

B: A handler pointer is modified to point to an interposer which does work before continuing to the original handler. C: Here the original handler has been entirely replaced with a custom handler.



## Introducing the Chrono-kernel: Daniel Bristol de Olivera Kernel Privilege for the People

Container Orchestrator

**Problem:** General purpose OSs like Linux meet the needs of many applications, yet some with extreme requirements (e.g. latency, throughput, security, etc.) slip through the cracks. Unfortunately, these apps are often the most valuable (e.g. high

**Deployment:** kElevate has wide applicability, it runs baremetal, in containers, and in virtualization. It can be invoked from any language via the syscall interface.

We have retrofitted Linux with the kElevate mechanism in an x86\_64 prototype, and we have demonstrated the same technique is viable on an ARM64 proof of concept.

**Context:** Like Kernel Modules, elevated threads are unlimited in their power to change the system, so they are similarly protected behind superuser access. Unlike Kernel Modules, elevated threads are built to run from the application context.

Elevated applications use standard build processes. ABI compatibility may be maintained if desired. When running in kernel privilege, an elevated thread can make structural changes to Linux. It can interpose on system structures like the interrupt descriptor table, modifying or replacing core system components and code.

execution.

A Use Case: When running in kernel privilege, an elevated thread can exploit the low level internal interfaces of the kernel, in addition to the usual syscall API, utilizing the kernel as a library. We demonstrate using this to shortcut standard syscall paths as well as to perform more aggressive shortcutting deep in the kernel.

Our Linux prototype uses the kElevate mechanism to turn Linux into a Chrono-kernel. At any point, an application can access the supervisor mode on the thread granularity. Any subset of threads in the system can be concurrently elevated. In these diagrams, elevated threads are marked with a badge.

こここ 333

privileged execution is exposed, you can script with kernel power. Now, changes that used to go into the kernel can instead be executables living in the filesystem.

You and the Chrono-kernel!

We have many more tasks than we have developers to implement them! We are looking for help from both types of developers: application and kernel. We are looking for applications with clear kernel bound latency/throughput/ real-time optimization objectives and kernel hackers who would be interested in helping us get through our to-do lists or bring their own perspectives to the work. If this sounds interesting to you please contact Tommy (tommyu@bu.edu).

## Contact: tommyu@bu.edu



licrobenchmark Experiments

Getting a Linux process's parent pid getppid() is one of the shortest syscalls. It approximates the total HW and SW costs of the syscall path. We call it in a loop and box plot the individual latencies.

"Linux" is a completely standard Fedora 35 environment.

the application in supervisor mode.

"Chrono-kernel" is the same environment plus the kElevate mechanism simply showing a comparable result; the system call path has not been

Approach: Instead of resorting to kernel bypass or custom operating system construction, the goal of this project is to marry the convenience of application programming on Linux with the full power of the kernel's mode of

A Chrono-kernel is a system that offers first-class support for application threads to access supervisor execution mode. We retrofit Linux with a new mechanism, kElevate, turning it into a Chrono-kernel. kElevate allows any set of application threads to toggle into the supervisor mode of hardware execution in  $\sim 10$  nanoseconds. We call these "elevated" threads.

> Our implementation introduces kElevate as a new system call. Interposition can be used to further optimize access to this primitive. For the X86 64 prototype, this involved a total of 357 LoC changes. We are able to keep changes to the kernel minimal because once the primitive of

because a slow **Iret** is used instead of the faster **sysret** instruction. We use **Iret** out of convenience to preserve existing kernel exit code; when higher performance is desired, interposition can be used to change the **lret** into a ret making elevated syscalls even faster than the Linux baseline. The time savings on the ShortCut case are mostly due to skipping system software on the ingress and egress paths, hardware overheads are relatively low. read() Latency write() Latency Linux Chrono-kern Chrono-kerne



B & C Small reads and writes can be accelerated by about 35% and 20% respectively. The relative savings decreases as a function of number of

bytes because cpu bound work takes up a relatively larger fraction of execution time. ShortCut makes use of the access to the kernel internal interfaces ksys\_read() and ksys\_write() respectively (see below).

The receive system call is a bit noisier than the others on most configurations. Interestingly not only does the ShortCut case improve on latency, the standard deviation is much smaller, which has positive implications on using kElevate in real-time contexts.

## Macrobenchmark Not Yet Peer Reviewed

		cop_seriariog.
	tcp_r	tcp_sendmsg
	inet s	sock_sendmsg
	sock_rea	sock_write_iter
	new_sync_r.	new_sync_write
	vfs_read	vfs_write
	ksys_read	ksys_write
	do_syscall_64	do_syscall_64
	entry_SYSCALL_64_a	[ entry_SYSCALL_64
read		write

Back traces from profiling Redis. In blue: skipping the syscall entry and exit code. In green: skipping the VFS and socket disambiguation

	Average Throughput (MB/s)	Percent Throughput Improvement wrt Linux
Linux	4.87	_
ksys_{read/write}	5.50	12.9%
tcp_{send/recv}	6.45	32.4%

We show baremetal latency plots demonstrating that a server using the kElevate mechanism improves significantly when using shortcutting directly into the ksys\_read and ksys\_write handlers. Further, we provide results for making shortcuts much deeper into the kernel, at tcp\_sendmsg and tcp\_recvmsg, cutting out the trip through the VFS layer to disambiguate the file descriptor.

Redis Throughput

Redis is an open source key value store, often used for in memory caching for microservices, it is one of the most popular applications by Docker's rankings. Redis makes read and write syscalls for its network traffic. The Redis server shortcuts significantly improve the throughput (by a third) and latency (a quarter) it achieves. A datacenter provider typically controls the software of both the Redis client and the server, so they have the opportunity to shortcut on both sides of the communication. We are currently quantifying the **power** savings as well.

**Ongoing work:** We are applying the same set of shortcuts to Memcached and Apache. This will demonstrate that our approach is general enough to be reused. We are interested in building tools that application developers can incorporate even if they are not kernel experts. These shortcuts may be used even when the application source code is not available (ABI compatible). Finally, we are also developing microkernel-like servers which can interface with unprivileged applications, securely accelerating their workloads.

Thomas Unger, Arlo Albelli, Ryan Sullivan, Orran Krieger,

Jonathan Appavoo



BOSTON

UNIVERSITY

## **Discussion:** syscalls take longer in the elevated case



0.2

Average Latency

(msec)

1.70

1.51

1.28

0.8

**Percent Latency** 

Improvement wrt Linux

11.2%

24.7%

0.6

Percentile