernels, especially seL4, which has pushed the L4 model furthest and was the first OS kernel to undergo a complete formal verification of its implementation as well as a sound analysis of worst-case execution times. We demonstrate that while much has changed, the fundamental principles of minimality and high IPC performance remain the main drivers of design and implementation decisions.

1 Introduction

Twenty years ago, Liedtke [1993a] demonstrated with his L4 kernel that microkernel IPC could be fast, a factor 10–20 faster than other contemporary microkernels.

Microkernels minimize the functionality that is provided by the kernel: the kernel provides a set of general mechanisms, while user-mode servers implement the actual operating system (OS) services [Levin et al., 1975]. User code obtains a system service by communicating with servers via an inter-process communication (IPC) mechanism, typically message passing. Hence, IPC is

on the critical path of any service invocation, and low IPC costs are essential.

By the early '90s, IPC performance had become the achilles heel of microkernels: typical cost for a one-way message was around 100 μ s, which was too high for building performant systems, with a resulting trend to move core services back into the kernel [Condict et al., 1994]. There were also arguments that high IPC costs were an (inherent?) consequence of the structure of microkernel-based systems [Chen and Bershad, 1993].

In this context, the order-of-magnitude improvement of IPC costs Liedtke demonstrated was quite remarkable. It was followed by work discussing the philosophy and mechanisms of L4 [Liedtke, 1995, 1996], the demonstration of a para-virtualized Linux on L4 with only a few percent overhead [Härtig et al., 1997], the deployment of L4 kernels on billions of mobile devices and in safety-critical systems, and, finally, the world's first functional correctness proof of an OS kernel [Klein et al., 2009]. It also had a strong influence on other research systems, such as Barrelfish [Baumann et al., 2009].

In this paper we examine the development of L4 over the last 20 years. Specifically we look at what makes modern L4 kernels tick, how this relates to Liedtke's original design and implementation, and which of his microkernel “essentials” have passed the test of time. We specifically examine how the lessons of the past have influenced the design of the latest generation of L4 microkernels, exemplified by *seL4* [Klein et al., 2009], but point out where other current L4 versions have made different design decisions.

2 Background

2.1 The L4 Microkernel Family

L4 developed out of an earlier system, called L3, developed by Liedtke [1993b] in the early 1980s on i386 platforms. L3 was a complete OS with built-in persistence, and it already featured user-mode drivers, still a

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the Owner/Author.

Copyright is held by the Owner/Author(s).
SOSP'13, Nov. 3–6, 2013, Farmington, Pennsylvania, USA.
ACM 978-1-4503-2388-8/13/11.
<http://dx.doi.org/10.1145/2517349.2522720>

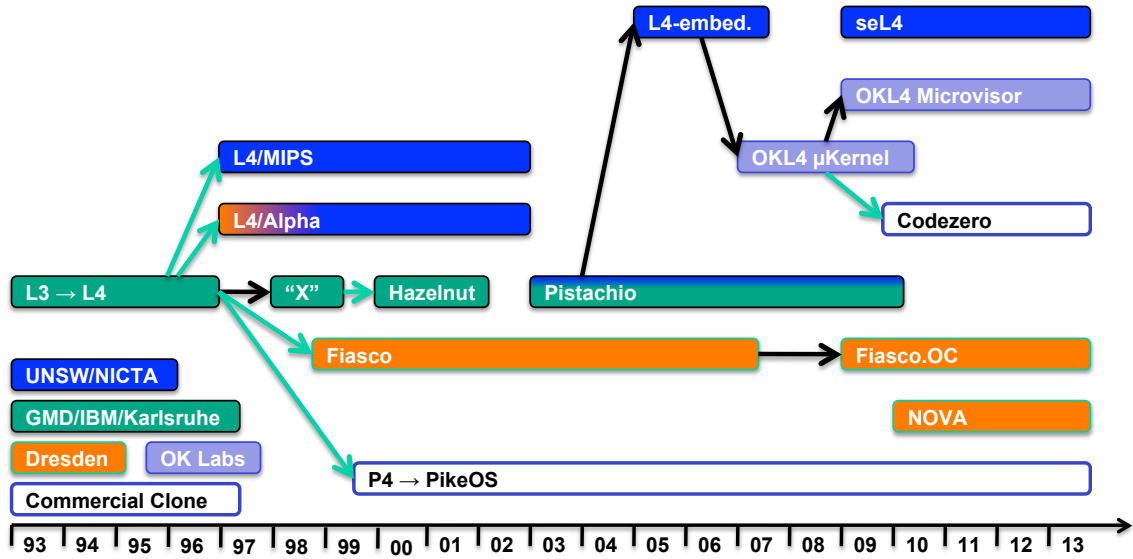


Figure 1: The L4 family tree (simplified). Black arrows indicate code, green arrows ABI inheritance. Box colours indicate origin as per key at the bottom left.

characteristic of L4 microkernels. It was commercially deployed in a few thousand installations (mainly schools and legal practices). Like all microkernels at the time, L3 suffered from IPC costs of the order of 100 μ s.

Liedtke initially used L3 to try out new ideas, and what he referred to as “L3” in early publications [Liedtke, 1993a] was actually an interim version of a radical re-design. He first used the name “L4” with the “V2” ABI circulated in the community from 1995.

In the following we refer to this version as the “original L4”. Liedtke implemented it completely in assembler on i486-based PCs and soon ported it to the Pentium.

This initial work triggered a twenty-year evolution, with multiple ABI revisions and from-scratch implementations, as depicted in Figure 1. It started with TU Dresden and UNSW re-implementing the ABI (with necessary adaptations) on 64-bit Alpha and MIPS processors, the latter implemented all longer-running operations in C. Both kernels achieved sub-microsecond IPC performance [Liedtke et al., 1997a] and were released as open source. The UNSW Alpha kernel was the first multiprocessor version of L4.

Liedtke, who had moved from GMD to IBM Watson, kept experimenting with the ABI in what became known as *Version X*. GMD and IBM imposed an IP regime which proved too restrictive for other researchers, prompting Dresden to implement a new x86 version from scratch, called *Fiasco* in reference to their experience in trying to deal with IP issues. The open-source *Fiasco* was the first L4 version written almost

completely in a higher-level language (C++) and is the oldest L4 codebase still actively maintained. It was the first L4 kernel with significant commercial use (estimated shipments up to 100,000).

After Liedtke’s move to Karlsruhe, he and his students did their own from-scratch implementation (in C), *Hazelnut*, which was the first L4 kernel that was ported (rather than re-implemented) to another architecture (from Pentium to ARM).

Karlsruhe’s experience with Version X and Hazelnut resulted in a major ABI revision, V4, aimed at improving kernel and application portability, multi-processor support and addressing various other shortcomings. After Liedtke’s tragic death in 2001, his students implemented the design in a new open-source kernel, *L4Ka::Pistachio* (“Pistachio” for short). It was written in C++, originally on x86 and PowerPC, and at UNSW/NICTA we soon after ported it to MIPS, Alpha, 64-bit PowerPC and ARM.¹ In most of these ports, less than 10% of the code changed.

At NICTA we soon re-targeted Pistachio for use in resource-constrained embedded systems, resulting in a fork called *NICTA::Pistachio-embedded* (“L4-embedded”). It saw massive-scale commercial deployment when Qualcomm adopted it as a protected-mode real-time OS for the firmware of their wireless modem processors. The NICTA spinout Open Kernel Labs took on the support and further development of the kernel, re-

¹There were also Itanium [Gray et al., 2005] and SPARC versions, but they were never completed.

Name	Year	Processor	MHz	Cycles	μs
Original	1993	i486	50	250	5.00
Original	1997	Pentium	160	121	0.75
L4/MIPS	1997	R4700	100	86	0.86
L4/Alpha	1997	21064	433	45	0.10
Hazelnut	2002	Pentium 4	1,400	2,000	1.38
Pistachio	2005	Itanium 2	1,500	36	0.02
OKL4	2007	XScale 255	400	151	0.64
NOVA	2010	Core i7 (Bloomfield) 32-bit	2,660	288	0.11
seL4	2013	Core i7 4770 (Haswell) 32-bit	3,400	301	0.09
seL4	2013	ARM11	532	188	0.35
seL4	2013	Cortex A9	1,000	316	0.32

Table 1: One-way IPC cost of various L4 kernels.

named the *OKL4 microkernel*.² Another deployed version is PikeOS, a commercial V2 clone by Sysgo, certified for use in safety-critical avionics and deployed in aircraft and trains.³

The influence of EROS [Shapiro et al., 1999] and increasing security focus resulted in a move to capabilities [Dennis and Van Horn, 1966] for access control, first with the 2.1 release of OKL4 (2008) and soon followed by Fiasco (renamed *Fiasco.OC*). Aiming for formal verification, which seemed infeasible for a code base not designed for the purpose, we instead opted for a from-scratch implementation for our capability-based *seL4* kernel.

The last few years also saw two new designs specifically aimed to support virtualization as the primary concern, *NOVA* [Steinberg and Kauer, 2010] from Dresden and the *OKL4 Microvisor* [Heiser and Leslie, 2010] from OK Labs.

A common thread throughout those two decades is the *minimality principle* introduced in Section 3.1, and a strong focus on the performance of the critical IPC operation; kernel implementors generally aim to stay close to the limits set by the micro-architecture, as shown in Table 1. Consequently, the L4 community tends to measure IPC latencies in cycles rather than microseconds, as this better relates to the hardware limits. In fact, Section 3.1 provides an interesting view of the context-switch-friendliness of the hardware (compare the cycle counts for Pentium 4 and Itanium, both from highly-optimised IPC implementations!)

²Total deployment is now in the billions, see Open Kernel Labs press release <http://www.ok-labs.com/releases/release/ok-labs-software-surpasses-milestone-of-1.5-billion-mobile-device-shipments>, January 2012.

³See Sysgo press releases <http://www.sysgo.com/>.

2.2 Modern representatives

We base our examination of the evolution of L4 design and evaluation on seL4, which we know well and which in many ways evolved furthest from the original design. We note where other recent versions ended up with different designs, and try to understand the reasons behind, and what this tells us about the degree of agreement about microkernel design in the L4 community.

Unlike any of the other systems, seL4 is designed from the beginning to support formal reasoning about security and safety, while maintaining the L4 tradition of minimality, performance and the ability to support almost arbitrary system architectures.

This led us to a radically new resource-management model, where all spatial allocation is explicit and directed by user-level code, including kernel memory [Elkaduwe et al., 2008]. It is also the first protected OS kernel in the literature with a complete and sound worst-case execution time (WCET) analysis [Blackham et al., 2011].

A second relevant system is Fiasco.OC, which is unique in that it is a code base that has lived through most of L4 history, starting as a clone of the original ABI (and not even designed for performance). It has, at some point in time, supported many different L4 ABI versions, often concurrently, and is now a high-performance kernel with the characteristics of the latest generation, including capability-based access control. Fiasco served as a testbed for many design explorations, especially with respect to real-time support [Härtig and Roitzsch, 2006].

Then there are two recent from-scratch designs: NOVA [Steinberg and Kauer, 2010], designed for hardware-supported virtualization on x86 platforms, and the OKL4 Microvisor [Heiser and Leslie, 2010] (“OKL4” for short), which was designed as a commercial platform for efficient para-virtualization on ARM processors.

Name	Architecture	Size (kLOC)		
		C/C++	asm	Total
Original	486	0	6.4	6.4
L4/Alpha	Alpha	0	14.2	14.2
L4/MIPS	MIPS64	6.0	4.5	10.5
Hazelnut	x86	10.0	0.8	10.8
Pistachio	x86	22.4	1.4	23.0
L4-embedded	ARMv5	7.6	1.4	9.0
OKL4 3.0	ARMv6	15.0	0.0	15.0
Fiasco.OC	x86	36.2	1.1	37.6
seL4	ARMv6	9.7	0.5	10.2

Table 2: Source lines of code (SLOC) of various L4 kernels.

3 Microkernel Design

Liedtke [1995] outlines principles and mechanisms which drove the design of the original L4. We examine how these evolved over time, and, specifically, how they compare with the current generation.

3.1 Minimality

The main drivers of Liedtke’s designs were minimality and IPC performance, with a strong belief that the former helps the latter. Specifically, he formulated the microkernel *minimality principle*:

A concept is tolerated inside the μ -kernel only if moving it outside the kernel, i.e. permitting competing implementations, would prevent the implementation of the system’s required functionality [Liedtke, 1995].

This principle, which is a more pointed formulation of “only minimal mechanisms and no policy in the kernel,” has continued to be the foundation of the design of L4 microkernels. The discussion in the following sections will demonstrate the community’s on-going efforts to remove features or replace them with more general (and powerful) ones.

The adherence to this principle can be seen from the comparison of source code sizes, shown in Table 2: while it is normal for systems to grow in size over time, seL4, the latest member of the family (and, arguably, the one that diverged strongest from the traditional model) is still essentially the same size as the early versions.⁴ Verification provided a particular strong motivation for

⁴In fact, seL4’s SLOC count is somewhat bloated as a consequence of the C code being mostly a “blind” manual translation from Haskell [Klein et al., 2009], together with generated bit-field accessor functions, resulting in hundreds of small functions. The kernel compiles into about 9 k ARM instructions.

minimality, as even 9,000 SLOC pushed the limits of what was achievable.

Retained: Minimality as key design principle.

Nevertheless, none of the designers of L4 kernels to date claim that they have developed a “pure” microkernel in the sense of strict adherence to the minimality principle. For example, all of them have a scheduler in the kernel, which implements a particular scheduling policy (usually hard-priority round-robin). To date, no-one has come up with a truly general in-kernel scheduler or a workable mechanism which would delegate all scheduling policy to user-level without imposing high overhead.

3.2 IPC

We mentioned earlier the importance of IPC performance, and that the design and implementation of L4 kernels consistently aimed at maximising it. However, the details have evolved considerably.

3.2.1 Synchronous IPC

The original L4 supported synchronous (rendezvous-style) IPC as the only communication, synchronisation, and signalling mechanism. Synchronous IPC avoids buffering in the kernel and the management and copying cost associated with it. It is also a prerequisite for a number of implementation tricks we will cover later, specifically the lazy scheduling (Section 4.2), direct process switch (Section 4.3), and temporary mapping (Section 3.2.2) optimisations.

While certainly minimal, and simple conceptually and in implementation, experience taught us significant drawbacks of this model: it forces a multi-threaded design onto otherwise simple systems, with the resulting synchronisation complexities. For example, the lack of functionality similar to UNIX `select()` required separate threads per interrupt source, and a single-threaded server could not wait for client requests and interrupts at the same time.

We addressed this in L4-embedded by adding *asynchronous notifications*, a very simple form of asynchronous IPC. We later refined this model in seL4 as *asynchronous endpoints* (AEPs, endpoints will be explained in Section 3.2.3): sending is non-blocking and asynchronous to the receiver, who has a choice of blocking or polling for a message. Logically, an asynchronous endpoint is similar to multiple binary semaphores allocated within a single word: Each AEP has a single-word *notification field*. A send operation specifies a mask of

bits (usually a single bit), which is OR-ed into the notification field. Waiting is effectively `select()` across the notification field.

In our present design, an asynchronous endpoint can be *bound* to a specific thread. If a notification arrives while the thread is waiting on a synchronous endpoint, the notification is delivered like a synchronous message (with an indication that it is really a notification).

In summary, seL4, like most other L4 kernels, retains the model of synchronous IPC but augments it with asynchronous notification. OKL4 has completely abandoned synchronous IPC, and replaced it by *virtual IRQs* (essentially asynchronous notifications). NOVA has augmented synchronous IPC with counting semaphores [Steinberg and Kauer, 2010], while Fiasco.OC has also augmented synchronous IPC with virtual IRQs.

Replaced: Synchronous IPC augmented with (seL4, NOVA, Fiasco.OC) or replaced by (OKL4) asynchronous notification.

Having two forms of IPC is a violation of minimality, so the OKL4 approach is more pure in this particular point (although its channel abstraction constitutes a different violation, see Section 3.2.2). Furthermore, the utility of synchronous IPC becomes more dubious in a multicore context: an RPC-like server invocation sequentialises client and server, which should be avoided if they are running on separate cores. We therefore expect communication protocols based on asynchronous notification to become more prevalent, and a move to asynchronous-only IPC sometime in the future remains a possibility.

3.2.2 IPC message structure

Original L4 IPC had rich semantics. Besides in-register (“short”) messages it supported messages of almost arbitrary size (a word-aligned “buffer” as well as multiple unaligned “strings”) in addition to in-register arguments. Coupled with the all-synchronous design, this approach avoids redundant copying.

Register arguments support zero-copy: the kernel always initiates the IPC from the sender’s context and switches to the receiver’s context without touching the message registers.

A drawback is the architecture-dependent and (especially on x86-32) small size of zero-copy messages. In fact, the number of available registers changed frequently with ABI changes, as changes to syscall arguments used or freed up registers.

Pistachio introduced the concept of a moderate-size (configurable in the range of 16–64) set of *virtual message registers*. The implementation mapped some of them to physical registers, the rest was contained in a

per-thread pinned part of the address space. The pinning ensures register-like semantics without the possibility of a page fault. Inlined access functions hide the distinction between physical and memory-backed registers from the user. seL4 and Fiasco.OC continue to use this approach.

The motivation is two-fold: Virtual message registers improve portability across architectures, and more importantly, they reduce the performance penalty for moderately-sized messages exceeding the number of physical registers: copying a small number of words is cheaper than establishing the temporary mapping involved in “long” IPC, as described below.

The benefits of in-register message transfers has diminished over time, as the architectural costs of context switching dominate IPC performance. For example, in-register message transfer on the ARM11 improves IPC performance by 10% (for a 4-word message) compared to passing via the kernel stack; on Cortex A9 this reduces to 4%. On x86-32, reserving any registers for message passing is detrimental to the compiler’s ability to optimise the code.

Replaced: Physical by virtual message registers.

“Long” messages could specify multiple buffers in a single IPC invocation to amortise the hardware mode- and context-switch costs. Long messages could be delivered with a single copy: executing in the sender’s context, the kernel sets up a temporarily mapped window into the receiver’s address space, covering (parts of) the message destination, and copies directly to the receiver.

This could trigger a page fault during copying in either the source or destination address space, which required the kernel to handle nested exceptions. Furthermore, the handling of such an exception required invoking a user-level page-fault handler. The handler had to be invoked while the kernel was still handling the IPC system call, yet the invocation had to pretend that the fault happened during normal user-mode execution. On return, the original system-call context had to be re-established. The result was significant kernel complexity, with many tricky corner cases that risked bugs in the implementation.

While long IPC provides functionality which cannot be emulated without some overhead, in practice it was rarely used: Shared buffers can avoid any explicit copying between address spaces, and are generally preferred for bulk data transfer. Additionally, asynchronous interfaces can be used for batching of transfers without resorting to explicit batching support in the kernel.

The main use of long IPC was for legacy POSIX read-write interfaces to servers, which require transferring the contents of arbitrary buffers to servers who do not necessarily have access to the client’s memory. However, the

rise of Linux as POSIX middleware, where Linux effectively shares memory with its applications, replaced this common use case with pass-by-reference. The remaining use cases either had interface flexibility or could be implemented with shared memory. Long IPC also violates the minimality principle (which talks about functionality, not performance).

As a consequence of this kernel complexity and the existence of user-level alternatives, we removed long IPC from L4-embedded, and NOVA and Fiasco.OC do not provide it either.

For seL4 there are even stronger reasons for staying away from supporting long messages: the formal verification approach explicitly avoided any concurrency in the kernel [Klein et al., 2009], and nested exceptions introduce a degree of concurrency. They also break the semantics of the C language by introducing additional control flow. While it is theoretically possible to formally reason about nested exceptions, it would make the already challenging verification task even harder. Of course, the in-kernel page faults could be avoided with extra checking, but that would introduce yet more complexity (besides penalising best-case performance), and would still require a more complex formal model to prove checking is complete and correct.

Abandoned: Long IPC.

OKL4 diverges at this point, by providing a new, asynchronous, single-copy, bulk-transfer mechanism called *channel* [Heiser and Leslie, 2010]. However, this is really a compromise retained for compatibility with the earlier OKL4 microkernel, which was aimed at memory-starved embedded systems. It was used to retrofit protection boundaries into a highly multithreaded (> 50 threads) real-time application, where a separate communication page per pair of communicating threads was too costly.

3.2.3 IPC destinations

Original L4 had threads as the targets of IPC operations. The motivation was to avoid the cache and TLB pollution associated with a level of indirection, although Liedtke [1993a] notes that ports could be implemented with an overhead of 12% (mostly 2 extra TLB misses). The model required that thread IDs were unique identifiers.

This model has a drawback of poor information hiding. A multi-threaded server has to expose its internal structure to clients, in order to spread client load, or use a gateway thread, which could become a bottleneck and would impose additional communication and synchronisation overhead. There were a number of proposals

to mitigate this but they all had drawbacks. Additionally, large-page support in modern CPUs has reduced the TLB pollution of indirection by increasing the likelihood of co-location on the same page. Last but not least, the global IDs introduced covert channels [Shapiro, 2003].

Influenced by EROS [Shapiro et al., 1999], *IPC endpoints* were adopted as IPC destinations by seL4 and Fiasco.OC [Lackorzynski and Warg, 2009]). seL4 (synchronous) endpoints are essentially ports: the root of the queue of pending senders or receivers is a now a separate kernel object, instead of being part of the recipient's thread control block (TCB). Unlike Mach ports, IPC endpoints do not provide any buffering.

Replaced: Thread IDs by port-like IPC endpoints as message destinations.

3.2.4 IPC timeouts

A blocking IPC mechanism creates opportunities for denial-of-service (DOS) attacks. For example, a malicious (or buggy) client could send a request to a server without ever attempting to collect the reply; owing to the rendezvous-style IPC, the sender would block indefinitely unless it implements a watchdog to abort and restart. L4's long IPC enables a slightly more sophisticated attack: A malicious client could send a long message to a server, ensure that it would page fault, and prevent its pager from servicing the fault.

To protect against such attacks, IPC operation in the original L4 had timeouts. Specifically, an IPC syscall specified 4 timeouts: one to limit blocking until start of the send phase, one to limit blocking in the receive phase, and two more to limit blocking on page faults during the send and receive phases (of long IPC).

Timeout values were encoded in a floating-point format that supported the values of zero, infinity, and finite values ranging from one millisecond to weeks. They added complexity for managing wakeup lists.

Practically, however, timeouts were of little use as a DOS defence. There is no theory, or even good heuristics, for choosing timeout values in a non-trivial system, and in reality, only the values zero and infinity were used: A client sends and receives with infinite timeouts, while a server waits for a request with an infinite but replies with a zero timeout. (The client uses an RPC-style `call` operation, consisting of a send followed by an atomic switch to a receive phase, guaranteeing that the client is ready to receive the server's reply.) Traditional watchdog timers represent a better approach to detecting unresponsive IPC interactions (e.g. resulting from deadlocks).

Having abandoned long IPC, in L4-embedded we replaced timeouts by a single flag supporting a choice

of polling (zero timeout) or blocking (infinite timeout). Only two flags are needed, one for the send and one for the receive phase. seL4 follows this model. A fully-asynchronous model, such as that of OKL4, is incompatible with timeouts and has no DOS issues that would require them.

Timeouts could also be used for timed sleeps by waiting on a message from a non-existing thread, a feature useful in real-time system. Dresden experimented with extensions, including absolute timeouts, which expire at a particular wall clock time rather than relative to the commencement of the system call. Our approach is to give userland access to a (physical or virtual) timer.

Abandoned: IPC timeouts in seL4, OKL4.

3.2.5 Communication Control

In the original L4, the kernel delivered the sender's unforgeable ID to the receiver. This allows a server to implement a form of discretionary access control, by ignoring undesirable messages. However, a server can be bombarded with spurious large messages by malicious clients. The time consumed by receiving such messages (even if copying is done by the kernel) prevents the server from performing useful work, and checking which ones to discard also costs time. Hence such messages can constitute a DOS attack, which can only be avoided by kernel support that prevents undesirable messages being sent in the first place [Liedtke et al., 1997b]. Mandatory access control policies also require a mechanism for mediating and authorising communication.

Original L4 provided this through a mechanism called *clans and chiefs*: Processes were organised in a hierarchy of "clans", each of which had a designated "chief". Inside the clan, all messages are transferred freely and the kernel guarantees message integrity. But messages crossing a clan boundary, whether outgoing or incoming, are redirected to the clan's chief, who can thus control the flow of messages. The mechanism also supports *confinement* [Lipner, 1975] of untrusted subsystems.

Liedtke [1995] argued that the clans-and-chiefs model only added two cycles per IPC operation, as clan IDs were encoded in thread IDs for quick comparison. However, the low overhead only applies where direct communication is possible. Once messages get re-directed, each such re-direction adds two messages to a (logically) single round-trip IPC, a significant overhead. Furthermore, the strict thread hierarchy was unwieldy in practice (and was probably the feature most cursed by people trying to build L4-based systems). For mandatory access control, the model quickly deteriorated into a chief per process. It is a prime example of kernel-enforced policy (address-space hierarchy) limiting the design space.

As a consequence of these drawbacks, many L4 implementations did not implement clans and chiefs (or disabled the feature at build time), but that meant that there was no way to control IPC. There were experiments with models based on a more general form of IPC redirection [Jaeger et al., 1999], but these failed to gain traction. The problem was finally resolved with flexible capability-mediated access control to endpoints.

Abandoned: Clans and chiefs.

3.3 User-level device drivers

A key consequence of the minimality principle, and maybe the most radical novelty of L4 (or, rather, its predecessor, L3 [Liedtke et al., 1991]) was to make all device drivers user-level processes.⁵ This is still a hallmark of all L4 kernels, and verification is a strong motivator for sticking with the approach: adding any unverified code, such as drivers, into the kernel would obliterate any guarantees, and verifying the large amount of driver code in real-world systems is out of reach for now.

A small number of drivers are still best kept in the kernel. In a modern L4 kernel this typically means a timer driver, used for preempting user processes at the end of their time slice, and a driver for the interrupt controller, which is required to safely distribute interrupts to user-level processes.

The user-level driver model is tightly coupled with modelling interrupts as IPC messages, which the kernel sends to the driver. Details of the model (IPC from a virtual thread vs upcall), as well as the association and acknowledgment protocol, have changed over the years (and at times changed back and back again) but the general approach still applies.

The most notable change was moving from synchronous to asynchronous IPC for interrupt delivery. This was driven by implementation simplification, as synchronous delivery required the emulation of virtual in-kernel threads as the sources of interrupt IPC.

User-level drivers have benefited from virtualisation-driven hardware improvements. I/O memory management units (IOMMUs) have enabled safe pass-through device access for drivers. User-level drivers have also benefited from hardware developments that reduce interrupt overheads, specifically interrupt coalescing support on modern network interfaces. In the x86 world, they have profited from dramatically decreased context-switching costs enabled by TLB tagging (among others).

Retained: User-level drivers as a core feature.

⁵Strictly speaking, this had been done before, in the Michigan Terminal system [Alexander, 1972] and the Monads OS [Keedy, 1979], but those designs had little direct impact on later ones and there is no information about performance.

Of course, user-level drivers have now become mainstream. They are supported (if not encouraged) on Linux, Windows and MacOS. Overheads in those systems are generally higher than in L4 with its highly optimised IPC, but we have shown in the past that low overheads are achievable even on Linux, at least on context-switch friendly hardware [Leslie et al., 2005]. In practice, though, only a tiny fraction of devices are performance-critical.

3.4 Resource management

Original L4's resource management, like its approach to communication control, was heavily based on process hierarchies. This applied to managing processes as well as virtual memory. Hierarchies are an effective way of managing and recovering resources, and provide a model of constraining sub-systems (where system mechanisms restrict children's privileges to be a subset of their parent's), but the cost is rigidity. However, the hierarchies are a form of policy, and as such a bad match for a microkernel (as discussed in Section 3.2.5).

Capabilities can provide a way out of the constraints of the hierarchy, which is one of several reasons all modern L4 kernels adopted capability-based access-control. Here we examine the most important resource-management issues arising from the original L4 model, and how we deal with them now.

3.4.1 Process hierarchy

A process (in L4 this is essentially a page table and a number of associated threads) consumes kernel resources, and unchecked allocation of TCBs and page tables could easily lead to denial of service. Original L4 dealt with that through a process hierarchy: "Task IDs" were essentially capabilities over address spaces, allowing creation and deletion.

There was a finite number of them (of the order of thousands), which the kernel handed out first-come, first-served. They could be delegated, but only up or down the hierarchy. (They were also closely tied to the thread hierarchy used for IPC control, see Section 3.2.5.) In a typical setup, the initial user process would grab all task IDs before creating any further processes.

Perhaps unsurprisingly, this model proved inflexible and restrictive; it was eventually replaced by fully-fledged capabilities.

Abandoned: Hierarchical process management.

3.4.2 Recursive page mappings

Original L4 tied authority over physical memory frames to existing page mappings. Having a valid mapping (of

a frame) in its address space gave a process the right to map this page into another address space. Instead of mapping, a process could *grant* one of its pages, which removed the page (and any authority over it) from the grantor. A mapping (but not a grant) could be revoked by an *unmap* operation. Address spaces were created empty, and were populated using the mapping primitive.

The recursive mapping model was anchored in a primordial address space σ_0 , which received a (one-on-one) mapping of all free frames left over after the kernel booted. σ_0 was the page-fault handler of all processes created at boot time, and would map each of its pages *once* to the first process that requested it (by faulting on an address in the page).

Note that, while the L4 memory model creates a hierarchy of mappings originating from each frame, it does not force a hierarchical view of address spaces: Mappings were established through IPC (similar to transferring a capability through an IPC message), and a process could map one of its pages to any other process it was allowed to send IPC to (provided the recipient agreed to receive mappings). Compared to Mach, L4 has no memory object semantics, only low-level address space management mechanisms that are closer to Mach's in-kernel pmap interface than its user-visible memory object abstraction [Rashid et al., 1988]. Memory objects, copy-on-write, and shadow-chains are all user-level created abstractions or implementation approaches.

The recursive mapping model was conceptually simple and elegant, and Liedtke was clearly proud of it – it figured prominently in many papers, including the first [Liedtke, 1993a], and in all his presentations. Yet, experience showed that there were significant drawbacks.

In order to support revocation at page granularity, the recursive address-space model requires substantial bookkeeping in the form of a *mapping database*. Moreover, the generality of the L4 memory model allows two colluding processes to force the kernel to consume large amounts of memory by recursively mapping the same frame to different pages in each other's address space, a potential DOS attack especially on 64-bit hardware, which can only be prevented by controlling IPC (via the dreaded clans-and-chiefs).

In L4-embedded we removed the recursive mapping model, after observing that for our real-world use cases, 25–50% of kernel memory use was consumed by the mapping database even without malicious processes. We replaced it by a model that more closely mirrors hardware, where mappings always originate from ranges of physical memory frames.

This approach comes at the expense of losing fine-grained delegation and revocation of memory (other than by brute-force scans of page tables), we therefore only considered it an interim pain relief. OKL4 somewhat

extends this minimal model, without achieving the generality and fine-grained control of the original model.

Mapping control is, of course, easily achieved in a capability-based system, using a variant of the standard grant-take model [Lipton and Snyder, 1977]. This is what seL4 does: the right to map is conveyed by a capability to a physical frame, not by having access to a virtual page backed by that frame, and thus seL4’s model is not recursive. Even with a frame capability, mapping is strictly limited by the explicit kernel memory model used to bookkeep the mappings, as described below.

Xen provides an interesting point of comparison. *Grant tables* allow the creation (based on possession of a valid mapping) of what is effectively a frame capability, which can be passed to another domain to establish shared mappings [Fraser et al., 2004]. A recent proposal extends grant tables to allow for revocation of frames [Ram et al., 2010]. The semantics of the memory mapping primitives is loosely similar to that of seL4, minus the propagation of page faults. In Xen’s case, the overhead for supporting fine-grained delegation and revocation is only paid in instances of sharing.

NOVA and Fiasco.OC both retain the recursive address space model, with authority to map determined by possession of a valid mapping. The consequent inability to restrict mapping and thus book keeping allocation is addressed by per-task kernel memory pools in Fiasco.OC.

The existing L4 address space models (specifically fine-grained delegation and revocation) represent different trade-offs between generality and minimality of mechanism, and potentially more space-efficient domain-specific approaches.

Multiple approaches: Some L4 kernels retain the model of recursive address-space construction, while seL4 and OKL4 originate mappings from frames.

3.4.3 Kernel memory

While capabilities provide a clean and elegant model for delegation, by themselves they do not solve the problem of resource management. A single malicious thread with grant right on a mapping can still use this to create a large number of mappings, forcing the kernel to consume large amounts of memory for meta-data, and potentially DOS-ing the system.

L4 kernels traditionally had a fixed-size heap from which the kernel allocated memory for its data structures. Original L4 had a *kernel pager*, called σ_1 , through which the kernel could request additional memory from userland. This does not solve the problem of malicious (or buggy) user code forcing unreasonable memory con-

sumption, it only shifts the problem. Consequently, σ_1 was not supported by most L4 kernels.

The fundamental problem, shared by most other OSes, is the insufficient isolation of user processes through the shared kernel heap. A satisfactory approach must be able to provide complete isolation. The underlying issue is that, even in a capability system, where authority is represented by capabilities, it is not possible to reason about the security state if there are resources outside the capability system.

Kernels that manage memory as a cache of user-level content only partially address this problem. While caching-based approaches remove the opportunity for DOS attacks based on memory exhaustion, they do not enable the strict isolation of kernel memory that is a prerequisite for performance isolation or real-time systems, and potentially introduce covert channels.

Liedtke et al. [1997b] examined this issue and proposed per-process kernel heaps together with a mechanism to donate extra memory to the kernel on exhaustion. NOVA, Fiasco and OKL4 all adopted variations of this approach. Per-process kernel heaps simplify user level (by removing control of allocation) at the expense of the ability to revoke allocations without destroying the process, and the ability to reason directly about allocated memory (as opposed to just bounding it). The trade-off is still being explored in the community.

We took a substantially different approach with seL4; its model for managing kernel memory is seL4’s main contribution to OS design. Motivated by the desire to reason about resource usage and isolation, we subject *all* kernel memory to authority conveyed by capabilities (except for the fixed amount used by the kernel to boot up, including its strictly bounded stack). Specifically, we completely remove the kernel heap, and provide userland with a mechanism to identify authorised kernel memory whenever the kernel allocates data structures. A side-effect is that this reduces the size and complexity of the kernel, a major bonus to verification.

The key is *making all kernel objects explicit* and subject to capability-based access control. This approach is inspired by hardware-based capability systems, specifically CAP [Needham and Walker, 1977] where hardware-interpreted capabilities directly refer to memory. HiStar [Zeldovich et al., 2011] also makes all kernel objects explicit, though it takes a caching approach to memory management.

Of course, user-visible kernel objects do not mean that someone with authority over a kernel object can directly read or write it. The capability provides the right to invoke (a subset of) object-specific methods, which includes destruction of the object. (Objects, once created, never change their size.) Crucially, the kernel object types include unused memory, called *Untyped* in seL4,

Object	Description
<i>TCB</i>	Thread control block
<i>Cnode</i>	Capability storage
<i>Synchronous Endpoint</i>	Port-like rendezvous object for synchronous IPC
<i>Asynchronous Endpoint</i>	Port-like object for asynchronous notification.
<i>Page Directory</i>	Top-level page table for ARM and IA-32 virtual memory
<i>Page Table</i>	Leaf page table for ARM and IA-32 virtual memory
<i>Frame</i>	4 KiB, 64 KiB, 1 MiB and 16 MiB objects that can be mapped by page tables to form virtual memory
<i>Untyped Memory</i>	Power-of-2 region of physical memory from which other kernel objects can be allocated

Table 3: seL4 kernel objects.

which can be used to create other objects.

Specifically, the only operation possible on *Untyped* is to *retype* part of it into some object type. The relationship of the new object to the original *Untyped* is recorded in a *capability derivation tree*, which also records other kinds of capability derivation, such as the creation of capability copies (with potentially reduced privileges). Once some *Untyped* has been retyped, the only operation possible on the (corresponding part of) the original *Untyped* is to revoke the derived object (see below).

Retyping is the only way to create objects in seL4. Hence, by limiting access to *Untyped* memory, a system can control resource allocation. Retyping can also produce smaller *Untyped* objects, which can then be independently managed – this is key to delegating resource management. The derivation from *Untyped* also ensures the kernel integrity property that no two typed objects overlap.

Table 3 gives the complete set of seL4 object types and their use. Userland can only directly access (load/store/fetch) memory corresponding to a *Frame* that is mapped in its address space (by inserting the *Frame* capability into a *Page Table*).

The resulting model has the following properties:

1. All authority is explicitly conferred (via capabilities).
2. Data access and authority can be confined.
3. The kernel itself (for its own data structures) adheres to the authority distributed to applications, including the consumption of physical memory.

4. Each kernel object can be reclaimed independently of any other kernel objects.
5. All operations execute, or are preemptible, in “short” time (constant or linear in the size of an object no bigger than a page).

Properties 1–3 ensure that it is possible to reason about system resources as well as security. Especially Property 3 was crucial to formally proving the kernel’s ability to ensure integrity, authority confinement and confidentiality [Murray et al., 2013; Sewell et al., 2011]. Property 5 ensures that all kernel latencies are bounded and thus supports its use for hard real-time systems [Blackham et al., 2011].

Property 4 ensures kernel integrity. Any holder of an appropriate capability can reclaim an object at any time (making the original *Untyped* again available for object creation). For example, page-table memory can be reclaimed without having to destroy the corresponding address space. This requires that the kernel is able to detect (and invalidate) any references to an object that is being reclaimed.

The requirement is satisfied with the help of the capability derivation tree. Objects are revoked by invoking the `revoke()` method on a *Untyped* object further up the tree; this will remove all capabilities to all objects derived from that *Untyped*. When the last capability to an object is removed, the object itself is deleted. This removes any in-kernel dependencies it may have with other objects, thus making it available for re-use. Removal of the last capability is easy to detect, as it cleans up the last leaf node in the capability tree referring to a particular memory location.

Revocation requires user-level book-keeping to associate *Untyped* capabilities with objects, often at the granularity of higher-level abstractions (such as processes) defined at user level. The precise semantics of *Untyped* and its relationship to user-level book-keeping is still being explored.

Added: User-level control over kernel memory in seL4, kernel memory quota in Fiasco.OC.

3.4.4 Time

Apart from memory, the other key resource that must be shared in a system is the CPU. Unlike memory, which can be sub-divided and effectively shared between multiple processes concurrently, the CPU can only be used by a single thread at a time, and must therefore be time-multiplexed.

All versions of L4 have achieved this multiplexing through a fixed-policy scheduler (pluggable in Fiasco.OC). The scheduling model of the original L4,

hard-priority round-robin, is still alive, despite being a gross heresy against the core microkernel religion of policy-freedom. All past attempts to export scheduling policy to user level have failed, generally due to intolerable overheads, or were incomplete or domain-specific.

Especially the Dresden group, which has a focus on real-time issues, experimented extensively with time issues, including absolute timeouts (see Section 3.2.4). They also explored several approaches to scheduling, as well as system structures suitable for real-time and analysed L4-based real time systems [Härtig and Roitzsch, 2006].

While able to address some specific problems, Dresden did not develop a policy-free and universal mechanism, and Fiasco.OC reverted to essentially the traditional L4 model. A more recent proposal for *scheduling contexts* allows mapping of hierarchical priority-based schedules onto a single priority scheduler [Lackorzynski et al., 2012]. While promising, this is not yet able to deal with the breadth of scheduling approaches used in the real-time community, especially earliest-deadline-first (EDF) scheduling.

One might argue that the notion of a single, general-purpose kernel suitable for all purposes may not be as relevant as it once was – these days we are used to environment-specific plugins. However, the formal verification of seL4 creates a powerful disincentive to changing the kernel, it strongly reinforces the desire to have a single platform for all usage scenarios. Hence, a policy-free approach to dealing with time is as desirable as it has ever been.

Unresolved: Principled, policy-free control of CPU time.

3.4.5 Multicore

Multiprocessor issues have been explored early-on in the L4 community. Most of the work, L4/Alpha and L4/MIPS notwithstanding, was done on x86 platforms, which were the earliest affordable multiprocessors. Early x86 multiprocessors and multicores had high inter-core communication cost and no shared caches. Consequently, the standard approach was to use per-processor scheduling queues (and minimal sharing of kernel data across cores), and thread migration only happening on explicit request by userland. Uhlig [2005] explored locking, synchronisation and consistency issues on platforms with many cores, and developed approaches for scalable concurrency control of kernel data structures based on RCU [McKenney et al., 2002]. NOVA and Fiasco.OC make extensive use of RCU.

With the shift of emphasis from high-end server to embedded and real-time platforms, multiprocessor issues took a back stage, and were only revived recently

with the advent of multicore versions of embedded processors. These are characterised by low inter-core communication cost and usually shared L2 caches, implying tradeoffs which differ from those on x86. The resulting low migration costs do typically not justify the overhead of a system call for migrating threads, and a global scheduler makes more sense.

Here, verification introduces new constraints. As discussed in Section 3.2.2, concurrency presents huge challenges for verification, and we kept it out of the kernel as much as possible. For multicores this means adopting either a big kernel lock or a multikernel approach [Baumann et al., 2009]. For a microkernel, where system calls are short, the former is not as silly as it may seem at first, as lock contention will be low, at least for a small number of cores sharing an L2.

The approach we are exploring at present is a *clustered multikernel*, a hybrid of a big-lock kernel (across cores which share an L2 cache), and a restricted variant of a multikernel (no memory migration is permitted between kernels) [von Tessin, 2012]. The clustered multikernel avoids concurrency in the majority of kernel code, which enables some of the formal guarantees to continue hold under some assumptions. The most significant assumptions are that (1) the lock itself, and any code outside of it, is correct and race free, and that (2) the kernel is robust to any concurrent changes to memory shared between the kernel and user-level (for seL4 is this is only a block of virtual IPC message registers).

The attraction of this approach is that it retains the existing uniprocessor proof with only small modifications. We have formally lifted a parallel composition of the uniprocessor automata and shown that refinement still holds. The disadvantage is that the formal guarantees no longer cover the entire kernel, and the large-step semantics used by the lifting framework preclude further extension of the formal framework to cover reasoning about the correctness of the lock, user-kernel concurrency, and any relaxation of resource migration restrictions.

A variation of a clustered multikernel may eventually be the best approach to obtaining full formal verification of a multiprocessor kernel, though we make no strong representations here. Much more work is required on the formal side to reason about fine-grained interleaving at the scale of a microkernel.

Unresolved: Handling of multicore processors in the age of verification.

4 Microkernel Implementation

Liedtke [1993a] list a set of design decisions and implementation tricks which helped making IPC fast in the

original i486 version, although a number of them smell of premature optimisation.

Some have already been mentioned, such as the temporary mapping window used in the (now obsolete) long IPC. Others are uncontroversial, such as the send-receive combinations in a single system call (the client-style call for an RPC-like invocation and the server-style reply-and-wait). We will discuss the remaining in more detail, including some traditional L4 implementation approaches which were less-publicised but long taken for granted in the community.

4.1 Strict process orientation and virtual TCB array

The original L4 had a separate kernel stack for each thread, allocated above its TCB on the same page. The TCB's base address was therefore a fixed offset from the stack base, and could be obtained by masking the least significant bits off the kernel stack pointer. Only a single TLB entry was required to cover both a thread's TCB and stack.

Furthermore, all TCBs were allocated in a sparse, virtually-addressed array, indexed by thread ID. During IPC, this enables a very fast lookup of the destination TCB, without first checking the validity of the ID: If the caller supplies an invalid ID, the lookup may access an unmapped TCB, triggering a page fault, which the kernel handles by aborting the IPC. If no fault happened, the validity of the thread ID can be established by comparing the caller-supplied value with the one found in the TCB. (Original L4's thread IDs had version numbers, which changed when the thread was destroyed and re-created. This was done to make thread IDs unique in time. Recording the current ID in the TCB allowed detecting stale thread IDs.)

Both features come at a cost: The many kernel stacks dominate the per-thread memory overhead, and they also increase the kernel's cache footprint. The virtual TCB array increases the kernel's virtual memory use and thus the TLB footprint, but avoids the additional cache footprint for the lookup table that would otherwise be required. Processors with a single page size and untagged TLBs left little opportunity to optimise beyond grouping data structures to minimise the number of pages touched. However, RISC processors had large-page sizes (or physical memory addressing) and tagged TLBs which changed the trade-offs.

The kernel's memory use became a significant issue when L4 was gaining traction in the embedded space, so the design needed revisiting.

Initial experiments with a single-stack kernel on a Pentium showed a reduction in kernel memory consumption, and improvements in IPC performance on

micro-benchmarks [Haeberlen, 2003]. Warton [2005] performed a thorough performance evaluation of the Pistachio process kernel vs an event-based (single-stack) kernel with continuations on an ARMv5 processor. He demonstrated comparable performance (generally within 1%) on micro-benchmarks, but a 20% performance advantage of the event kernel on a multi-tasking workload (AIM7). He also found that the event kernel's per-thread memory use was a quarter of that of the process kernel, despite the event kernel requiring more than twice the TCB size of the process kernel (to store the continuations).

Concurrently, Nourai [2005] analysed the trade-offs of virtual vs physical addressing of TCBs. He implemented physical addressing, also in Pistachio, although on a MIPS64 processor. He found little if any differences in IPC performance in micro-benchmarks, but significantly better performance of the physically-addressed kernel on workloads that stressed the TLB. MIPS is somewhat anomalous in that it supports physical addressing even with the MMU enabled, while on most other architectures "physical" addressing is simulated by idempotent large-page mappings, potentially in conjunction with "global" mappings. Still Nourai's results convincingly indicate that there is no significant performance benefit from the virtually-addressed TCBs.

An event-based kernel that avoids in-kernel page-fault exceptions preserves the semantics of the C language. As discussed earlier in Section 3.2.2, remaining within the semantics of C reduces the complexity of verification.

Together, these results made us choose an event-based design with physically-addressed kernel data for L4-embedded, and seL4 followed suit. While this decision was driven initially by the realities of resource-starved embedded systems and later the needs of verification, the approach's benefits are not restricted to those contexts, and we believe it is generally the best approach on modern hardware.

Replaced: Process kernel by event kernel in seL4, OKL4 and NOVA.

Abandoned: Virtual TCB addressing.

4.2 Lazy scheduling

In the rendezvous model of IPC, a thread's state frequently alternates between runnable and blocked. This implies frequent queue manipulations, moving a thread into and out of the ready queue, often many times within a time slice.

Liedtke's *lazy scheduling* trick minimises these queue manipulations: When a thread blocks on an IPC operation, the kernel updates its state in the TCB but leaves

the thread in the ready queue, with the expectation it will unblock soon. When the scheduler is invoked upon a time-slice preemption, it traverses the ready queue until it finds a thread that is really runnable, and removes the ones that are not. The approach was complemented by lazy updates of wakeup queues.

Lazy scheduling moves work from a high-frequency operation to the less frequently invoked scheduler. We observed the drawback when analysing seL4's worst-case execution time (WCET) for enabling its use in hard real-time systems [Blackham et al., 2012]: The execution time of the scheduler is only bounded by the number of threads in the system!

To address the issue, we adopted an alternative optimisation, referred to as *Benno scheduling*, which does not suffer from pathological timing behaviour: Here, the ready queue contains all runnable threads *except* the currently executing one. Hence, the ready queue usually does not get modified when threads block or unblock during IPC. At preemption time, the kernel inserts the (still runnable but no longer executing) preempted thread into the ready queue. The removal of timeouts means that there are no more wakeup queues to manipulate. Endpoint wait queues must be strictly maintained, but in the common case (of a server responding to client requests received via a single endpoint) they are hot in cache, so the cost of those queue manipulations is low. This approach has similar average-case performance as lazy scheduling, while also having a bounded WCET.

Replaced: Lazy scheduling by Benno scheduling.

4.3 Direct process switch

L4 traditionally tries to avoid running the scheduler during IPC. If a thread gets blocked during an IPC call, the kernel switches to a readily-identifiable runnable thread, which then executes on the original thread's time slice, generally ignoring priorities. This approach is called *direct process switch*.

It makes more sense than one might think at first, especially when assuming that servers have at least the same priority as clients. On the one hand, if a (client) thread performs a `call` operation (to a server), the caller will obviously block until the callee replies. Having been able to execute the syscall, the thread must be the highest-priority runnable thread, and the best way to observe its priority is to ensure that the callee completes as quickly as possible (and the callee is likely of higher priority anyway).

On the other hand, if a server replies (using *reply-and-wait*) to a waiting client, and the server has a request waiting from another client, it makes sense

to continue the server to take advantage of the primed cache by executing the receive phase of its IPC.

Modern L4 versions, concerned about correct real-time behaviour, retain direct-process switch where it conforms with priorities, and else invoke the scheduler. In fact, direct-process switch is a form of time-slice donation, and Steinberg et al. [2005] showed that can be used to implement priority-inheritance and priority-ceiling protocols. Fiasco.OC and NOVA support this by allowing the user to specify donation on a per-call basis.

Replaced: Direct process switch subject to priorities in seL4 and optional in Fiasco.OC and NOVA.

4.4 Preemption

Traditionally L4 implementations had interrupts disabled while executing within the kernel, although some (like L4/MIPS) contained preemption points in long-running operations, where interrupts were briefly enabled. Such an approach significantly simplifies kernel implementation, as most of the kernel requires no concurrency control, and generally leads to better average-case performance.

However, the original L4 ABI had a number of long-running system calls, and early Fiasco work made the kernel fully preemptive in order to improve real-time performance [Hohmuth and Härtig, 2001]. Later ABI versions removed most of the long-running operations, and Fiasco.OC reverted to the original, mostly non-preemptible approach.

In the case of seL4, there is an additional reason for a non-preemptible kernel: avoiding concurrency to make formal verification tractable [Klein et al., 2009]. Given seL4's focus on safety-critical systems, many of which are of a hard real-time nature, we need hard bounds on the latency of interrupt delivery. It is therefore essential to avoid long-running kernel operations, and use preemption points where this is not possible (e.g. the practically unbounded object deletion). We put significant effort into placement of preemption points, as well as on data structures and algorithms that minimises the need for them [Blackham et al., 2012].

Note that a continuation-based event kernel provides natural support for preemption points (by making them continuation points).

Retained: Mostly non-preemptible design with strategic preemption points.

4.5 Non-portability

Liedtke [1995] makes the point that a microkernel implementation should not strive for portability, as a

hardware abstraction introduces overheads and hides hardware-specific optimisation opportunities. He cites subtle architectural changes between the “compatible” i486 and Pentium processors resulting in shifting trade-offs and implying significant changes in the optimal implementation.

This argument was debunked by Liedtke himself, with the high-performance yet portable Hazelnut kernel and especially Pistachio. Careful design and implementation made it possible to develop an implementation that was 80–90% architecture-agnostic.

In seL4, the architecture-agnostic code (between x86 and ARM) only accounts for about 50%. About half the code deals with virtual memory management which is necessarily architecture-specific. The lower fraction of portable code is a result of seL4’s overall smaller size, with most (architecture-agnostic) resource-management code moved to userland. There is little architecture-specific optimisation except for the IPC fastpath. Steinberg [2013] similarly estimates a 50% rewrite for porting NOVA to ARM.

Replaced: Non-portable implementation by significant portion of architecture-agnostic code.

4.6 Non-standard calling convention

The original L4 kernel was completely implemented in assembler, and therefore the calling convention for functions was irrelevant inside the kernel. At the ABI, all registers which were not needed as syscall parameters were designated as message registers. The library interface provided inlined assembler stubs to convert the compiler’s calling convention to the kernel ABI (in the hope the compiler would optimise away any conversion overhead).

The next generation of L4 kernels, starting with L4/MIPS, were all written at least partially in C. At the point of entering C code, these kernels had to re-establish the C compiler’s calling convention, and revert to the kernel’s convention on return. This made calling C functions relatively expensive, and therefore discouraged the use of C except for inherently expensive operations.

Later kernels were written almost exclusively in C (Hazelnut) or C++ (Fiasco, Pistachio). The cost of the calling-convention mismatch (and the lack of Liedtke-style masochism required for micro-optimising every bit of code) meant that the C code did not exhibit performance that was competitive to the old assembler kernel. The implementors of those kernels therefore started to introduce hand-crafted assembler *fast paths*. These led to IPC performance comparable to the original L4 (see Table 1).

The traditional approach was unsuitable for seL4, as the verification framework could only deal with C code [Klein et al., 2009], and we wanted to verify the kernel’s functionality as completely as feasible. This requires restricting assembler code to the bare minimum, and rules out calling-convention conversions, forcing us to adopt the tool chain’s standard calling conventions.

Abandoned: Non-standard calling conventions.

4.7 Implementation language

seL4 is also highly dependent on fast-path code to obtain competitive IPC performance, but the fast paths must now be implemented in C. The assembler fast path had already been abandoned in the commercial OKL4 kernel because of the high maintenance cost of assembler code, which in the commercial environment outweighed any performance degradation. For seL4 we were willing to tolerate no more than a 10% degradation in IPC performance, compared to the fastest kernels on the same architecture.

Fortunately, it turned out that by carefully hand-crafting the fast path, we can achieve highly-competitive IPC latencies [Blackham and Heiser, 2012]. Specifically this means manually re-ordering statements, making use of (verified) invariants that the compiler is unable to determine by static analysis,

In fact, the finally achieved latency of 188 cycles for a one-way IPC on an ARM11 processor is about 10% *better* than the fastest IPC we had measured on any other kernel on the same hardware! This is partially a result of the simplified seL4 ABI and IPC semantics, and the fact that the event-based kernel no longer requires saving and restoring the C calling convention on a stack switch. We also benefit from improved compilers, especially their support for annotating condition branches for the common case, which helps code locality.

In any case, this result demonstrates that assembler implementations are no longer justified by performance arguments. The Dresden team in fact found that they could achieve highly-competitive performance without any fastpathing.

Abandoned: Assembler code for performance.

The first L4 kernel written completely in a high-level language was Fiasco, which chose C++ rather than C (which had been used for parts of the MIPS kernel a few years earlier). Given the state of C++ compilers at the time, this may seem a courageous decision, but is at least partially explained by the fact that Fiasco was not initially designed with performance in mind. This changed later, and the recent Fiasco experience demonstrates that,

when used correctly, there is no performance penalty from C++ code.

The Karlsruhe team also chose C++ for Pistachio, mostly to support portability. Despite a high degree of enthusiasm about C++ in Dresden and Karlsruhe, we never saw any convincing advantages offered by C++ for microkernel implementation. Furthermore, OK Labs found that the availability of good C++ compilers was a real problem in the embedded space, and they converted their version of the microkernel back to straight C.

For seL4, the requirements of verifications forced the choice of C. While Dresden's VFiasco project attempted to verify the C++ kernel [Hohmuth and Tews, 2005], it never completed formalising the semantics of the C++ subset used by Fiasco. In contrast, by using C to implement seL4, we could build on an existing formalisation of C [Norris, 1998], a key enabler for the verification.

Abandoned: C++ for seL4 and OKL4.

5 Conclusions

It is rare that a research operating system has both a significant developer community, significant commercial deployment, as well as a long period of evolution. L4 is such a system, with 20 years of evolution of the API, of design and implementation principles, and about a dozen from-scratch implementations. We see this as a great opportunity to reflect on the principles and know-how that has stood the test of time, and what has failed to survive increased insights, changed deployment scenarios and the evolution of CPU architectures.

Design choices and implementation tricks came and went (including some which were close to the original designer's heart). However, the most general principles behind L4, minimality (including running device drivers at user level) and a strong focus on performance, still remain relevant and foremost in the minds of developers. Specifically we find that the key microkernel performance metric, IPC latency, has remained essentially unchanged (in terms of clock cycles), as far as comparisons across vastly different ISAs and micro architectures have any validity, in stark contrast to the trend identified by Ousterhout [1990] just a few years before L4 was created. Furthermore, and maybe most surprisingly, the code size has essentially remained constant, a rather unusual development in software systems.

Formal verification increased the importance of minimality, and also increased pressure for simplification of the implementation. Several design decisions, such as the simplified message structure, user-level control of kernel memory and the approach to multicores are strongly influenced by verification. It also impacted a number of implementation approaches, such as the use

of an event-oriented kernel, adoption of standard calling convention, and the choice of C as the implementation language. However, we do not think that this has led to tradeoffs which we would consider inferior when ignoring verification.

With formal verification, L4 has convincingly delivered on one of the core promises microkernels made many years ago: robustness. We think it is a great testament to the brilliance of Liedtke's original L4 design that this was achieved while, or maybe due to, staying true to the original L4 philosophy. It may have taken an awfully long time, but time has finally proved right the once radical ideas of Brinch Hansen [1970].

There is one concept that has, so far, resisted any satisfactory abstraction: *time*. L4 kernels still implement a specific scheduling policy – in most cases priority-based round-robin – the last major holdout of policy in the kernel. This probably represents the largest limitation of generality of L4 kernels. There is work underway at Dresden and NICTA that indicates that a single, parameterised kernel scheduler may actually be able to support all standard scheduling policies, and we expect it will not take another 20 years to get there.

Acknowledgements

NICTA is funded by the Australian Government as represented by the Department of Broadband, Communications and the Digital Economy and the Australian Research Council through the ICT Centre of Excellence program.

L4 would not exist without its inventor, Jochen Liedtke, and we pay tribute to his brilliance. We are also greatly indebted to the many people who contributed to L4 over two decades, generations of staff and students at IBM Watson, TU Dresden, University of Karlsruhe, the University of New South Wales and NICTA; there are too many to name them all.

We are specifically grateful to members of the L4 community who provided feedback on drafts of this paper: Andrew Baumann, Ben Leslie, Chuck Gray, Hermann Härtig. We thank Adam Lackorzynski for digging out the original L4 sources and extracting SLOCcounts, and Adrian Danis for some last-minute seL4 optimisations and measurements. We thank the anonymous SOSP reviewers for their insightful comments, and our shepherd John Ousterhout for his feedback.

References

Michael T. Alexander. Organization and features of the Michigan terminal system. In *AFIPS Conference Pro-*

- ceedings, 1972 Spring Joint Computer Conference, pages 585–591, 1972.
- Andrew Baumann, Paul Barham, Pierre-Evariste Dagand, Tim Harris, Rebecca Isaacs, Simon Peter, Timothy Roscoe, Adrian Schüpbach, and Akhilesh Singhanian. The multikernel: A new OS architecture for scalable multicore systems. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles*, Big Sky, MT, USA, October 2009. ACM.
- Bernard Blackham and Gernot Heiser. Correct, fast, maintainable – choose any three! In *Proceedings of the 3rd Asia-Pacific Workshop on Systems (AP-Sys)*, pages 13:1–13:7, Seoul, Korea, July 2012. doi: 10.1145/2349896.2349909.
- Bernard Blackham, Yao Shi, Sudipta Chattopadhyay, Abhik Roychoudhury, and Gernot Heiser. Timing analysis of a protected operating system kernel. In *Proceedings of the 32nd IEEE Real-Time Systems Symposium*, pages 339–348, Vienna, Austria, November 2011. doi: 10.1109/RTSS.2011.38.
- Bernard Blackham, Yao Shi, and Gernot Heiser. Improving interrupt response time in a verifiable protected microkernel. In *Proceedings of the 7th EuroSys Conference*, pages 323–336, Bern, Switzerland, April 2012. doi: 10.1145/2168836.2168869.
- Per Brinch Hansen. The nucleus of a multiprogramming operating system. *Communications of the ACM*, 13: 238–250, 1970.
- J. Bradley Chen and Brian N. Bershad. The impact of operating system structure on memory system performance. In *Proceedings of the 14th ACM Symposium on Operating Systems Principles*, pages 120–133, Asheville, NC, USA, December 1993.
- Michael Condict, Don Bolinger, Dave Mitchell, and Eamonn McManus. Microkernel modularity with integrated kernel performance. Technical report, OSF Research Institute, June 1994.
- Jack B. Dennis and Earl C. Van Horn. Programming semantics for multiprogrammed computations. *Communications of the ACM*, 9:143–155, 1966.
- Dharmika Elkaduwe, Philip Derrin, and Kevin Elphinstone. Kernel design for isolation and assurance of physical memory. In *1st Workshop on Isolation and Integration in Embedded Systems*, pages 35–40, Glasgow, UK, April 2008. ACM SIGOPS. doi: 10.1145/1435458.
- Keir Fraser, Steven Hand, Rolf Neugebauer, Ian Pratt, Andrew Warfield, and Mark Williamson. Safe hardware access with the Xen virtual machine monitor. In *Proceedings of the 1st Workshop on Operating System and Architectural Support for the On-Demand IT Infrastructure (OASIS)*, 2004.
- Charles Gray, Matthew Chapman, Peter Chubb, David Mosberger-Tang, and Gernot Heiser. Itanium — a system implementor’s tale. In *Proceedings of the 2005 USENIX Annual Technical Conference*, pages 264–278, Anaheim, CA, USA, April 2005.
- Andreas Haeberlen. Managing kernel memory resources from user level. Diploma thesis, Dept of Computer Science, University of Karlsruhe, April 2003. URL http://os.ibds.kit.edu/english/97_639.php.
- Hermann Härtig and Michael Roitzsch. Ten years of research on L4-based real-time systems. In *Proceedings of the 8th Real-Time Linux Workshop*, Lanzhou, China, 2006.
- Hermann Härtig, Michael Hohmuth, Jochen Liedtke, Sebastian Schönberg, and Jean Wolter. The performance of μ -kernel-based systems. In *Proceedings of the 16th ACM Symposium on Operating Systems Principles*, pages 66–77, St. Malo, France, October 1997.
- Gernot Heiser and Ben Leslie. The OKL4 Microvisor: Convergence point of microkernels and hypervisors. In *Proceedings of the 1st Asia-Pacific Workshop on Systems (APSys)*, pages 19–24, New Delhi, India, August 2010.
- Michael Hohmuth and Hermann Härtig. Pragmatic non-blocking synchronization for real-time systems. In *Proceedings of the 2001 USENIX Annual Technical Conference*, Boston, MA, USA, 2001.
- Michael Hohmuth and Hendrik Tews. The VFiasco approach for a verified operating system. In *Proceedings of the 2nd Workshop on Programming Languages and Operating Systems (PLOS)*, Glasgow, UK, July 2005.
- Trent Jaeger, Kevin Elphinstone, Jochen Liedtke, Vsevolod Panteleenko, and Yoonho Park. Flexible access control using IPC redirection. In *Proceedings of the 7th Workshop on Hot Topics in Operating Systems*, Rio Rico, AZ, USA, March 1999.
- J. Leslie Keedy. On the programming of device drivers for in-process systems. Monads Report 5, Dept. of Computer Science, Monash University, Clayton VIC, Australia, 1979.

- 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. doi: 10.1145/1629575.1629596.
- Adam Lackorzynski and Alexander Warg. Taming subsystems: capabilities as universal resource access control in L4. In *2nd Workshop on Isolation and Integration in Embedded Systems*, pages 25–30, Nuremberg, Germany, March 2009.
- Adam Lackorzynski, Alexander Warg, Marcus Völpl, and Hermann Härtig. Flattening hierarchical scheduling. In *International Conference on Embedded Software*, pages 93–102, Tampere, Finland, October 2012.
- Ben Leslie, Peter Chubb, Nicholas FitzRoy-Dale, Stefan Götz, Charles Gray, Luke Macpherson, Daniel Potts, Yueting (Rita) Shen, Kevin Elphinstone, and Gernot Heiser. User-level device drivers: Achieved performance. *Journal of Computer Science and Technology*, 20(5):654–664, September 2005.
- Roy Levin, Ellis S. Cohen, William M. Corwin, Fred J. Pollack, and William A. Wulf. Policy/mechanism separation in HYDRA. In *Proceedings of the 5th ACM Symposium on Operating Systems Principles*, pages 132–140, 1975.
- Jochen Liedtke. Improving IPC by kernel design. In *Proceedings of the 14th ACM Symposium on Operating Systems Principles*, pages 175–188, Asheville, NC, USA, December 1993a.
- Jochen Liedtke. A persistent system in real use: Experience of the first 13 years. In *Proceedings of the 3rd IEEE International Workshop on Object Orientation in Operating Systems (IWOOS)*, pages 2–11, Asheville, NC, USA, December 1993b. IEEE.
- Jochen Liedtke. On μ -kernel construction. In *Proceedings of the 15th ACM Symposium on Operating Systems Principles*, pages 237–250, Copper Mountain, CO, USA, December 1995.
- Jochen Liedtke. Towards real microkernels. *Communications of the ACM*, 39(9):70–77, September 1996.
- Jochen Liedtke, Ulrich Bartling, Uwe Beyer, Dietmar Heinrichs, Rudolf Ruland, and Gyula Szalay. Two years of experience with a μ -kernel based OS. *ACM Operating Systems Review*, 25(2):51–62, April 1991.
- Jochen Liedtke, Kevin Elphinstone, Sebastian Schönberg, Herrman Härtig, Gernot Heiser, Nayeem Islam, and Trent Jaeger. Achieved IPC performance (still the foundation for extensibility). In *Proceedings of the 6th Workshop on Hot Topics in Operating Systems*, pages 28–31, Cape Cod, MA, USA, May 1997a.
- Jochen Liedtke, Nayeem Islam, and Trent Jaeger. Preventing denial-of-service attacks on a μ -kernel for WebOSes. In *Proceedings of the 6th Workshop on Hot Topics in Operating Systems*, pages 73–79, Cape Cod, MA, USA, May 1997b. IEEE.
- Steven. B. Lipner. A comment on the confinement problem. In *Proceedings of the 5th ACM Symposium on Operating Systems Principles*, pages 192–196. ACM, 1975.
- Richard J. Lipton and Lawrence Snyder. A linear time algorithm for deciding subject security. *Journal of the ACM*, 24(3):455–464, 1977. ISSN 0004-5411. doi: <http://doi.acm.org/10.1145/322017.322025>.
- Paul E. McKenney, Dipankar Sarma, Andrea Arcangeli, Andi Kleen, Orran Krieger, and Rusty Russell. Read copy update. In *Proceedings of the Ottawa Linux Symposium*, 2002. URL <http://www.rdrop.com/users/paulmck/rclock/rcu.2002.07.08.pdf>.
- 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 and Privacy*, pages 415–429, San Francisco, CA, May 2013. ISBN 10.1109/SP.2013.35.
- Roger M. Needham and R.D.H. Walker. The Cambridge CAP computer and its protection system. In *Proceedings of the 6th ACM Symposium on Operating Systems Principles*, pages 1–10. ACM, November 1977.
- Michael Norrish. *C formalised in HOL*. PhD thesis, University of Cambridge Computer Laboratory, 1998.
- Abi Nourai. A physically-addressed L4 kernel. BE thesis, School of Computer Science and Engineering, University of NSW, Sydney 2052, Australia, March 2005. Available from publications page at <http://ssrg.nicta.com.au/>.
- John K. Ousterhout. Why aren't operating systems getting faster as fast as hardware? In *Proceedings of the 1990 Summer USENIX Technical Conference*, pages 247–56, June 1990.

- Kaushik Kumar Ram, Jose Renato Santos, and Yoshio Turner. Redesigning Xen's memory sharing mechanism for safe and efficient I/O virtualization. In *Proceedings of the 2nd Workshop on I/O Virtualization*, Pittsburgh, PA, USA, 2010. USENIX.
- Richard Rashid, Avadis Tevanian, Jr., Michael Young, David Golub, Robert Baron, David Black, William J. Bolosky, and Jonathan Chew. Machine-independent virtual memory management for paged uniprocessor and multiprocessor architectures. *IEEE Transactions on Computers*, C-37:896–908, 1988.
- Thomas Sewell, Simon Winwood, Peter Gammie, Toby Murray, June Andronick, and Gerwin Klein. seL4 enforces integrity. In Marko C. J. D. van Eekelen, Herman Geuvers, Julien Schmaltz, and Freek Wiedijk, editors, *2nd International Conference on Interactive Theorem Proving*, volume 6898 of *Lecture Notes in Computer Science*, pages 325–340, Nijmegen, The Netherlands, August 2011. Springer. doi: http://dx.doi.org/10.1007/978-3-642-22863-6_24.
- Jonathan S. Shapiro. Vulnerabilities in synchronous IPC designs. In *Proceedings of the IEEE Symposium on Security and Privacy*, Oakland, CA, USA, May 2003. URL citeseer.ist.psu.edu/shapiro03vulnerabilities.html.
- Jonathan S. Shapiro, Jonathan M. Smith, and David J. Farber. EROS: A fast capability system. In *Proceedings of the 17th ACM Symposium on Operating Systems Principles*, pages 170–185, Charleston, SC, USA, December 1999. URL <http://www.eros-os.org/papers/sosp99-eros-preprint.ps>.
- Udo Steinberg. Personal communication, 2013.
- Udo Steinberg and Bernhard Kauer. NOVA: A microhypervisor-based secure virtualization architecture. In *Proceedings of the 5th EuroSys Conference*, Paris, France, April 2010.
- Udo Steinberg, Jean Wolter, and Hermann Härtig. Fast component interaction for real-time systems. In *EuroMicro Conference on Real-Time Systems*, pages 89–97, Palma de Mallorca, Spain, July 2005.
- Volkmar Uhlig. *Scalability of Microkernel-Based Systems*. PhD thesis, University of Karlsruhe, Karlsruhe, Germany, June 2005.
- Michael von Tessin. The clustered multikernel: An approach to formal verification of multiprocessor OS kernels. In *Proceedings of the 2nd Workshop on Systems for Future Multi-core Architectures*, Bern, Switzerland, April 2012.
- Matthew Warton. Single kernel stack L4. BE thesis, School of Computer Science and Engineering, University of NSW, Sydney 2052, Australia, November 2005.
- Nickolai Zeldovich, Silas Boyd-Wickizer, Eddie Kohler, and David Mazières. Making information flow explicit in HiStar. *Communications of the ACM*, 54(11): 93–101, November 2011.

译文：

摘要

L4 微内核系统经历了 20 年的发展与演进。它拥有活跃的用户和开发者社区，并且已有商业版本被大规模部署在各种系统中，包括安全性关键系统。本文总结了这 20 年来关于微内核设计与实现的经验教训。我们回顾了早期的 L4 设计论文，并分析了从最初的 L4 到最新一代 L4 微内核，尤其是 seL4 的设计与实现演变。seL4 将 L4 模型推进到了最远的程度，是第一个操作系统内核，其实现通过了完整的形式化验证，并完成了最坏情况执行时间（WCET）的可靠分析。我们将展示，尽管发生了许多变化，最小化原则和高性能的进程间通信（IPC）仍然是设计与实现决策的主要驱动力。

第1章简介

20 年前，Liedtke [1993a] 通过其 L4 内核展示了微内核的 IPC 可以非常快速，其性能比其他同时代微内核快了 10 到 20 倍。微内核将内核的职责压缩到最低限度：仅保留一组通用机制，而将真正的操作系统服务交由用户态的服务程序来实现 [Levin et al., 1975]。用户代码通过进程间通信（IPC）机制（通常是消息传递）与服务器通信以获取系统服务。因此，IPC 处于任何服务调用的关键路径上，降低 IPC 的开销至关重要。

在 90 年代初，IPC 性能成为了微内核的致命弱点：一次单向消息传递的典型成本大约为 100 微秒，这对于构建高性能系统来说代价太高，进而导致了将核心服务迁回内核的趋势 [Condict et al., 1994]。也有观点认为高 IPC 成本是微内核系统结构的固有结果 [Chen and Bershad, 1993]。

在这种背景下，Liedtke 展示的 IPC 成本数量级的改进显得尤为突出。随后，一系列工作探讨了 L4 的设计哲学与机制 [Liedtke, 1995, 1996]，展示了在 L4 上运行的半虚拟化 Linux，其性能开销仅为几个百分点 [Härtig et al., 1997]；L4 微内核被部署在数十亿台移动设备和安全性关键系统中，并最终实现了世界上第一个

操作系统内核的功能正确性形式化验证 [Klein et al., 2009]。并且，它还对其他操作系统的研究产生了深远影响，例如 Barrelfish [Baumann et al., 2009]。

本文考察了 L4 在过去 20 年的发展历程。我们重点分析现代 L4 微内核的核心特征，它们与 Liedtke 最初的设计和实现之间的联系，以及哪些微内核“基本要素”经受住了时间的检验。我们特别关注过去的经验教训如何影响了以 seL4 [Klein et al., 2009] 为代表的最新一代 L4 微内核的设计，同时指出当前其他的 L4 版本在设计上所做的不同选择。

第2章2 背景

2.1 L4 微内核家族

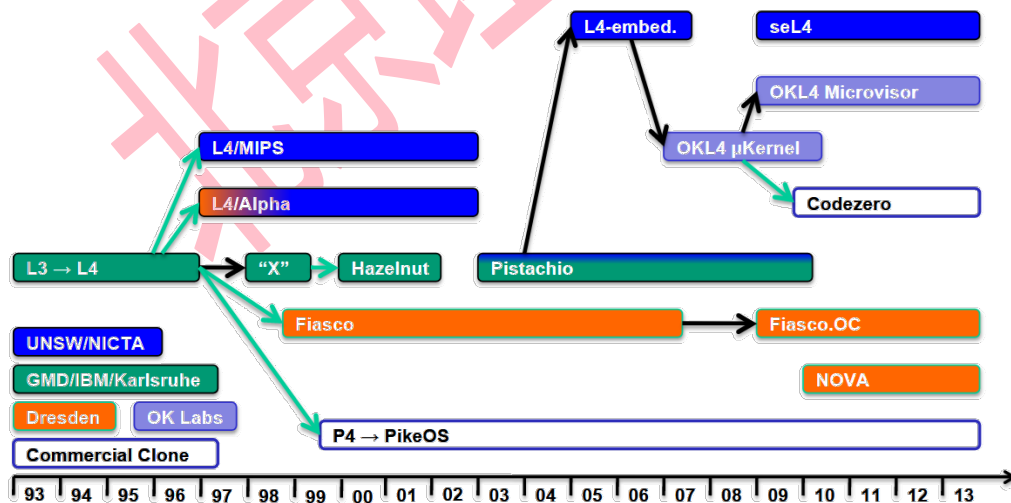


图 2-1 简化的 L4 家谱

L4 起源于 Liedtke 在 1980 年代初期在 i386 平台上开发的早期系统 L3 [Liedtke, 1993b]。L3 是一个具备内建持久性支持的完整操作系统，并且已经具备用户态驱动程序，这一点至今仍是 L4 微内核的典型特征。L3 被商业化部署在几千个安装场景中，主要是学校和律师事务所等机构。像当时所有微内核一样，L3 的进程间通信（IPC）开销大约为 100 微秒。

Liedtke 最初利用 L3 进行新思想的尝试，在其早期出版物中所称的“L3”实际上是一次彻底重新设计的中间版本。他从 1995 年起将“L4”这一名称用于社区中广泛传播的“V2” ABI 版本。本文中我们将这一版本称为“初代 L4”。Liedtke 在基于 i486 的 PC 上使用汇编语言完全实现了该内核，并很快将其移植到 Pentium 上。

这项初始工作启动了一个为期二十年的发展历程，包括多个 ABI 修订版本和从头开始的重新实现，如图 1 所示。最初是德累斯顿工业大学（TU Dresden）和新南威尔士大学（UNSW）分别在 64 位 Alpha 和 MIPS 处理器上重新实现了该 ABI（并作了必要调整），MIPS 实现中所有运行时间较长的操作均用 C 语言完成。这两个内核都实现了亚微秒级的 IPC 性能 [Liedtke et al., 1997a]，并以开源方式发布。UNSW 的 Alpha 版本是第一个支持多处理器的 L4 内核。

当时已从 GMD 转至 IBM Watson 的 Liedtke，继续在所谓的 Version X 中探索 ABI 的演进。GMD 和 IBM 实施的知识产权限制对于其他研究者来说过于严格，因此促使德累斯顿团队从头实现了一个新的 x86 版本，命名为 Fiasco，以此表达他们在处理 IP 问题上的挫折经历。Fiasco 是第一个几乎完全用高级语言（C++）编写的 L4 内核版本，也是至今仍在维护的最早 L4 代码库。它也是第一个获得较大规模商业部署的 L4 内核版本（估计出货量达 10 万台）。

Liedtke 后来转至卡尔斯鲁厄理工学院，他和他的学生基于 C 语言从零开发了 Hazelnut，这是第一个被移植（而非重新实现）到另一架构（从 Pentium 到 ARM）的 L4 内核。

卡尔斯鲁厄团队基于 Version X 和 Hazelnut 的经验提出了一个主要的 ABI 修订版本 V4，旨在提升内核和应用程序的可移植性、支持多处理器，并解决其它已知问题。在 Liedtke 于 2001 年不幸去世后，他的学生们将这一设计落实为一个新的开源内核 L4Ka::Pistachio（简称 Pistachio）。该内核使用 C++ 实现，最初支持 x86 和 PowerPC，UNSW /NICTA 的研究团队随后将其移植到了 MIPS、Alpha、64 位 PowerPC 和 ARM。在大多数移植中，修改的代码量不到 10%。

表 2-1 各种 L4 内核的单向 IPC 成本。

Name	Year	Processor	MHz	Cycles	μs
Original	1993	i486	50	250	5.00
Original	1997	Pentium	160	121	0.75
L4/MIPS	1997	R4700	100	86	0.86
L4/Alpha	1997	21064	433	45	0.10
Hazelnut	2002	Pentium 4	1,400	2,000	1.38
Pistachio	2005	Itanium 2	1,500	36	0.02
OKL4	2007	XScale 255	400	151	0.64
NOVA	2010	Core i7 (Bloomfield) 32-bit	2,660	288	0.11
seL4	2013	Core i7 4770 (Haswell) 32-bit	3,400	301	0.09
seL4	2013	ARM11	532	188	0.35
seL4	2013	Cortex A9	1,000	316	0.32

在 NICTA，我们很快将 Pistachio 重新定位到用于资源受限的嵌入式系统，形成了一个分支版本 NICTA::Pistachio-embedded（又称 L4-embedded）。当高通公司将其作为受保护模式的实时操作系统，用于无线调制解调器处理器的固件后，该版本得到了大规模商业部署。NICTA 的衍生公司 Open Kernel Labs 继续支持并开发该内核，并将其更名为 OKL4 微内核。另一个被广泛部署的版本是由 Sysgo 公司开发

的商业 V2 克隆版——PikeOS。它通过了用于安全关键航空电子设备的认证，并部署于飞机和列车中。

在日益增长的安全需求和 EROS [Shapiro et al., 1999] 的广泛影响下，访问控制机制开始转向能力模型 [Dennis and Van Horn, 1966]。OKL4 在 2008 年发布的 2.1 版本中首次引入能力机制，随后 Fiasco(后改名为 Fiasco.OC)也采用了这一机制。为了实现形式化验证，我们选择了从零开发一个专为能力机制设计的 seL4 内核，因为已有的代码库并未考虑可验证性。

近年来还出现了两个主要以虚拟化为核心目标的新设计：来自德累斯顿的 NOVA [Steinberg and Kauer, 2010] 和 OK Labs 的 OKL4 Microvisor [Heiser and Leslie, 2010]。

贯穿这二十年的主线是第 3.1 节中提出的“最小化”原则，以及对关键 IPC 操作性能的高度关注。如表 1 所示，内核开发者通常致力于尽可能的接近微架构所能实现的性能极限。因此，L4 社区倾向于以 CPU 周期而非微秒来衡量 IPC 延迟，因为这更贴近硬件限制。事实上，第 3.1 节还提供了关于硬件对上下文切换友好程度的有趣观点(可以比较 Pentium 4 和 Itanium 在高度优化 IPC 实现中的周期开销差异)。

2.2 现代代表

我们将分析 L4 设计演进的基础建立在 seL4 上——这是我们最熟悉的版本，也是在多个方面偏离原始设计最远的系统。我们会指出其他现代版本在设计上的不同取向，并尝试理解背后的原因，以此反映 L4 社区在微内核设计理念上的一致程度。

与其他系统不同的是，seL4 是唯一一个从一开始就为了支持安全性与可靠性形式化验证而设计的系统，同时仍秉承 L4 在最小性、高性能以及支持几乎所有系统架构

方面的传统。这一目标促使我们设计出一种全新的资源管理模型，所有空间资源的分配均由用户态代码显式控制，包括内核内存 [Elkaduwe et al., 2008]。它也是目前文献中第一个具有完整且可靠的最坏情况执行时间（WCET）分析的受保护的操作系统内核 [Blackham et al., 2011]。

另一个值得关注的系统是 Fiasco.OC，它是唯一跨越了 L4 大部分发展历史的代码库，最初是原始 ABI 的一个克隆（甚至并未面向性能优化），但后来逐步发展为一款高性能内核，具备新一代内核的各种特征，包括基于能力的访问控制。Fiasco 也曾作为许多设计探索的测试平台，特别是在实时支持方面 [Härtig and Roitzsch, 2006]。

表 2-2 各种 L4 内核的源代码行（SLOC）

Name	Architecture	Size (kLOC)		
		C/C++	asm	Total
Original	486	0	6.4	6.4
L4/Alpha	Alpha	0	14.2	14.2
L4/MIPS	MIPS64	6.0	4.5	10.5
Hazelnut	x86	10.0	0.8	10.8
Pistachio	x86	22.4	1.4	23.0
L4-embedded	ARMv5	7.6	1.4	9.0
OKL4 3.0	ARMv6	15.0	0.0	15.0
Fiasco.OC	x86	36.2	1.1	37.6
seL4	ARMv6	9.7	0.5	10.2

Table 2: Source lines of code (SLOC) of various L4 kernels.

此外还有两个从零开发的较新设计：专为 x86 平台上的硬件辅助虚拟化而设计的 NOVA [Steinberg and Kauer, 2010]；以及面向 ARM 处理器的高效准虚拟化商业平台 OKL4 Microvisor [Heiser and Leslie, 2010]（简称 OKL4）。

第3章3 微内核设计

Liedtke [1995] 概述了驱动最初 L4 设计的原则与机制。我们将探讨这些原则是如何随着时间演进的，并重点将他们与当前一代微内核的进行比较。

3.1 最小化

Liedtke 设计的主要驱动力是最小性与 IPC 性能，并且他坚信前者有助于后者。具体来说，他提出了微内核最小化原则如下：

仅当将某个概念保留在微内核中是系统实现所必需的，即如果将其移出内核从而允许多种实现方式将无法实现系统所需功能时，才可接受其存在于内核中 [Liedtke, 1995]。

该原则是“仅在内核中保留最基本机制，而不包含任何策略”这一思想的更明确表述，并持续成为 L4 微内核设计的基石。接下来的章节的讨论将展示社区持续努力地去去除功能或将其替换为更通用（且功能更强大）机制的过程。

如表 2 所示，这一原则的坚持可以从源代码规模的对比中看出：尽管系统随着时间推移往往会变得更大，但作为家族中最新成员的 seL4（尽管在设计上与传统模型的偏离最为显著）其代码规模仍基本与早期版本相当。形式化验证对最小化具有特别强的推动力，因为即使是 9,000 行代码的 SLOC 也已经接近可验证性的极限。

尽管如此，迄今为止尚无任何 L4 微内核设计者声称其开发了严格遵循最小化原则的“纯粹”微内核。例如，所有 L4 内核都包含一个调度器，并在内核中实现了一定的

调度策略（通常为固定优先级轮转调度）。目前还没有人提出一种真正通用的内核内调度器，或者一种能够将所有调度策略委托给用户态、且不会带来高开销的可行机制。

3.2 IPC

3.2.1 同步 IPC

最初的 L4 仅支持同步（回合式）IPC 作为唯一的通信、同步与信号传递机制。同步 IPC 避免了内核中的缓冲以及缓冲所附带的管理与拷贝开销，同时也是后续若干实现优化的前提条件，如惰性调度（第 4.2 节）、直接进程切换（第 4.3 节）以及临时映射（第 3.2.2 节）等。

尽管这一机制在概念与实现上都极其简洁且符合最小化原则，我们仍在实践中发现该模型也存在显著缺陷：它将多线程设计强加于原本可以保持简单的系统，进而带来了同步方面的复杂性。例如，缺乏类似 UNIX `select()` 的功能，导致每个中断源都需要一个独立线程，单线程服务器无法同时等待客户端的请求与中断。

在 L4-embedded 中，我们通过引入异步通知机制来解决上述问题，这是一种非常简洁的异步 IPC 形式。我们在 seL4 中进一步发展了这一模型，引入了异步端点（AEP，端点将在第 3.2.3 节中详细介绍）：发送操作是非阻塞的，且对接收者而言是异步的，接收者可以选择阻塞等待或轮询接收消息。逻辑上，异步端点类似于将多个二进制信号量打包进一个字中：每个异步端点具有一个单字的通知字段，发送操作指定一个掩码位（通常为一个比特位），该比特位将与通知字段进行按位“或”操作。等待操作实际上相当于对通知字段的 `select()`。

在我们当前的设计中，一个异步端点可以绑定到某个特定线程。如果线程正在同步端点上等待，而此时有通知到达，该通知将像同步消息一样被交付（但带有标识表明他实际上是一个通知）。

总而言之，seL4 与大多数其他 L4 内核一样，保留了同步 IPC 模型，但通过异步通知进行了增强。OKL4 完全放弃了同步 IPC，转而使用虚拟中断（本质上即异步通知）。NOVA 通过计数信号量机制扩展了同步 IPC [Steinberg and Kauer, 2010]，而 Fiasco.OC 也通过引入虚拟中断增强了同步 IPC。

引入两种 IPC 机制违反了“最小化”原则，因此在这一点上 OKL4 的做法更为“纯粹”（尽管其通道抽象在另一个层面上违背了该原则，详见第 3.2.2 节）。此外，在多核环境中，同步 IPC 的效用也变得更加可疑：类似 RPC 的服务器调用会将客户端与服务器强行串行化，而如果它们运行在不同核心上，这种串行化应尽量避免。因此，我们认为基于异步通知的通信协议将会更加普及，而且未来转向完全异步 IPC 的可能性仍然存在。

3.2.2 IPC 消息结构

最初的 L4 IPC 拥有丰富的语义。除了寄存器内（“短”）消息外，还支持几乎任意大小的消息（包括字对齐的缓冲区以及多个非对齐的字符串）。结合完全同步的设计，该机制避免了冗余拷贝。

寄存器参数支持零拷贝：内核始终在发送方上下文中启动 IPC，然后切换到接收方上下文，期间不触碰消息寄存器。这一机制的缺点在于其具有架构依赖性，且在某些平台（如 x86-32）上支持的零拷贝消息大小较小。实际上，随着 ABI 的演变，可用寄存器数量频繁变化，因为系统调用参数的改变会占用或释放寄存器。

Pistachio 引入了虚拟消息寄存器的概念，其数量在 16 到 64 之间可配置。在实现上，一部分虚拟寄存器映射到物理寄存器，其余部分存储于每个线程的地址空间中一段“固定”的内存区域。该区域不会触发页面错误，因此具有与寄存器相似的语义。内联访问函数隐藏了物理寄存器与内存寄存器之间的差异。seL4 和 Fiasco.OC 延续了这一做法。

此方案的动机主要有两个：其一，虚拟消息寄存器提升了跨架构的可移植性；其二，对于超过物理寄存器数量的中等大小消息，它们的性能开销较小，仅需少量拷贝操作而不需要“长消息”IPC 所需的临时映射。

随着时间推移，上下文切换的架构开销成为 IPC 性能的主导因素，因此寄存器内消息传输的优势逐渐减弱。例如，在 ARM11 上，使用寄存器传输 4 字消息可带来约 10% 的性能提升；而在 Cortex A9 上，该提升下降至约 4%。在 x86-32 上，为消息传输保留寄存器反而会妨碍编译器的优化能力。

“长消息”可在一次 IPC 中指定多个缓冲区以均摊模式切换与上下文切换的开销。在发送方上下文中，内核建立一个临时映射窗口指向接收方地址空间，并将消息复制过去。该过程可能在源或目的地址空间中触发页面错误，进而要求内核处理嵌套异常。内核需调用用户级页面错误处理程序，同时必须模拟异常发生于普通用户态执行过程中。恢复后，需重建原始系统调用上下文。这一流程导致内实现极为复杂，且充满边界情况，容易引入错误。

尽管“长消息 IPC”提供了一些用户态难以仿真的功能，实际上其使用场景极少：共享缓冲区几乎总是更优选择，适用于大多数批量数据传输。同时，异步接口也可实现批量传输，而无需内核提供显式支持。“长消息”主要用于传统 POSIX 的 read/write 接口，需要将任意缓冲区数据传给无法访问客户端内存的服务器。然而，随着 Linux 成为 POSIX 中间件，其与应用程序共享内存，这一用例已被按引用传递

（pass-by-reference）所取代。其他剩余用例也要么接口灵活、要么可通过共享内存实现。

“长消息 IPC”违背了“最小化原则”（此原则关注的是功能，而非性能）。因此，由于内核复杂性以及存在用户态替代方案，我们在 L4-embedded 中移除了长消息 IPC，NOVA 与 Fiasco.OC 也不再提供此功能。对于 seL4，放弃长消息的理由更加充分：seL4 的形式化验证明确避免了任何形式的内核并发 [Klein et al., 2009]，而嵌

套异常将引入并发。同时，这还打破了 C 语言的语义，引入额外控制流。虽然理论上可以对嵌套异常进行形式化推理，但这会显著增加验证工作的难度。当然，内核可以通过增加检查来避免页错误，但这不仅降低了最佳性能，还需要更复杂的形式模型来证明其完整性与正确性。

3.2.3 IPC 目标

最初的 L4 将线程作为 IPC 操作的目标。这种设计的目的是为了避免间接机制而导致的缓存和 TLB 污染，尽管 Liedtke [1993a] 指出，端口可以以大约 12% 的开销实现（主要是额外的 2 次 TLB 未命中）。该模型要求线程 ID 是唯一标识符。这个模型的缺点是信息隐藏性差。

一个多线程服务器必须向客户端暴露其内部结构，以实现客户端负载均衡，或者使用一个网关线程，而这可能成为瓶颈，并引入额外的通信与同步开销。虽然曾提出一些缓解方法，但都存在不足。此外，现代 CPU 对大页的支持降低了由于间接机制带来的 TLB 污染，因为大页更容易使不同实体驻留在同一页中。而且，全局 ID 还引入了隐蔽信道 [Shapiro, 2003]。

受 EROS [Shapiro et al., 1999] 的影响，seL4 和 Fiasco.OC [Lackorzynski and Warg, 2009] 采用了 IPC 端点作为 IPC 的目标。seL4 中的（同步）端点本质上类似端口：挂起发送者或接收者队列的根现在是一个独立的内核对象，而不再是接收线程的线程控制块（TCB）的一部分。与 Mach 的端口不同，IPC 端点不提供任何缓冲。

3.2.4 IPC 超时

阻塞式 IPC 机制可能导致拒绝服务（DOS）攻击。例如，一个恶意（或有缺陷的）客户端可能向服务器发送请求但永远不接收回复；由于是合式 IPC，发送者将会无限期阻塞，除非实现看门狗机制进行中止和重启。L4 的长 IPC 还可能实现更复杂的

攻击方式：恶意客户端可以向服务器发送一条长消息，确保其触发页面错误并阻止其分页器处理该错误。

为了防范此类攻击，原始的 L4 IPC 操作中包含超时机制。具体来说，一次 IPC 系统调用中可以指定 4 个超时：一个用于限制进入发送阶段前的阻塞时间，一个用于接收阶段的阻塞时间，以及两个用于发送和接收阶段页面错误时的阻塞限制（适用于长 IPC）。

超时值使用浮点格式编码，支持零、无限大以及从 1 毫秒到数周的有限值。它们增加了管理唤醒列表的复杂性。

然而，实际上超时机制在防御 DOS 攻击中作用微乎其微。对于非平凡系统，并不存在理论或良好的启发式方法来选择超时时间。实际上，系统中实际上只使用了零和无限两种超时值：客户端使用无限超时进行发送和接收，服务器在接收请求时使用无限超时，在发送回复时使用零超时。（客户端使用一种类似 RPC 的调用操作，即先发送后原子切换到接收阶段，从而保证客户端在服务器回复时处于可接收状态。）传统的看门狗计时器代表了一种更好的方式，用于检测未响应的 IPC 交互（如死锁造成的）。

L4-embedded 已经放弃了长 IPC，因此我们用一个标志位取代了原来的超时机制，该标志位支持轮询（零超时）或阻塞（无限超时）两种选择。仅需两个标志位，分别用于发送和接收阶段。seL4 延用了这一模型。OKL4 的完全异步模型与超时机制不兼容，也不需要它来防御 DOS 攻击。

超时机制还可以通过等待来自一个不存在线程的消息来实现定时休眠，这对实时系统很有用。Dresden 曾尝试扩展模型，包括绝对超时机制，即超时在某个墙上时钟的时间点而不是相对系统调用开始时间到。我们的做法是让用户态访问（物理或虚拟）定时器。

3.2.5 通信控制

在原始的 L4 中，内核会将发送者无法伪造的 ID 传递给接收者。这使得服务器可以实现一种可自由支配的访问控制方式，即忽略不希望接收的消息。然而，恶意客户端仍可以用大量虚假消息（尤其是大消息）轰炸服务器。即使复制操作由内核执行，接收这些消息所耗费的时间也会阻止服务器执行有用的工作，而判断是否丢弃这些消息本身也消耗时间。因此，这种消息构成了一种 DOS 攻击，只有通过内核支持阻止不希望的消息发送才能避免这种攻击 [Liedtke et al., 1997b]。强制访问控制策略同样需要一种机制作为中介并授权通信。

原始的 L4 通过一种叫做“氏族与首领”（clans and chiefs）的机制来实现这一点：进程被组织成“氏族”层级结构，每个氏族有一个指定的“首领”。氏族内部的消息可以自由传递，并由内核保证消息完整性。而跨越氏族边界的消息（无论是发出还是接收）都会被重定向到该氏族的首领，从而由其控制消息流。该机制还支持对不可信子系统的隔离 [Lipner, 1975]。

Liedtke [1995] 认为该模型每次 IPC 操作仅增加两个周期开销，因为氏族 ID 可以被编码进线程 ID 以实现快速比较。然而，这种低开销仅适用于可直接通信的情况。一旦发生重定向，每次重定向都会为一次（逻辑上）单轮 IPC 添加两个消息，造成显著开销。此外，严格的线程层级结构在实践中非常笨拙（大概是构建基于 L4 系统时最令人痛苦的功能之一）。在强制访问控制场景中，该模型最终迅速演变为“每个进程一个首领”。这正是内核强制策略（地址空间层级）限制设计空间的典型例子。

由于这些缺点，许多 L4 实现根本没有实现氏族与首领机制（或在构建时将其禁用），这也意味着无法控制 IPC。后来也曾尝试使用更通用的 IPC 重定向模型 [Jaeger et al., 1999]，但并未能获得广泛采用。最终该问题通过灵活的、基于能力的端点访问控制机制得以解决。

3.3 用户级设备驱动

L4（或者更准确地说，它的前身 L3）最具创新性的理念之一，是将所有设备驱动程序都作为用户态进程来实现。这一设计是其“最小性原则”的核心体现，也成为所有 L4 内核的标志性特征之一。将驱动移出内核的一个重要原因是为了支持形式化验证：一旦将未经验证的代码（如驱动）加入内核，将彻底破坏内核的安全性与正确性保证，而现实系统中庞大的驱动代码量又使得对其进行形式化验证在当前阶段几乎不可能实现。

当然，仍有极少数驱动程序最好保留在内核中。在现代 L4 内核中，这通常是定时器驱动，用于在时间片结束时抢占用户进程、中断控制器驱动，用于将硬件中断安全地分发给用户态驱动程序。

L4 的用户态驱动模型与将中断建模为 IPC 消息这一机制紧密耦合：当发生中断时，内核会向相应的驱动程序发送一条 IPC 消息。这个机制的具体细节（比如虚拟线程发出的 IPC 与上行调用的差异），以及中断消息的绑定与确认协议，多年来曾多次调整，甚至反复变化，但整体思路始终一致。最显著的变化是将中断传递方式从同步 IPC 转为异步 IPC，原因在于同步方式需要模拟内核中的虚拟线程作为中断消息的源头，这在实现上比较复杂，而异步传递大大简化了设计。

用户态驱动程序受益于硬件虚拟化带来的改进。IOMMU 的出现允许用户态驱动程序安全地访问硬件设备。现代网络接口还支持中断合并，进一步降低了中断处理带来的系统开销。在 x86 架构下，TLB 标记等机制大幅降低了上下文切换的成本，也让用户态驱动变得更高效。

如今，用户态驱动程序已经成为主流做法，在 Linux、Windows 和 macOS 等操作系统中都得到了支持，尽管这些系统中并未特别鼓励这种做法。相比之下，由于 L4 拥有高度优化的 IPC 机制，其用户态驱动模型的系统开销要小得多。我们过去也在

Linux 上展示过，即使在较为通用的系统中，只要硬件支持良好，用户态驱动程序依然可以实现极低的系统开销。但在实际应用中，真正对性能有极高要求的设备只是极少数，因此大多数驱动放在用户态对整体系统性能影响并不大。

3.4 资源管理

初代 L4 的资源管理方式，与其通信控制策略类似，严重依赖于进程层级结构。这种方式既适用于进程管理，也适用于虚拟内存管理。层级结构在资源管理和回收方面是一种有效手段，并且能够约束子系统（系统机制限制子进程的权限为其父进程的子集），但其缺点是灵活性差。然而，这种层级结构本质上是一种策略，因此与微内核的理念（如第 3.2.5 节所述）不相匹配。能力机制可以打破层级结构的限制，这也是所有现代 L4 内核都采用基于能力的访问控制的原因之一。以下我们将分析初代 L4 模型中最重要的资源管理问题，以及我们现在是如何处理这些问题的。

3.4.1 进程层级结构

一个进程（在 L4 中本质上就是一个页表加上一些关联线程）会消耗内核资源，而不受控制地分配 TCB 和页表很容易导致拒绝服务攻击。初代 L4 通过进程层级结构来处理这个问题：“任务 ID”本质上是对地址空间的能力，允许创建和删除。它们的数量是有限的（数量级为几千），由内核按先到先得的方式分配。它们可以被委托，但只能沿着层级向上或向下传递。（它们还与用于 IPC 控制的线程层级密切相关，见第 3.2.5 节。）在典型的设置中，初始用户进程会在创建任何子进程之前先占用所有的任务 ID。可以预料，这种模型被证明是僵化且受限的；最终它被完整的能力机制所取代。

3.4.2 递归页映射

初代 L4 将对物理内存帧的权限与现有的页映射绑定在一起。一个进程如果在其地址空间中拥有一个页帧的有效映射，就拥有将该页映射到另一个地址空间的权利。进程也可以将某个页面“授予”出去，这会将该页（以及对其的权限）从授予者那里移除。映射（但不是授予）可以通过 `unmap` 操作撤销。地址空间在创建时空，需要通过映射原语填充。

递归映射模型以原始地址空间 σ_0 为基础，该空间在内核启动后获得所有剩余空闲帧的一对一映射。 σ_0 是所有启动时创建进程的缺页处理器，并会将其页面首次映射给第一个请求该页（通过地址缺页）的进程。

值得注意的是，尽管 L4 的内存模型从每个帧出发构建了一个映射层级，它并不强制地址空间视图是层级化的：映射是通过 IPC 建立的（类似于通过 IPC 消息传递能力），一个进程可以将其某页映射给任何它允许发送 IPC 的进程（前提是接收方同意接收该映射）。与 Mach 不同，L4 没有内存对象语义，只有更接近 Mach 内核内部 `pmap` 接口的底层地址空间管理机制，而不是用户可见的内存对象抽象 [Rashid et al., 1988]。内存对象、写时复制（copy-on-write）和影子链（shadow-chain）都只是用户级构造的抽象或实现方式。

递归映射模型在概念上简洁优雅，Liedtke 对此非常自豪——它在许多论文中都占据重要位置，包括最早的那篇 [Liedtke, 1993a]，以及他所有的演讲中。然而，实践表明它存在显著缺陷。

为了支持页粒度的撤销，递归地址空间模型需要以映射数据库的形式进行大量记录和管理。而且，L4 内存模型的普适性也使得两个串通的进程可以通过将相同的帧递归映射到彼此地址空间中的不同页面来迫使内核消耗大量内存，这是在 64 位硬件上可能发生的拒绝服务攻击，仅能通过控制 IPC（依赖令人诟病的 `clans-and-chiefs`）来阻止。

在 L4-embedded 中，我们移除了递归映射模型，因为我们发现即使在没有恶意进程的情况下，映射数据库就占用了 25–50% 的内核内存。我们用一种更贴近硬件的模型取而代之，其中映射总是来源于物理内存帧的区间。

这种方式的代价是失去了内存的细粒度委托和撤销能力（除非对页表进行暴力扫描），因此我们只将它视为一个过渡性的缓解手段。OKL4 在这个最小模型的基础上做了一些扩展，但并未达到原模型的通用性和细粒度控制。

映射控制在基于能力的系统中非常容易实现，使用标准的 grant-take 模型的变体 [Lipton and Snyder, 1977]。这正是 seL4 的做法：映射权限是通过对物理帧的能力来传递的，而不是通过对以该帧为后备的虚拟页的访问权限来传递的，因此 seL4 的模型不是递归的。即使拥有帧能力，映射仍然严格受限于内核用于记录映射的显式内存模型，如下所述。

Xen 提供了一个有趣的对比点。其 grant table 允许在拥有有效映射的基础上创建一种类似帧能力的机制，可传递给其他域以建立共享映射 [Fraser et al., 2004]。最近有一个扩展 proposal，允许撤销帧 [Ram et al., 2010]。其内存映射原语的语义与 seL4 大致相似，但不包括缺页传播。在 Xen 的场景下，只有在共享时才会承担细粒度委托和撤销所带来的开销。

NOVA 和 Fiasco.OC 仍然保留递归地址空间模型，映射权限由有效映射的持有决定。Fiasco.OC 通过每个任务的内核内存池机制，解决了无法限制映射和记录开销分配的问题。

现有的 L4 地址空间模型（尤其是细粒度委托与撤销）在机制的通用性与最小性之间，以及可能更节省空间的领域特定方法之间做出了不同权衡。

3.4.3 内核内存

虽然能力机制为权限委托提供了一种简洁而优雅模型，但单靠它们无法解决资源管理的问题。一个拥有映射授予权限的恶意线程，仍然可以利用这一点创建大量映射，迫使内核为元数据消耗大量内存，从而可能导致系统遭受拒绝服务（DoS）攻击。L4 内核传统上使用一个固定大小的堆为其数据结构分配内存。初代 L4 使用一个称为 $\sigma 1$ 的内核分页器，允许内核通过它向用户空间请求额外内存。然而，这并不能解决恶意（或存在缺陷的）用户代码强迫系统消耗不合理内存的问题，它只是将问题转移了。因此，大多数 L4 内核并不支持 $\sigma 1$ 。

这个根本问题，也存在于多数操作系统中，其本质是通过共享内核堆导致的用户进程隔离不足。一个令人满意的方案必须能够提供完全的隔离。即使在一个以能力表示权限的系统中，如果存在不受能力系统控制的资源，那么也无法对系统的安全状态进行完整推理。那些将内存仅作为用户级内容缓存来管理的内核，仅能部分缓解这个问题。虽然基于缓存的方法可以避免因内存耗尽引发的 DoS 攻击，但它们并不支持内核内存的严格隔离，而这种隔离是性能隔离或实时系统的前提条件，而且还可能引入隐蔽通道（covert channels）。

Liedtke 等人 [1997b] 研究了这个问题，并提出了每个进程独立内核堆的机制，并辅以在内存耗尽时向内核捐赠额外内存的方案。NOVA、Fiasco 和 OKL4 都采用了这一方案的不同变体。每进程内核堆简化了用户级别的管理（因为用户不再控制分配），代价是失去了在不销毁整个进程的情况下撤销分配的能力，以及无法直接推理已分配内存的状态（只能对其进行限制）。这一权衡仍在学术界持续探索中。

我们在 seL4 中采用了一种根本不同的做法，其内核内存管理模型是 seL4 对操作系统设计的主要贡献之一。我们希望能够对资源使用和隔离进行精确推理，因此我们将所有内核内存（除了内核启动时使用的固定数量内存，包括其严格受限的栈）

都纳入由能力控制的权限系统中。具体而言，我们彻底移除了内核堆，并提供给用户空间一种机制，使其能在内核分配数据结构时识别获得授权的内存。

这一机制的副作用是显著减少了内核的体积和复杂性，这对形式化验证是一个重大优势。关键在于：使所有内核对象都显式化，并受能力机制控制访问。这种设计灵感来源于基于硬件的能力系统，特别是 CAP 系统 [Needham 和 Walker, 1977]，其中硬件解释的能力直接对应内存。HiStar [Zeldovich 等, 2011] 也使所有内核对象显式化，不过它采用的是基于缓存的内存管理方式。

当然，使内核对象对用户可见，并不意味着拥有相关能力的用户可以直接读写该对象。能力仅提供调用该对象特定方法（的子集）的权限，其中包括销毁该对象。（对象一旦创建，其大小就不会再改变。）尤其重要的是，seL4 中的内核对象类型包括“未类型化内存”（Untyped），它可以被用于创建其他对象。对 Untyped 唯一允许的操作是将其部分区域重新类型化为某种对象类型。新对象与原始 Untyped 的关系被记录在能力派生树（capability derivation tree）中，该树还记录其他能力派生形式，如创建权限减少的能力副本。

一旦某块 Untyped 被重新类型化，对应区域的唯一允许操作就是撤销（revoke）所派生的对象（见下文）。在 seL4 中，重新类型化是唯一创建对象的方式。因此，系统可以通过限制对 Untyped 内存的访问来控制资源分配。重新类型化也可以产生更小的 Untyped 对象，这些对象可以独立管理——这是资源管理委托的关键所在。Untyped 的派生过程还确保了内核的完整性属性，即任何两个已类型化对象不会重叠。

表 3 给出了 seL4 所有对象类型及其用途的完整列表。用户空间只能直接访问（加载/存储/获取）映射到其地址空间中的 Frame 所对应的内存（这是通过将 Frame 能力插入页表实现的）。

这个模型具有以下属性：

1. 所有权限都是通过能力显式赋予的；
2. 数据访问和权限可以被限制（confined）；
3. 内核本身（用于其数据结构）也遵循分配给应用程序的权限，包括对物理内存的消耗；
4. 每个内核对象都可以独立于其他对象被回收；
5. 所有操作的执行时间，或可被抢占的时间，都是“短”的（常数时间，或与对象大小线性相关，但不超过一个页）。

属性 1-3 保证了我们能够对系统资源及安全性进行推理。尤其是属性 3，在形式化证明内核可确保完整性、权限限制和保密性时至关重要 [Murray 等, 2013; Sewell 等, 2011]。属性 5 保证了所有内核操作的延迟有界，从而支持其用于硬实时系统 [Blackham 等, 2011]。属性 4 确保了内核完整性。任何持有合适能力的实体都可以随时回收一个对象（从而使原始 Untyped 区域可再次用于创建对象）。例如，在不销毁整个地址空间的情况下，可以单独回收页表所占用的内存。为了实现这一点，内核必须能够检测并使对该对象的所有引用失效。

这一需求通过能力派生树得以满足。撤销对象是通过在派生树上较高位置的 Untyped 对象上调用 `revoke()` 方法实现的；这将移除所有从该 Untyped 派生出的对象的所有能力。当一个对象的最后一个能力被移除时，该对象也会被删除。这将移除它在内核中与其他对象之间的所有依赖，使该内存区域可以重新使用。移除最后一个能力的行为很容易检测，因为这会清除能力树中指向某个特定内存位置的最后一个叶节点。

撤销操作需要用户空间在更高层抽象（如进程）上进行 Untyped 能力与对象的关联管理。这一机制的具体语义，以及 Untyped 与用户级账本的关系，目前仍在研究探索中。

3.4.4 多核

L4 社区很早就开始研究多处理器相关问题。除 L4/Alpha 和 L4/MIPS 外，大多数工作都是在 x86 平台上完成的，这些平台是最早可负担的多处理器系统。早期的 x86 多处理器和多核系统具有较高的核间通信成本，并且没有共享缓存。因此，标准做法是使用每个处理器独立的调度队列（核心之间的内核数据共享最小化），线程迁移只在用户空间显式请求时发生。Uhlig [2005] 探讨了在多核平台上的锁机制、同步与一致性问题，并基于 RCU (Read-Copy-Update) [McKenney 等, 2002] 提出了可扩展的内核数据结构并发控制方法。NOVA 和 Fiasco.OC 广泛采用了 RCU。

随着研究重点从高端服务器转向嵌入式和实时平台，多处理器问题一度被搁置，直到近期嵌入式处理器开始支持多核，这一问题才被重新关注。这类处理器通常具有较低的核间通信成本和共享的 L2 缓存，因此与 x86 平台的取舍不同。在这类平台上，线程迁移的成本通常不足以支撑额外的系统调用开销，因此全局调度器具有更大的意义。

在此背景下，形式化验证引入了新的限制。如第 3.2.2 节所述，并发性给验证带来巨大挑战，因此我们尽可能地将并发排除在内核之外。对于多核系统，这意味着要么采用大内核锁 (big kernel lock)，要么采用多内核结构 (multikernel) [Baumann 等, 2009]。对于微内核而言，由于系统调用执行时间较短，大锁策略并不像看起来那么荒谬，至少对于共享 L2 缓存的少量核心而言，其锁争用是可以接受的。

我们目前正在探索的策略是“簇状多内核” (clustered multikernel) ——这是一种混合方式，在共享 L2 缓存的核心之间使用大锁，在不同簇之间则采取限制性多内核方式（不允许在内核之间迁移内存）[von Tessin, 2012]。簇状多内核避免了大

部分内核代码中的并发问题，因此可以在一定假设条件下继续维持形式化保证。主要的假设包括：（1）锁机制本身以及其之外的代码都是正确且无竞争的；（2）内核对与用户空间共享内存（对于 seL4，仅为虚拟 IPC 消息寄存器块）中的并发修改是健壮的。

表 3-1 seL4 内核对象

Object	Description
<i>TCB</i>	Thread control block
<i>Cnode</i>	Capability storage
<i>Synchronous Endpoint</i>	Port-like rendezvous object for synchronous IPC
<i>Asynchronous Endpoint</i>	Port-like object for asynchronous notification.
<i>Page Directory</i>	Top-level page table for ARM and IA-32 virtual memory
<i>Page Table</i>	Leaf page table for ARM and IA-32 virtual memory
<i>Frame</i>	4 KiB, 64 KiB, 1 MiB and 16 MiB objects that can be mapped by page tables to form virtual memory
<i>Untyped Memory</i>	Power-of-2 region of physical memory from which other kernel objects can be allocated

该方法的优势在于只需做出较小的修改即可保留现有的单核形式化证明。我们已经在形式上完成了多个单核自动机的并行组合，并证明该组合仍然符合系统的行为精化。然而缺点是：形式化保证不再涵盖整个内核，并且 lifting 框架所使用的大步语义（large-step semantics）限制了我们进一步扩展该形式框架，使其无法对锁的正确性、用户-内核并发以及资源迁移放宽等方面进行更细致的形式化推理。

尽管我们不作强烈断言，但某种变体的簇状多内核可能最终会成为实现多处理器微内核完全形式化验证的最佳方式。不过，在形式化方面仍需大量工作，才能对微内核级别的细粒度交错行为进行严密的推理。

第4章 微内核实现

Liedtke [1993a] 总结了一系列设计决策与实现技巧，使得最初 i486 版本中的 IPC 十分高效，尽管其中一些带有过早优化的气息。

其中一些技巧前文已提到，比如（现已废弃的）long IPC 中使用的临时映射窗口。还有一些是无争议的，例如将发送-接收操作组合为一次系统调用（用于类似 RPC 的客户端调用和服务端的“回复-等待”模式）。我们将在下文详细讨论剩余内容，包括一些传统的 L4 实现方法，这些方法虽较少公开讨论，但在社区中早已成为理所当然的默认实践。

4.1 严格的进程导向与虚拟 TCB 数组

最初的 L4 为每个线程分配了独立的内核栈，该栈位于该线程 TCB（线程控制块）之上的同一页内。TCB 的基地址因此可以通过内核栈指针进行掩码运算获得。这种设计只需要一个 TLB 项即可覆盖线程的 TCB 和栈。

此外，所有 TCB 被分配在一个稀疏的虚拟地址数组中，通过线程 ID 进行索引。在 IPC 期间，这允许非常快速地查找目标线程的 TCB，而无需事先验证 ID 的合法性：如果调用者提供了无效 ID，查找操作可能会访问未映射的 TCB，从而触发页错误，由内核处理并终止该 IPC。如果未发生页错误，则通过将调用者提供的 ID 与 TCB 中存储的实际 ID 进行比较来验证其合法性。（原始 L4 的线程 ID 包含版本号，每次线程销毁并重建时版本号都会变化，从而保证线程 ID 在时间上唯一。将当前 ID 存储在 TCB 中可用于检测过时的线程 ID。）

这两种机制都存在代价：大量的内核栈导致每个线程的内存开销较大，也增加了内核的缓存占用。虚拟 TCB 数组则扩大了内核的虚拟内存使用以及 TLB 占用，但避免了额外查找表带来的缓存开销。在仅支持单页大小和非标记 TLB 的处理器上，优化空间有限，只能尽量将数据结构分组以减少页面使用。然而，RISC 处理器通常支持大页（或物理内存寻址）和带标签的 TLB，使得取舍方案有所不同。

当 L4 开始在嵌入式领域获得应用时，内核的内存使用成为一个重大问题，因此设计必须重新审视。

在 Pentium 上采用单栈内核的初步实验显示内核内存使用减少，且在微基准测试中 IPC 性能有所提升 [Haeberlen, 2003]。Warton [2005] 在 ARMv5 处理器上对 Pistachio 进程式内核与基于事件（单栈）内核进行了全面性能评估。他发现两者在微基准中性能相当（相差一般在 1% 以内），但在多任务负载（如 AIM7）下，事件式内核性能高出 20%。尽管事件内核的 TCB 大小是进程内核的两倍多（因为需要存储 continuation），其每线程内存开销却只有进程内核的四分之一。

与此同时，Nourai [2005] 分析了虚拟地址和物理地址方式的 TCB 的权衡。他在 Pistachio（MIPS64 架构）中实现了物理寻址版本。尽管在微基准测试中几乎看不出差异，但在 TLB 压力较大的工作负载下，物理寻址版本表现更优。MIPS 较为特殊，它即便启用了 MMU 也支持物理寻址，而大多数其他架构通常通过幂等大页映射以及“全局”映射来模拟物理寻址。不过 Nourai 的结果足以表明虚拟寻址的 TCB 并无明显性能优势。

事件驱动内核避免了内核中的页错误异常，因而可以保持 C 语言语义。如第 3.2.2 节所讨论，保持在 C 的语义范围内可以简化形式化验证的复杂度。

综合这些结果，我们在 L4embedded 中选择了事件驱动设计与物理地址的内核数据结构，seL4 也采用了这一策略。尽管这一决定最初是出于嵌入式系统资源紧张和形

式化验证的需要，但其优势并不局限于这些场景，我们认为这在现代硬件上是更优的通用方法。

4.2 延迟调度

在会合 IPC 模型中，线程状态频繁在“可运行”与“阻塞”之间切换。这导致频繁的队列操作：将线程加入或移出就绪队列，可能在一个时间片内反复多次发生。

Liedtke 提出了“延迟调度”的技巧以最小化这些队列操作：当线程在执行 IPC 时阻塞，内核只更新其 TCB 中的状态，但不立即从就绪队列中移除，假设它很快就会被唤醒。当因为时间片耗尽而调用调度器时，调度器会遍历就绪队列，直到找到真正处于可运行状态的线程，并将不可运行的线程移除。该方法还辅以对唤醒队列的延迟更新。

延迟调度将原本在高频操作中的工作，延迟到较少发生的调度阶段来处理。然而，我们在分析 seL4 的最坏情况执行时间（WCET）时发现了它的一个问题 [Blackham 等，2012]：调度器的执行时间仅受系统中线程数量的限制！

为了解决这个问题，我们采用了另一种优化方法，称为 Benno 调度，它不会出现上述异常的时序行为：该策略中，就绪队列中包含所有“可运行”的线程（除了当前正在执行的那个）。因此，在 IPC 过程中线程被阻塞或解除阻塞时，通常无需修改就绪队列。在时间片用尽时，内核只需将仍然“可运行但已停止执行”的线程重新插入就绪队列中。由于移除了超时机制，因此也不再需要管理唤醒队列。尽管端点等待队列仍需严格维护，但在常见场景中（例如服务器通过一个端点响应客户端请求），这些数据结构通常处于缓存中，因此操作开销很低。

这种方法具有与延迟调度相似的平均性能，同时也具有可界定的最坏执行时间。

4.3 直接进程切换

L4 传统上尽量避免在 IPC(进程间通信)期间运行调度器。如果一个线程在发起 IPC 调用时被阻塞,内核会切换到一个可以明确识别的可运行线程,并让该线程在原线程的时间片上运行,通常忽略优先级。这种方式称为直接进程切换(direct process switch)。

这种方法比一开始看起来更合理,尤其是在假设服务器线程优先级不低于客户端线程的情况下。一方面,如果一个(客户端)线程发起一个调用操作(到服务器),调用方显然会阻塞直到被回复。既然线程能够执行系统调用,说明它是当前优先级最高的可运行线程,而保持该优先级最好的方式是让被调用方尽快完成操作(而被调用方通常本身也具有较高优先级)。

另一方面,如果服务器对一个等待的客户端执行“回复并等待”(reply-and-wait)操作,并且服务器还有来自其他客户端的请求在等待,那么继续运行服务器线程是合理的,可以利用缓存的预热优势执行 IPC 的接收阶段。

现代的 L4 版本为了保证实时行为的正确性,在直接进程切换符合优先级时保留该机制,否则调用调度器。事实上,直接进程切换是一种时间片捐赠(time-slice donation)的形式,而 Steinberg 等人[2005]表明这可以用于实现优先级继承和优先级上限协议。Fiasco.OC 和 NOVA 通过允许用户按调用设置是否捐赠时间片来支持这一点。

4.4 抢占

传统的 L4 实现会在内核执行期间禁用中断,尽管有些(如 L4/MIPS)在长时间运行的操作中设置了抢占点,在这些点会短暂启用中断。这种做法能显著简化内核实现,因为大多数内核部分无需并发控制,通常能带来更好的平均性能。

然而，原始的 L4 ABI 包含一些运行时间较长的系统调用，早期的 Fiasco 工作将内核做成了完全可抢占，以提升实时性能 [Hohmuth 和 Härtig, 2001]。后来版本的 ABI 移除了大部分长时间运行的操作，Fiasco.OC 又回到了最初的、基本不可抢占的方式。

对于 seL4 来说，不可抢占内核还有一个额外的理由：避免并发使形式化验证更可行 [Klein 等人, 2009]。鉴于 seL4 针对的是许多属于硬实时系统的安全关键系统，我们需要对中断响应延迟有严格的上限。因此必须避免长时间运行的内核操作，并在无法避免的情况下（例如对象删除这种几乎无法界定时间的操作）使用抢占点。我们在抢占点的设置以及最小化其需求的数据结构和算法设计上投入了大量精力 [Blackham 等人, 2012]。

需要注意的是，基于 continuation（延续）的事件内核天然支持抢占点（通过将其作为 continuation 点来实现）。

4.5 不可移植性

Liedtke [1995] 指出，微内核的实现不应追求可移植性，因为硬件抽象层会带来额外开销，并掩盖对特定硬件的优化机会。他还举例说即便是“兼容”的 i486 和 Pentium 处理器之间也有微妙的架构差异，导致最优实现方式出现显著变化。

但这个论点被 Liedtke 自己推翻了，他开发了高性能但可移植的 Hazelnut 内核，尤其是 Pistachio。通过仔细的设计和实现，使得最终版本的实现有 80–90% 与架构无关。

在 seL4 中，架构无关的代码（在 x86 和 ARM 之间）仅占约 50%。大约一半的代码用于虚拟内存管理，这部分本质上就是架构相关的。架构无关代码比例较低的原因是 seL4 整体体积较小，大部分（架构无关的）资源管理代码被移到了用户态。

除了 IPC 快路径，架构相关的优化也很少。Steinberg [2013] 估算，将 NOVA 移植到 ARM 也需重写大约 50% 的代码。

4.6 非标准调用约定

原始的 L4 内核完全由汇编语言实现，因此函数的调用约定在内核内部并不重要。在 ABI 层面，所有非系统调用参数使用的寄存器都被指定为消息寄存器。库接口使用内联汇编生成函数封装，将编译器的调用约定转换为内核的 ABI（希望编译器能优化掉这些转换的开销）。

下一代 L4 内核（从 L4/MIPS 开始）至少部分使用 C 语言实现。在进入 C 代码时，这些内核必须重新建立 C 编译器的调用约定，返回时再切换回内核约定。这使得调用 C 函数代价较高，因此只在处理本身开销较大的操作时才使用 C。

后来的内核几乎完全用 C（Hazelnut）或 C++（Fiasco、Pistachio）编写。调用约定不匹配的成本（以及缺乏 Liedtke 式地“虐待自己”以微优化每一行代码的意愿）导致这些 C 代码的性能比不上旧的汇编内核。因此这些内核的开发者开始引入手写的汇编快速路径，从而获得与原始 L4 可比的 IPC 性能（见表 1）。

这种传统方式并不适用于 seL4，因为其验证框架只能处理 C 代码 [Klein 等人, 2009]，而我们希望尽可能全面地验证内核功能。这要求将汇编代码限制到最小，并避免调用约定转换，迫使我们采用工具链的标准调用约定。

4.7 实现语言

seL4 也极度依赖 fast-path 代码来获得具有竞争力的 IPC 性能，但这些 fast-path 现在必须用 C 语言实现。由于汇编代码的高维护成本，汇编 fast-path 已经在商业的 OKL4 内核中被废弃，因为在商业环境中，这种成本超过了任何性能

上的收益。对于 seL4，我们容忍的 IPC 性能下降不超过 10%，相对于在相同架构上最快的内核。

幸运的是，事实证明，只要精心手工构建 fast-path，就能实现高度竞争力的 IPC 延迟 [Blackham 和 Heiser, 2012]。具体而言，这意味着手动重新排列语句，并利用编译器无法通过静态分析得出的（已验证的）不变量。

实际上，在 ARM11 处理器上，一个单向 IPC 的最终延迟为 188 个周期，比我们在同一硬件上测量到的任何其他内核的最快 IPC 性能还要好 10%。这部分归功于简化的 seL4 ABI 和 IPC 语义，以及事件驱动内核不再需要在堆栈切换时保存和恢复 C 调用约定。此外，我们还得益于改进的编译器，特别是它们对常见情况条件分支进行注解的支持，这有助于提升代码局部性。

无论如何，这一结果表明，汇编实现已不再具有性能优势。德累斯顿的团队事实上发现，即使不使用 fast-path，他们也能实现具有竞争力的性能。

第一个完全用高级语言编写的 L4 内核是 Fiasco，它选择了 C++ 而非 C（后者早些年用于 MIPS 内核的部分实现）。鉴于当时 C++ 编译器的状况，这一决定看似大胆，但至少部分原因是 Fiasco 最初并未以性能为目标。后来情况发生了变化，最近的 Fiasco 经验证明，只要使用得当，C++ 代码不会造成性能损失。

卡尔斯鲁厄团队也为 Pistachio 选择了 C++，主要是为了支持可移植性。尽管德累斯顿和卡尔斯鲁厄对 C++ 抱有高度热情，但我们从未看到 C++ 在微内核实现中提供了任何令人信服的优势。此外，OK Labs 发现，在嵌入式领域中，优秀的 C++ 编译器难以获得，因此他们将微内核版本转回为纯 C。

对于 seL4，形式化验证的需求迫使我们选择 C。尽管德累斯顿的 VFiasco 项目试图验证用 C++ 编写的内核 [Hohmuth 和 Tews, 2005]，但它从未完成对 Fiasco 所

用 C++ 子集语义的形式化。相反，通过用 C 实现 seL4，我们能够基于已有的 C 语言形式化 [Norrish, 1998] 构建，这是实现验证的关键推动力。

第5章 结论

很少有一个研究型操作系统同时拥有庞大的开发者社区、显著的商业部署，并经历了长期的演进。而 L4 正是这样一个系统：其 API、设计与实现原则已演进了 20 年，期间还诞生了十多个从零开发的实现版本。我们认为这是一个极佳的机会，去反思那些经受住时间考验的设计原则与技术诀窍，以及哪些理念因更深入的理解、应用场景的变化和 CPU 架构的演进而被淘汰。

在这期间，许多设计选择与实现技巧曾被采纳，也有不少被放弃，其中一些甚至是最初设计者尤为珍视的。然而，L4 背后最核心的通用原则仍然是最小化设计（包括让设备驱动在用户态运行）和对性能的高度关注。这两个原则至今仍在开发者心中占据首要地位。尤其值得注意的是，作为微内核关键性能指标的 IPC 延迟（进程间通信延迟），在不同的指令集架构和微架构之间对比时，其时钟周期基本保持不变，这与 Ousterhout 在 L4 出现前不久所观察到的性能趋势形成鲜明对比。更令人惊讶的是，L4 的代码规模基本保持稳定，这在软件系统中是极其罕见的现象。

形式化验证的引入进一步强化了“最小化”这一核心理念，并推动了实现的简化。许多设计决策，如简化的消息结构、用户级控制内核内存、多核处理策略等，都深受可验证性的影响。同时，它也改变了一些实现方法，例如使用事件驱动的内核结构、采用标准的调用约定、选择 C 语言作为实现语言。然而，我们并不认为这些为了可验证性而做出的取舍在忽略验证需求的前提下是更差的选择。

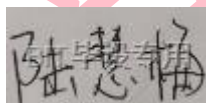
正是通过形式化验证，L4 兑现了微内核当年最核心的承诺之一：健壮性。我们认为，L4 在保持其初衷的同时实现了这一目标，这本身就是对 Liedtke 原始设计思想的最

大肯定。这也许花费了非常漫长的时间，但最终，Brinch Hansen 在 1970 年提出的那些曾被视为激进的想法被时间所证实。

不过，有一个领域至今仍未找到令人满意的抽象方式：时间。当前的 L4 内核仍然实现了一种具体的调度策略——通常是基于优先级的轮转调度，这是目前内核中唯一剩下的策略性内容，也正是 L4 通用性受限的最大因素。然而，Dresden 和 NICTA 的研究表明，一个可参数化的通用调度器有可能支持所有标准的调度策略。我们相信，这一次不需要再等上 20 年了。

指导教师意见：符合要求！

指导教师签字：



日期：2025 年 4 月 17 日