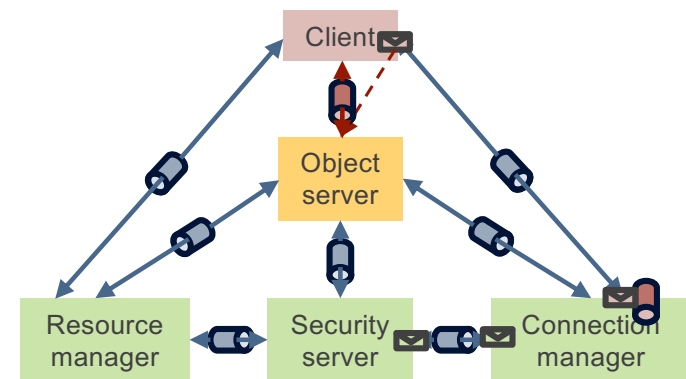


School of Computer Science & Engineering
COMP9242 Advanced Operating Systems

2025 T3 Week 10 Part 2

seL4 Research at TS@UNSW
@GernotHeiser



Copyright Notice

These slides are distributed under the Creative Commons Attribution 4.0 International (CC BY 4.0) License

- You are free:
 - to share—to copy, distribute and transmit the work
 - to remix—to adapt the work
- under the following conditions:
 - **Attribution:** You must attribute the work (but not in any way that suggests that the author endorses you or your use of the work) as follows:

“Courtesy of Gernot Heiser, UNSW Sydney”

The complete license text can be found at
<http://creativecommons.org/licenses/by/4.0/legalcode>

Today's Lecture

- More on LionsOS performance
- Secure General-purpose OS
- Pancake: System language for verified systems
- Time Protection: Prevention of μ arch timing channels

More ob LionsOS

Fast – secure – adaptable!



Microkernel Overheads

Chen et al, OSDI'24

High syscall rate = 61k/s

seL4 round-trip address-space switch = 1k cy

Assume average 2 R-T AS switches / syscall:

$$\text{Switch O/H} = 2 \times 61\text{k}/\text{k} \times 1\text{kcy} = 122\text{M cy/s}$$

Assume 3GHz clock:

$$\text{O/H} = 122\text{M cy/s} / 3\text{Gcy/s} = 122/3\text{k} = 4\%$$

Assume 4-core CPU:

$$\text{O/H} = 4\%/4 = 1\% \text{ of CPU!}$$

Assume Linux max CPU load = 25%

$$\text{relative O/H} = 4 \times 1\% = 4\%$$

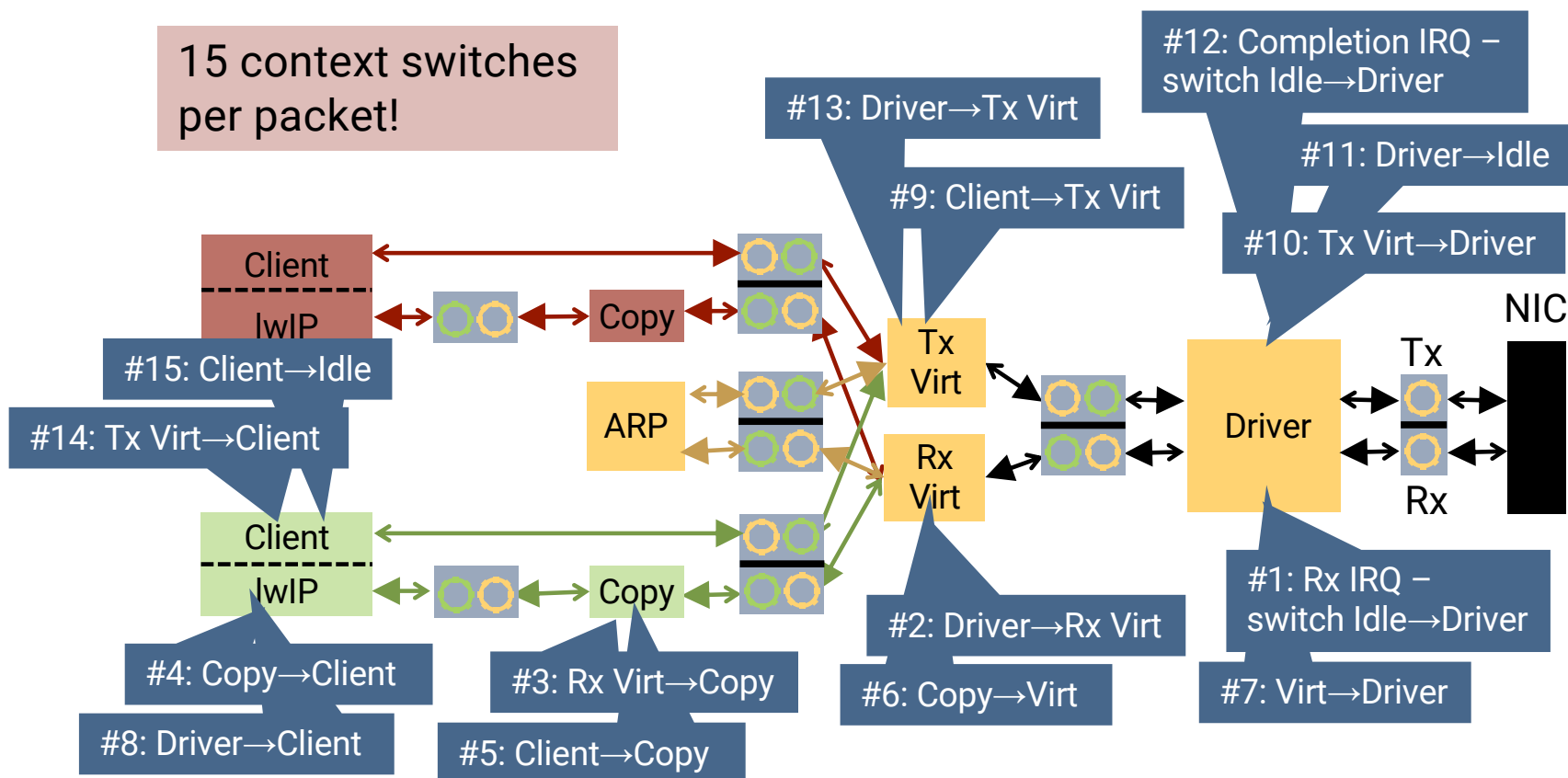
Conservative
IMHO

Why would
anyone care?



Packet Round-Trip Context Switches

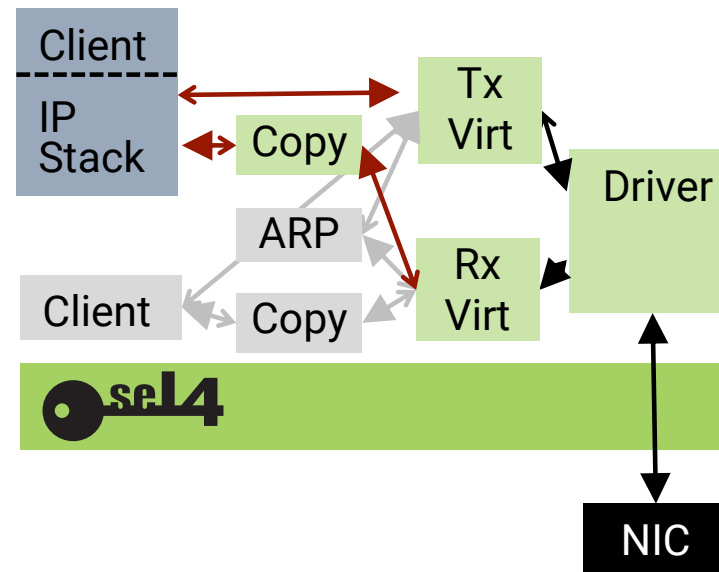
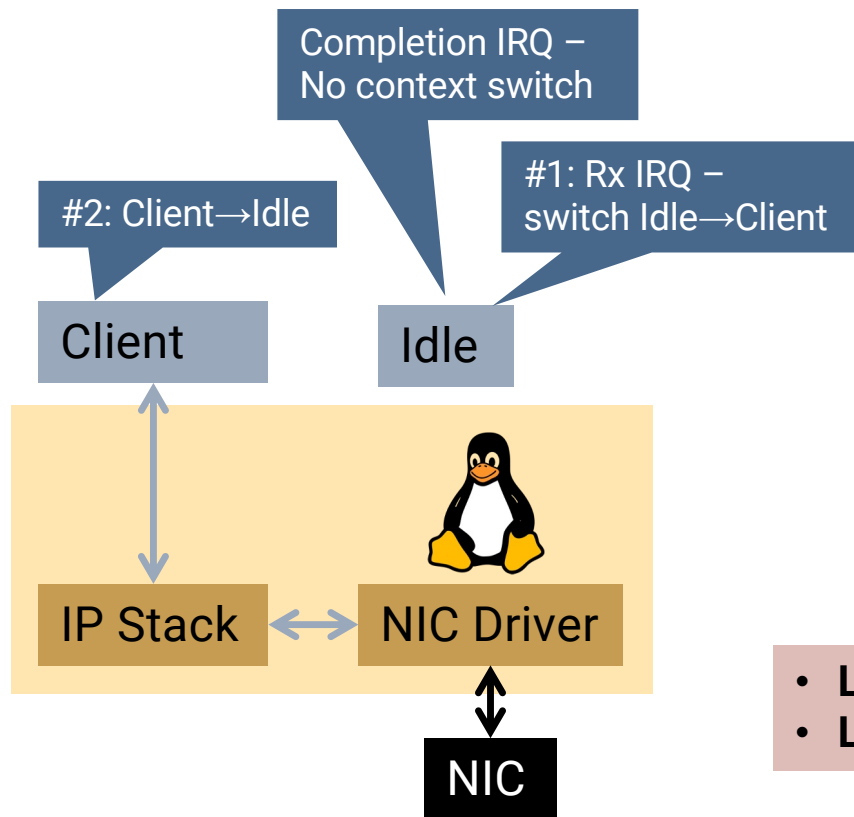
15 context switches per packet!



Must return free buffers!



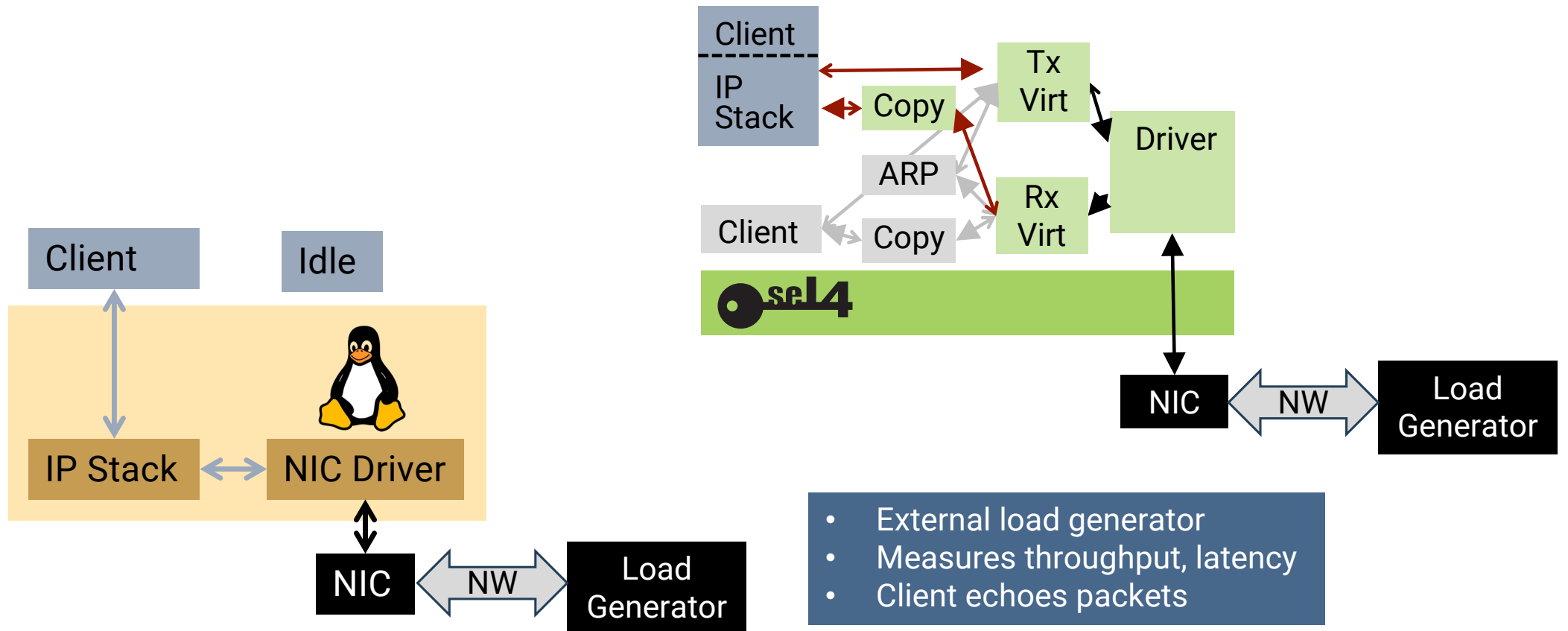
Comparing to Linux



- **LionsOS**: 30 mode switches, 15 context switches
- **Linux**: 6 mode switches, 2 context switches



Comparing Performance: Setup



What Do We Expect?

Ethernet packet size = 1.5kB

Assume Linux mode switch = half context switch

LionsOS O/H = 12/pk \times 0.5k cy = 6k cy/pkt

Max packet rate for 1Gb/s NIC:

$$\text{rate} = 1\text{Gb/s} / 1.5\text{kB} = 1\text{Gb/s} / 12\text{kb} = 833\text{k/s}$$

Worst-case O/H for 1Gb/s NIC:

$$\text{O/H} = 6\text{k cy/pkt} * 833\text{k pkt/s} = 0.5\text{G cy/s}$$

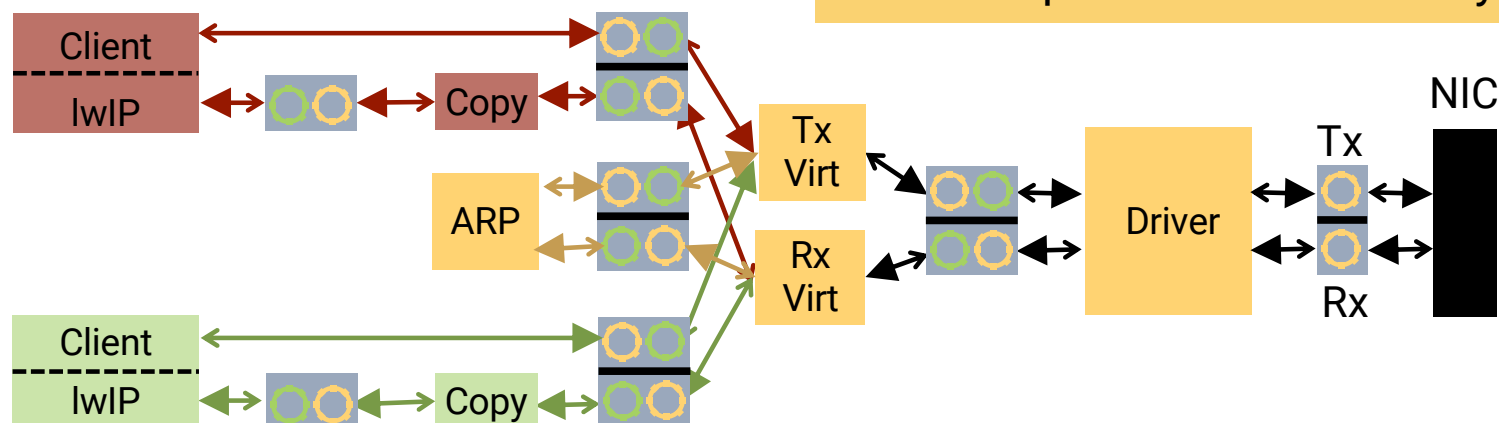
Assume 3GHz clock:

$$\text{rel O/H} = 0.5\text{G cy/s} / 3\text{G cy/s} = 17\% \text{ of core}$$



However, There's Batching

- Each component will process everything in its queue before signalling another component
- No component will ever busy-poll!



- Dramatically reduces context switches under load!
- Measure 5–10 pkt/IRQ!

What Do We Expect?

Ethernet packet size = 1.5kB

Assume Linux mode switch = half context switch

LionsOS O/H = $12/\text{pk} \times 0.5\text{k cy} = 6\text{k cy/pkt}$

Max packet rate for 1Gb/s NIC:

$$\text{rate} = 1\text{Gb/s} / 1.5\text{kB} = 1\text{Gb/s} / 12\text{kb} = 833\text{k/s}$$

Worst-case O/H for 1Gb/s NIC:

$$\text{O/H} = 6\text{k cy/pkt} * 833\text{k pkt/s} = 0.5\text{G cy/s}$$

Assume 3GHz clock:

$$\text{rel O/H} = 0.5\text{G cy/s} / 3\text{G cy/s} = 17\% \text{ of core}$$

At 100Mb/s, packet spacing = $1/(83\text{k/s}) = 12\mu\text{s}$

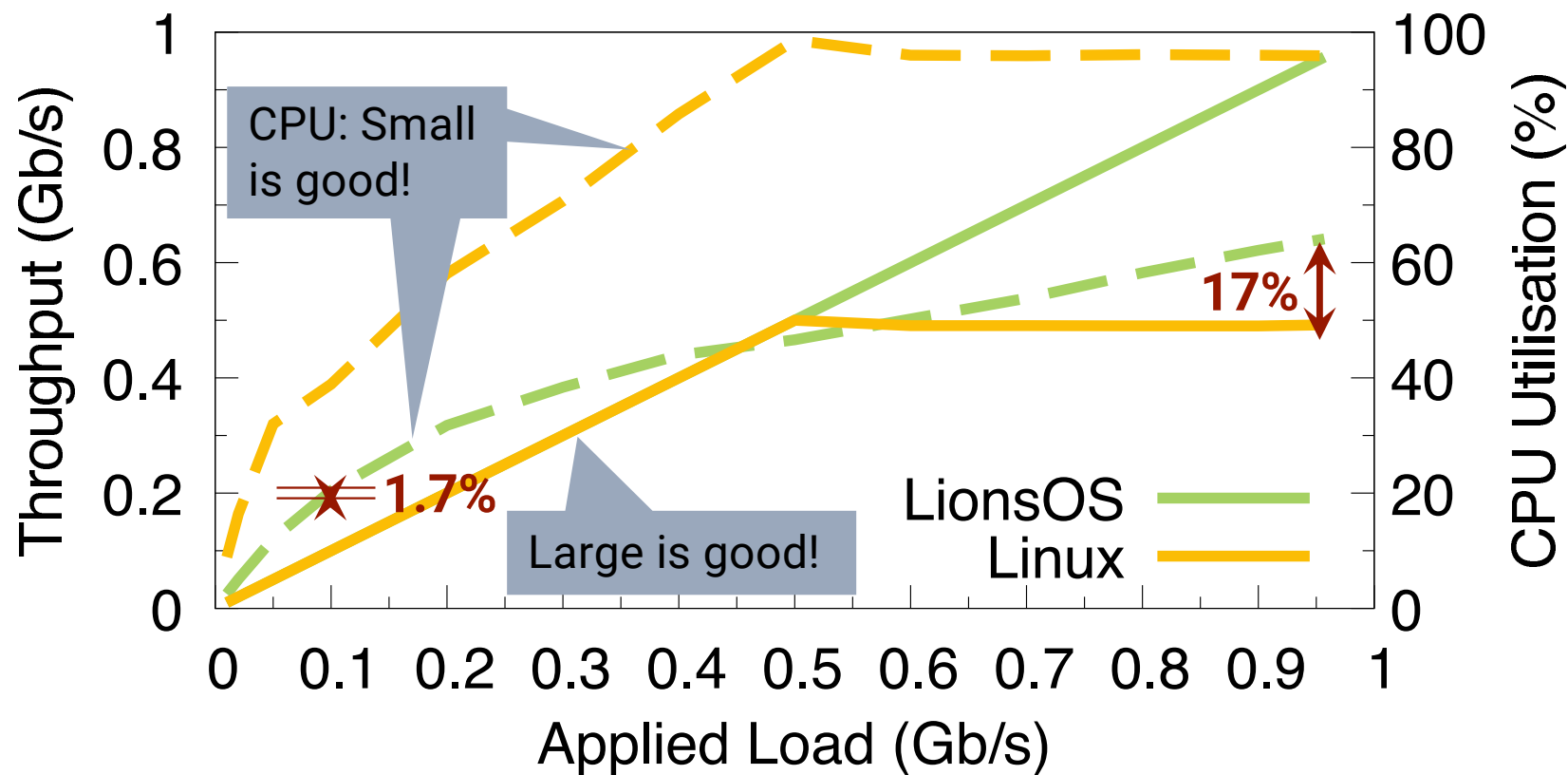
$$\text{rel O/H} = 17\%/10 = 1.7\% \text{ of core}$$

Highly pessimistic
due to natural
batching!

Avoid batching by
spacing packets!



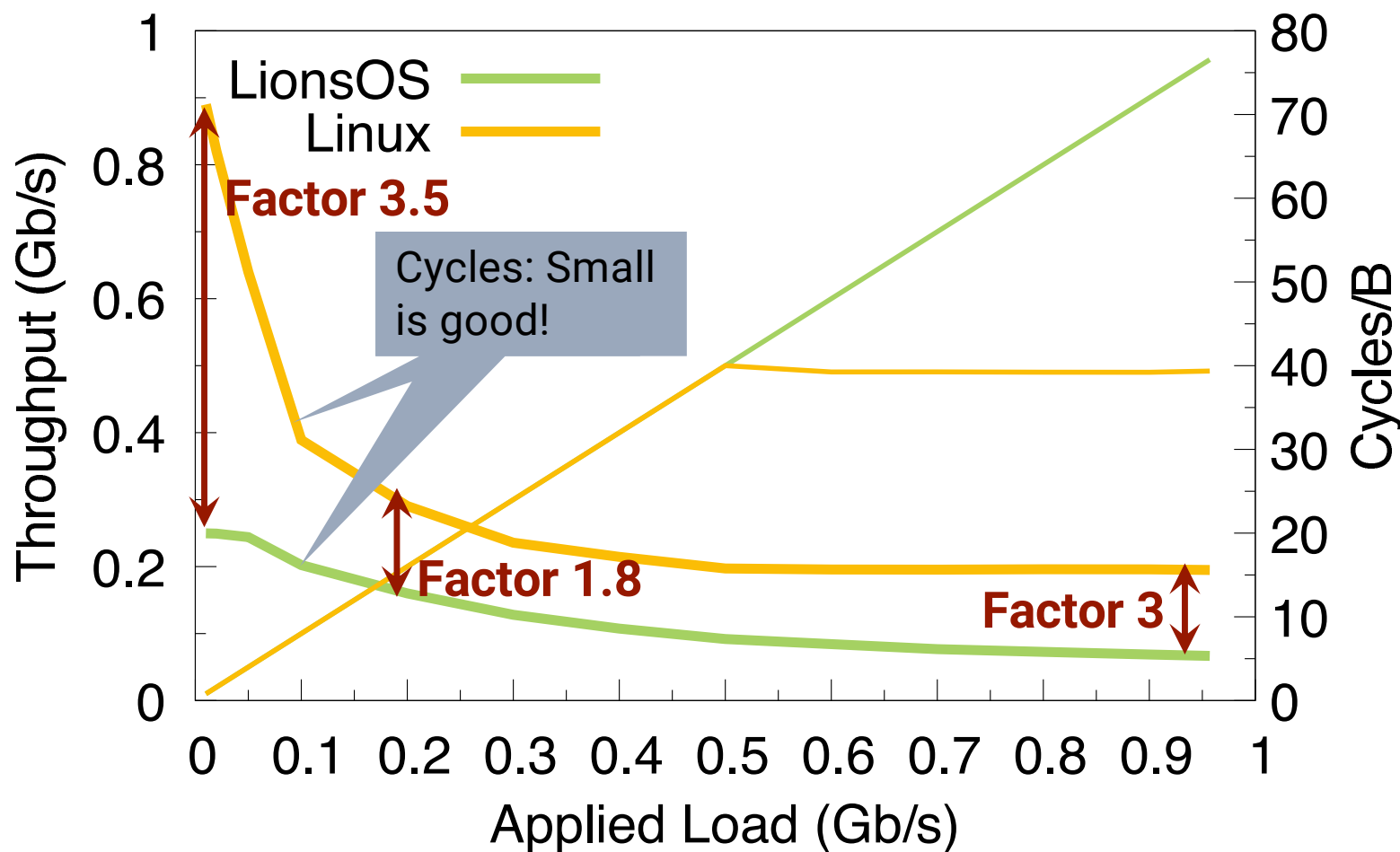
Performance: i.MX8M, 1Gb/s Eth, UDP



Single-core configuration

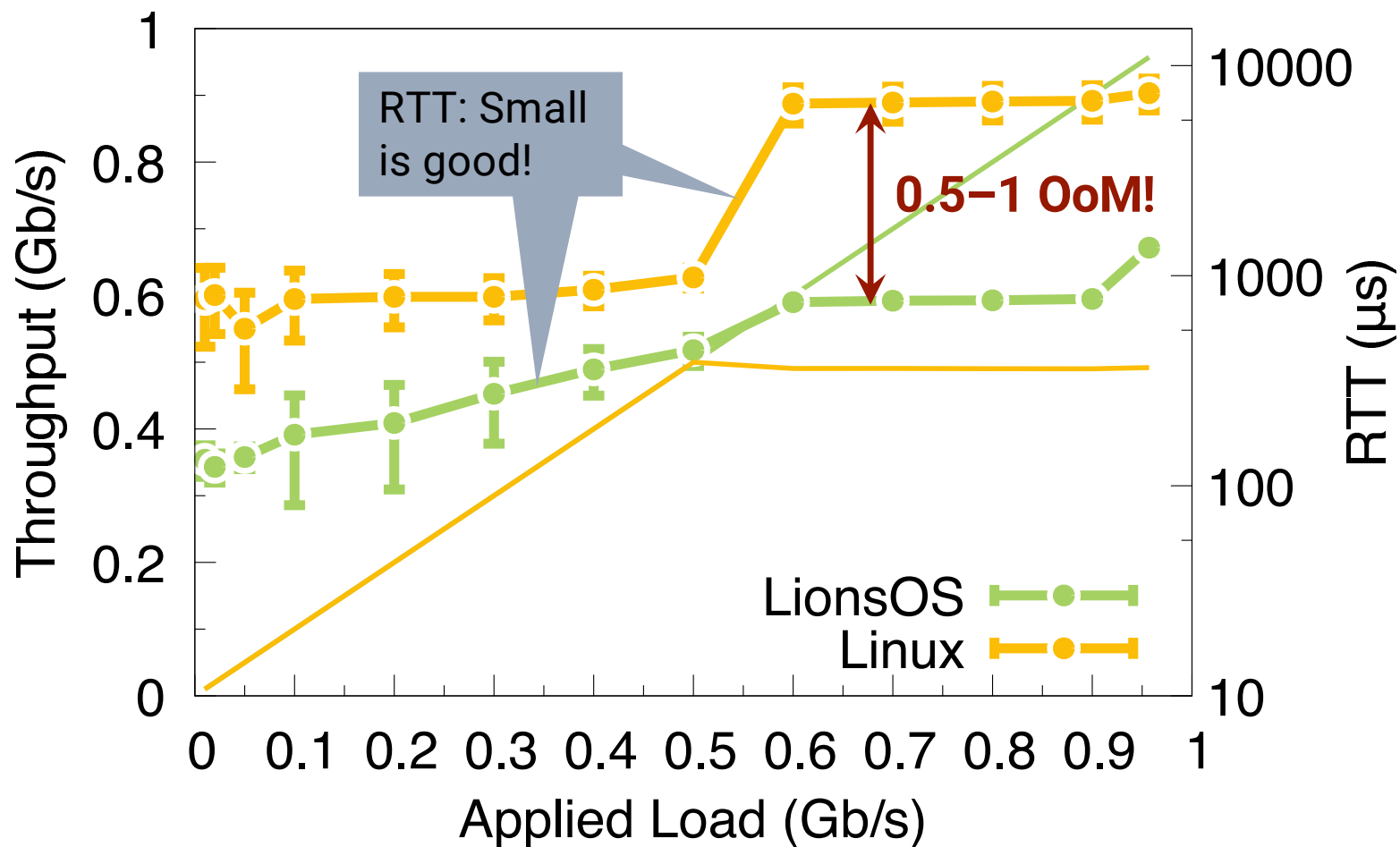


Performance: Processing Cost per Byte



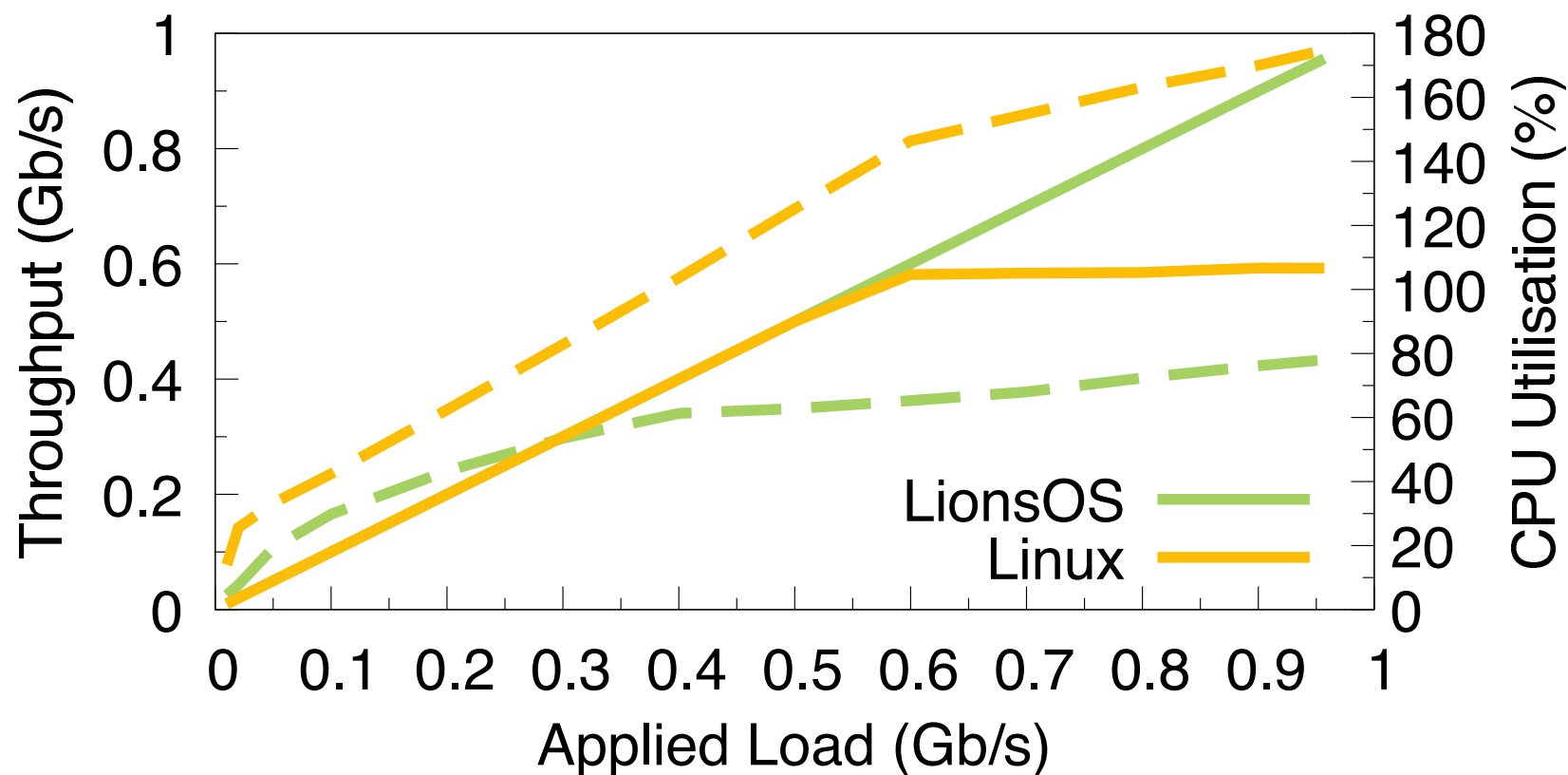


Performance: Round-Trip Times





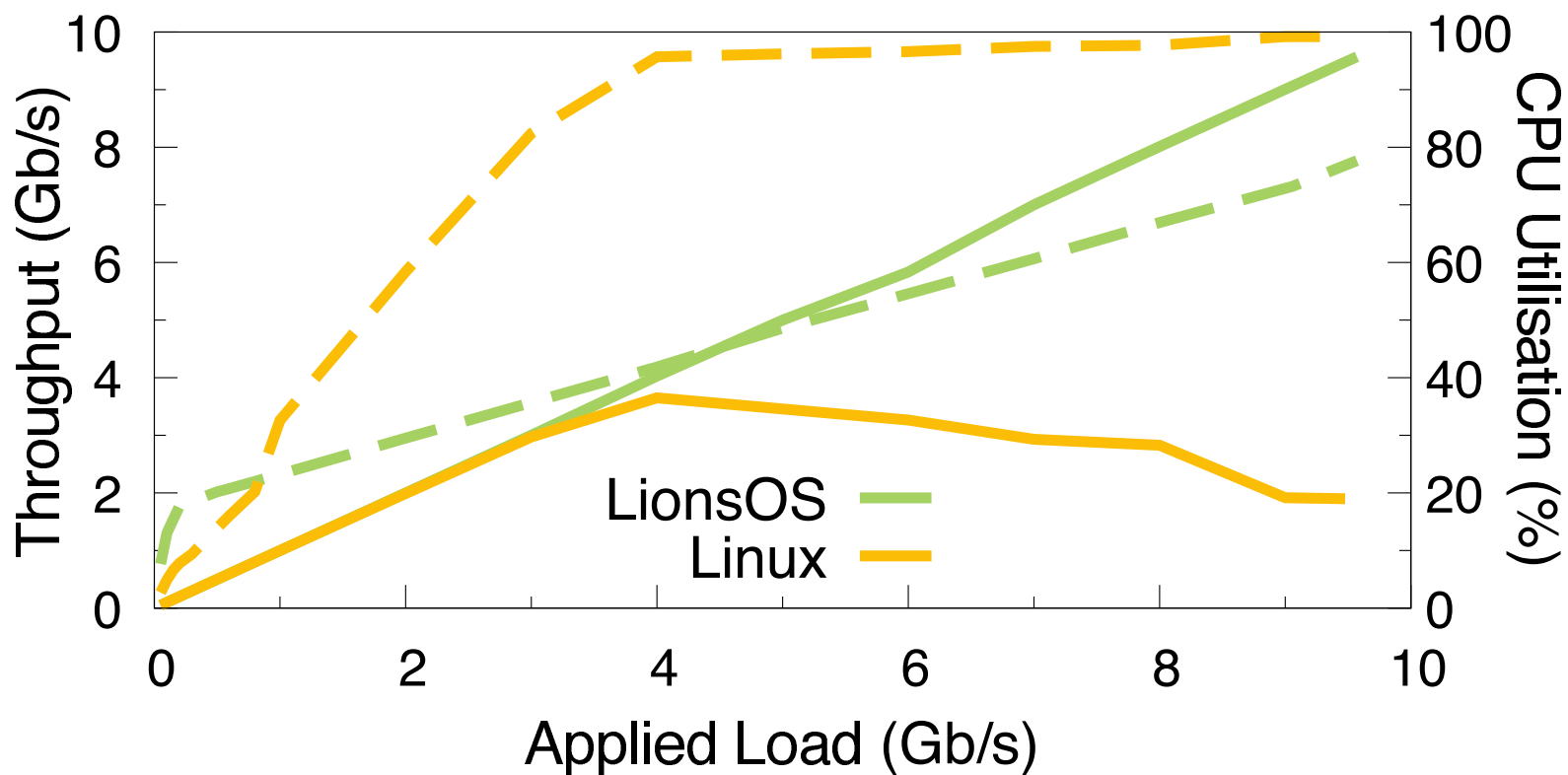
Performance: i.MX8M, 1Gb/s Eth, UDP



Multicore configuration



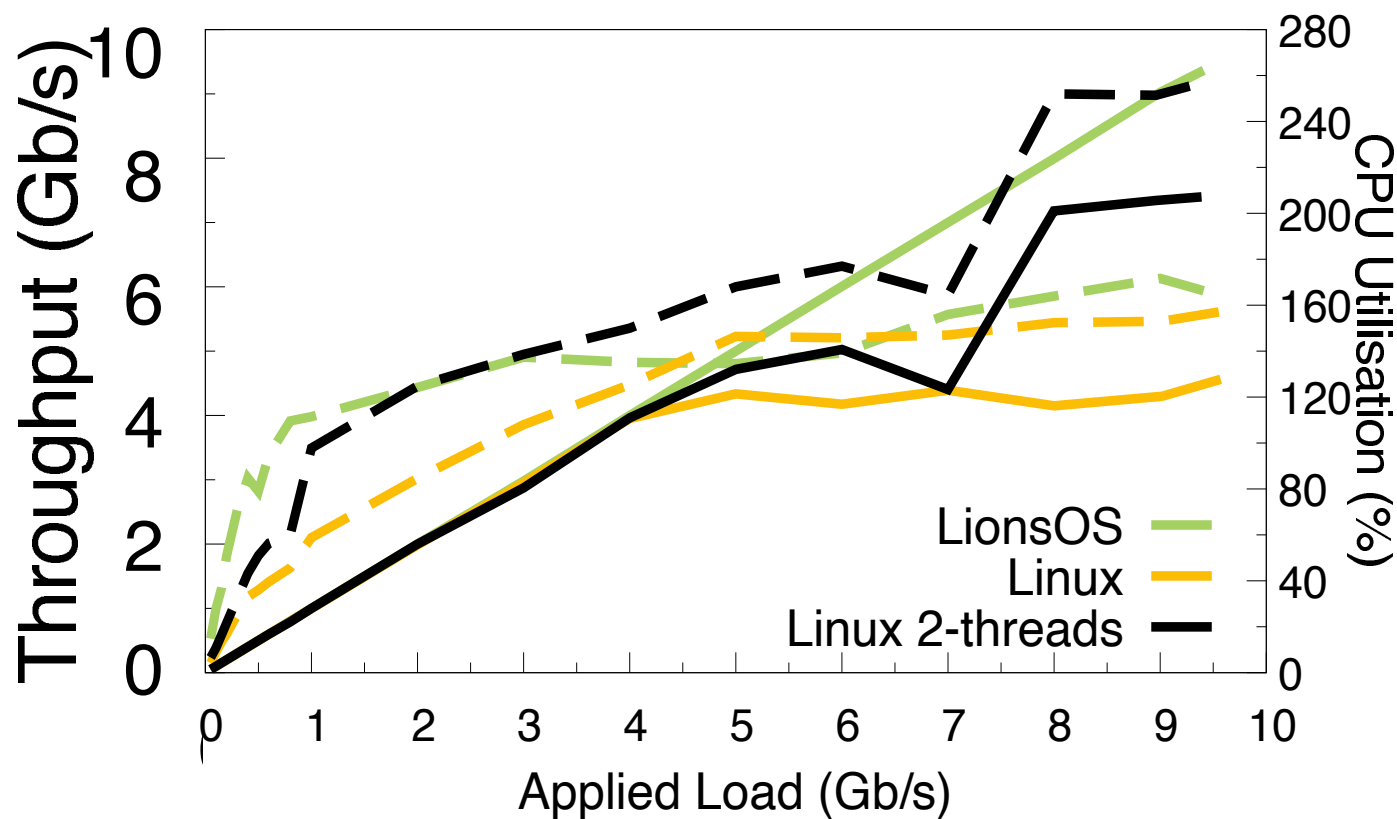
Performance: x86, 10Gb/s Eth, UDP



Single-core configuration



Performance: x86, 10Gb/s Eth, UDP

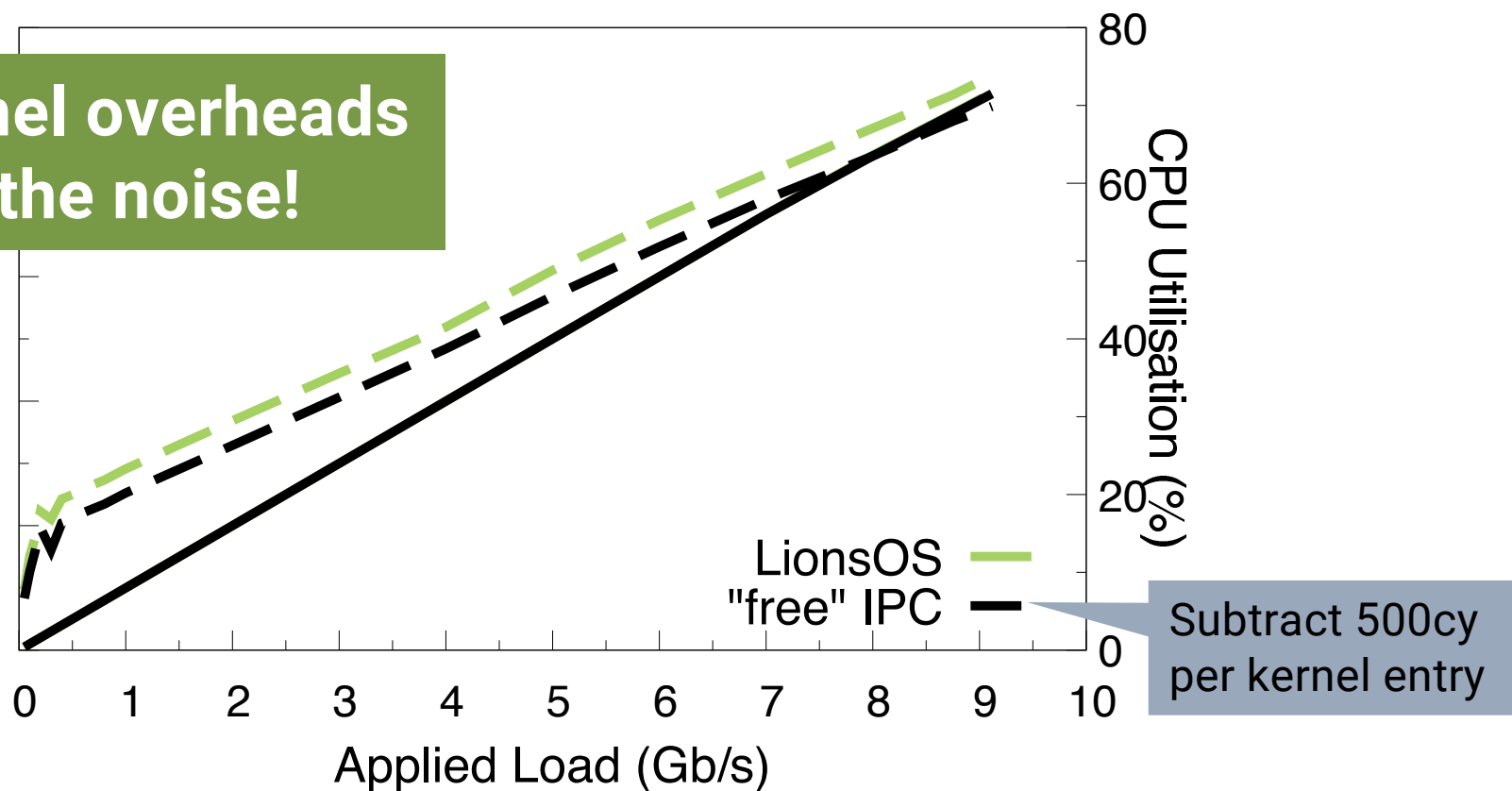


Multicore configuration



Syscall cost simulations (x86)

Microkernel overheads
are in the noise!



Single-core configuration



Why Is LionsOS Faster Than Linux?

Linux:

- NW driver: 3k lines
- NW system total: 1M lines

LionsOS executes less code!

LionsOS:

- NW driver: 400 lines
- Virtualiser: 160 lines
- Copier: 80 lines
- IP stack: much simpler, client library
- shared NW system total < 1,000 lines



Why Is LionsOS So Simple?

Provide **exactly** the functionality needed, not more

Simple programming model:

- strictly sequential code (Microkit)
- event-based (Microkit)
- single-producer, single-consumer queues
- location transparency
- ...

Static **architecture**, mostly static resource management

LionsOS Status

- Runs sel4.systems web site
- LionsOS-based firewall released as a community project
- Doom demo in TS lab 😊
- Runs on Arm, RISC-V, x86
- Native networking, storage, I2C, SPI, ...
- Driver VMs for audio, 2G graphics, ...
- Verification in progress



But I Want A **Real** OS!

Cost Of A Dynamic OS

- More complexity, larger code size

Might affect cache footprint?

- Double book-keeping, multiple server invocations

IPC overheads in the noise

- Higher startup times due to dynamic resource allocation

fork() will be the test!

- Resource revocation may require indirection

Do we need them?

seL4 caps can be revoked without

- “Universal” policies are complex & costly

Do We Need “Universal” Policies?

Claim:

- Systems rarely change policies on-the-fly
- Can change policy by replacing policy module

Keep configuration complexity off-line!

LionsOS experiment:

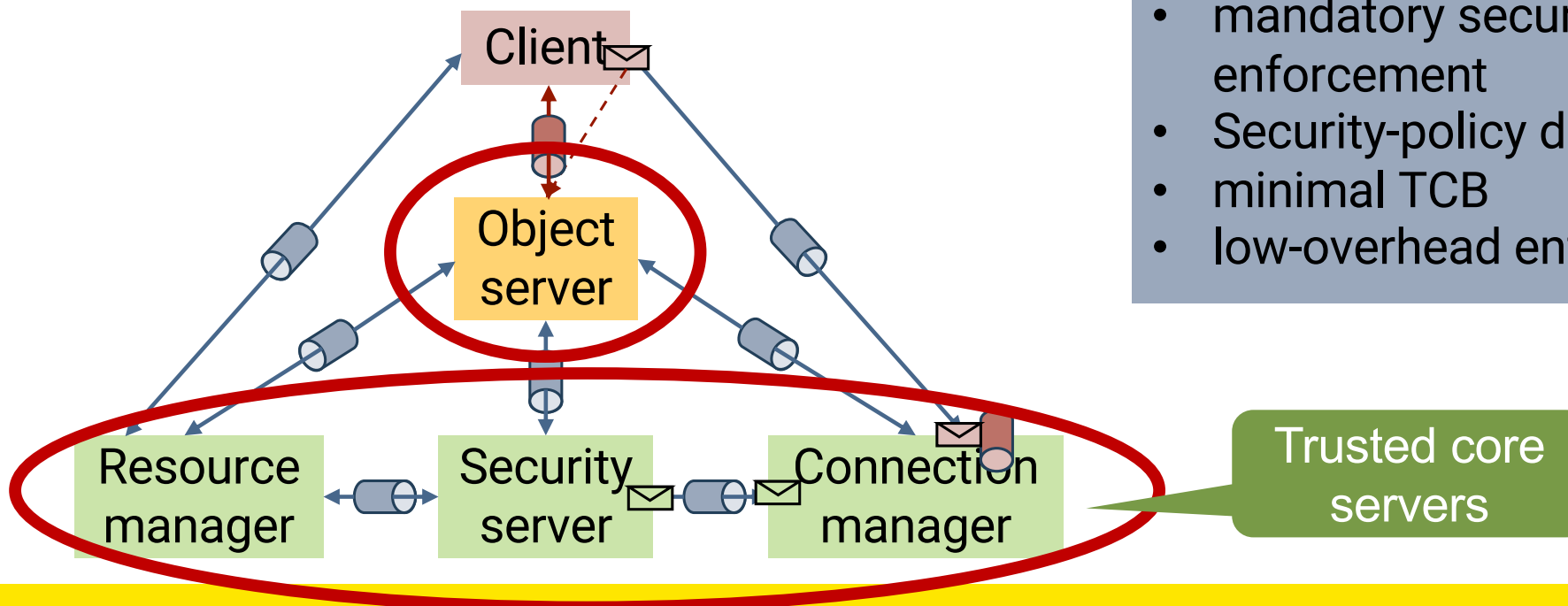
- Reload component with new policy implementation
- Cost: **17 μ s** on i.MX8M

Djawula: PoC Of General-Purpose OS

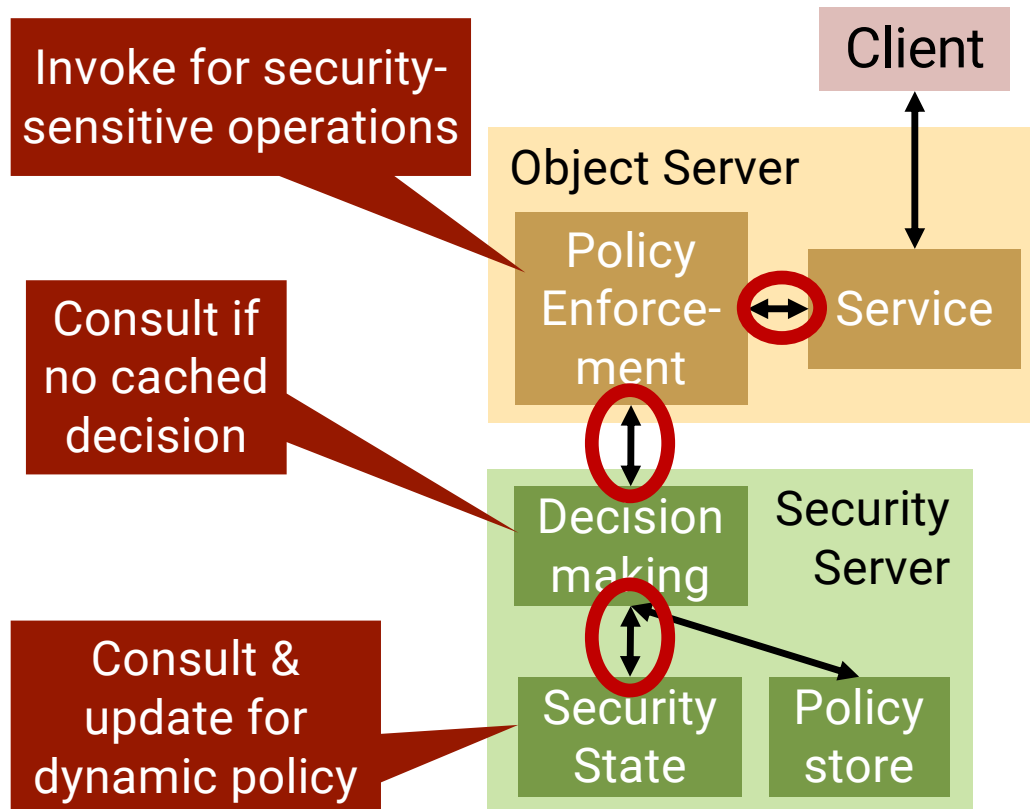
Aim: General-purpose OS that **provably** enforces a general security policy

Requires:

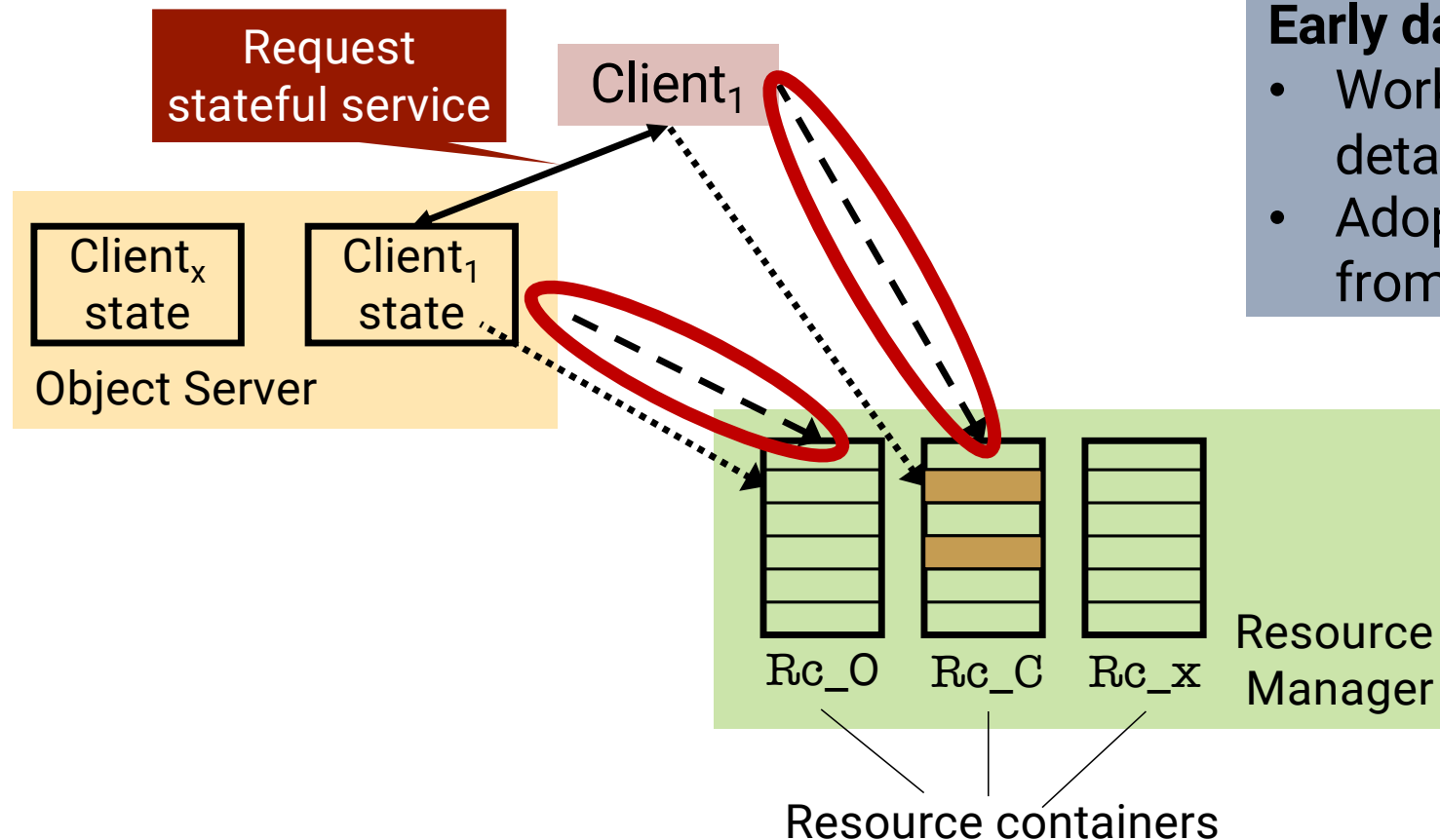
- mandatory security-policy enforcement
- Security-policy diversity
- minimal TCB
- low-overhead enforcement



Core Ideas: Dynamic Enforcement



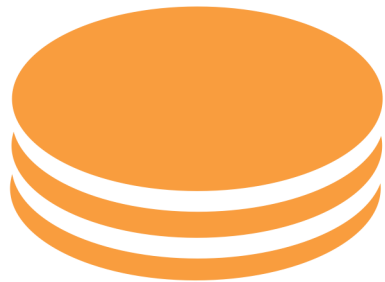
Core Ideas: Resource Donation



Early days:

- Working on framework, details of model
- Adopt components from LionsOS

Scaling Verification



PANCAKE

A Language for
Verified Systems Programming

Driver Dilemma

seL4 is one-off,
justifies cost

High seL4 verification
costs partially due to
C language

Better language
would reduce cost

Drivers are
commodity,
must be cheap!

Drivers are low-
level, need C-like
language

Lions OS

Idea:

1. Simplify drivers
2. Design verification-friendly systems language: Pancake
3. Automate (part of) verification

- Verified compiler
- de-compilation

- Well-defined semantics
- Memory-safe

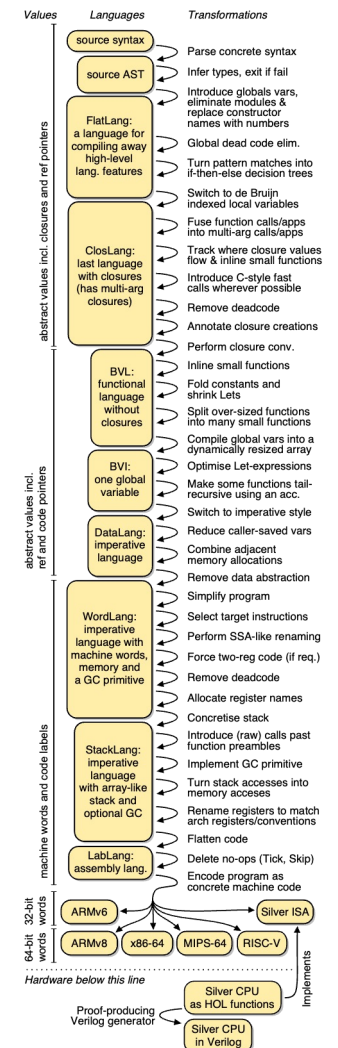
CakeML: Verified Implementation of ML



Re-use framework
for new systems
language: Pancake

- ✓ Mature functional language
- ✓ Large and active ecosystem of developers and users
- ✓ Code generation from abstract specs
- ❑ Managed \Rightarrow not suitable for systems code
- ✓ Used for verified application code

<https://cakeml.org>



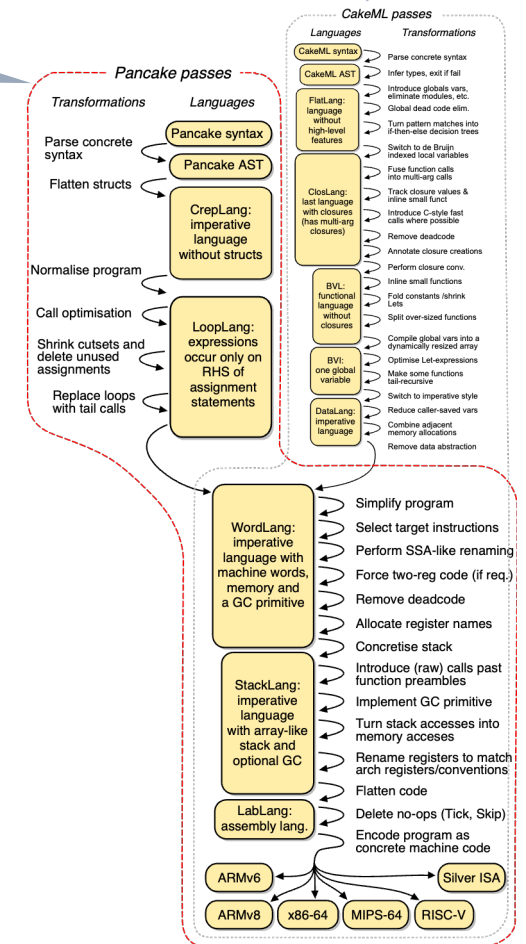
Pancake: New Systems Language

Approach:

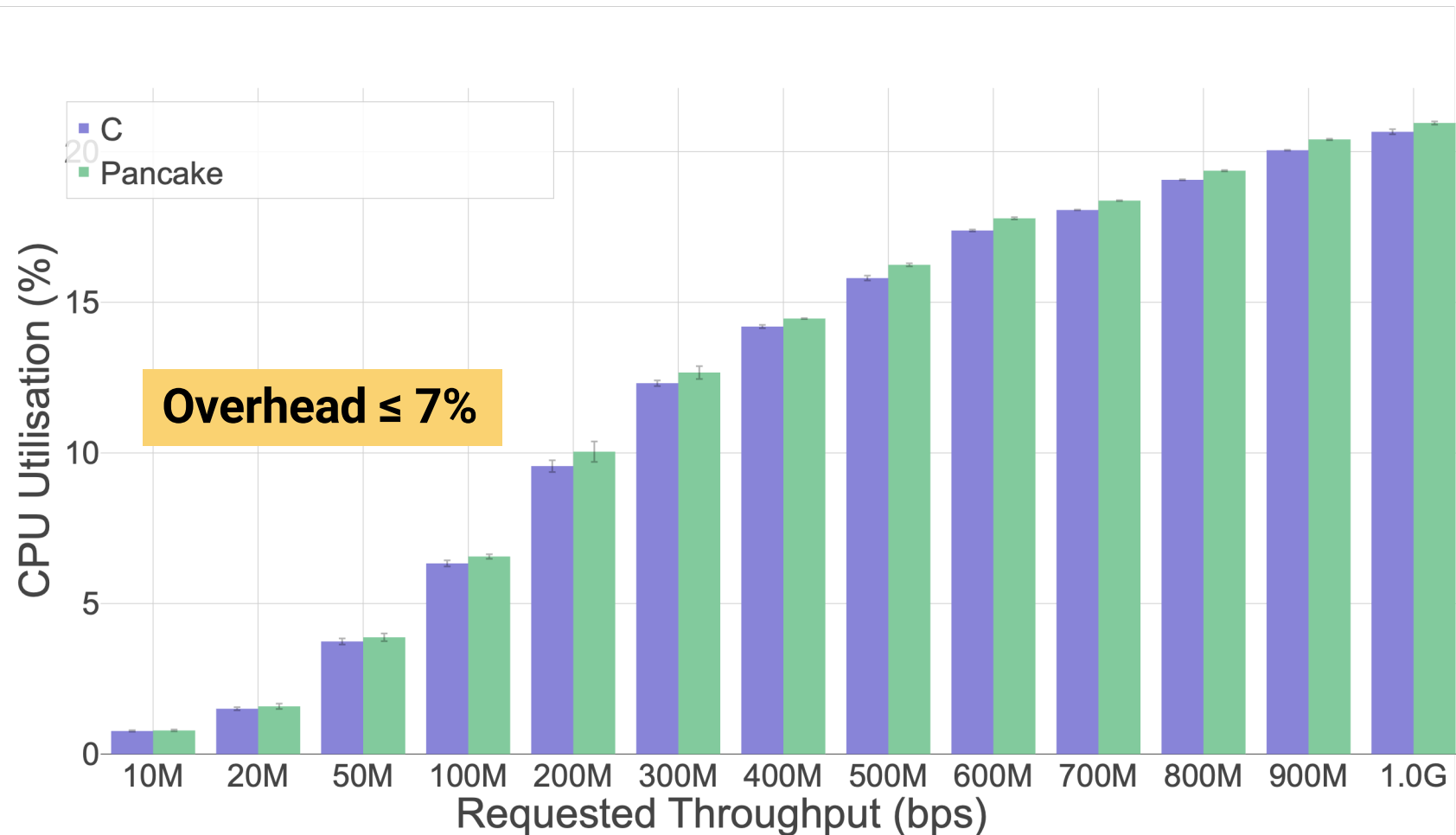
- Re-use lower part of CakeML compiler stack
- Get verified Pancake compiler quickly
- Retain mature framework/ecosystem

Pancake

CakeML



Performance: LionsOS Ethernet Driver



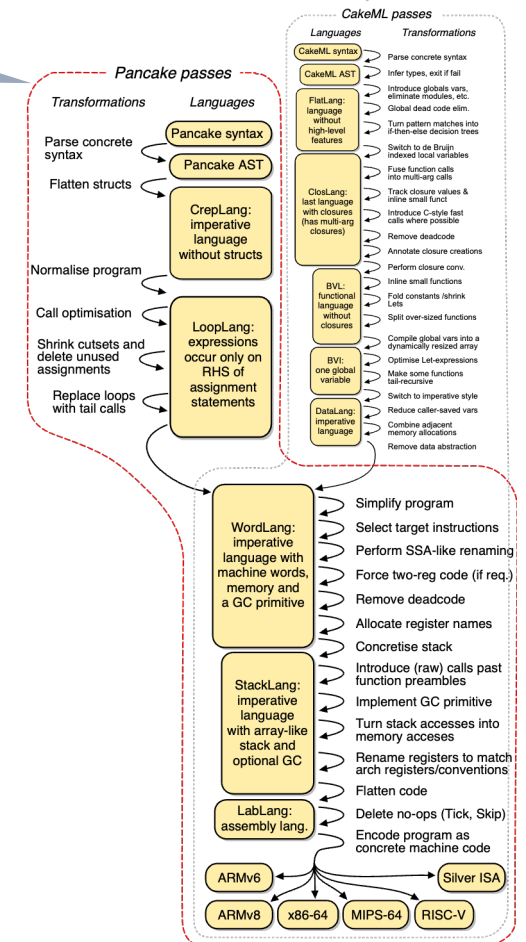
Pancake: New Systems Language

Pancake

CakeML

Status:

- “Usable” rump language, sufficient for LionsOS
 - Implemented drivers & other LionsOS components
 - Need C-escapes for cache management instructions
- **Verified compiler!**
- In progress:
 - Verification of components
 - More performance work
 - Decompilation from Pancake to HOL
 - Semantics for non-terminating programs
 - Efficient verification framework (Hoare logic)



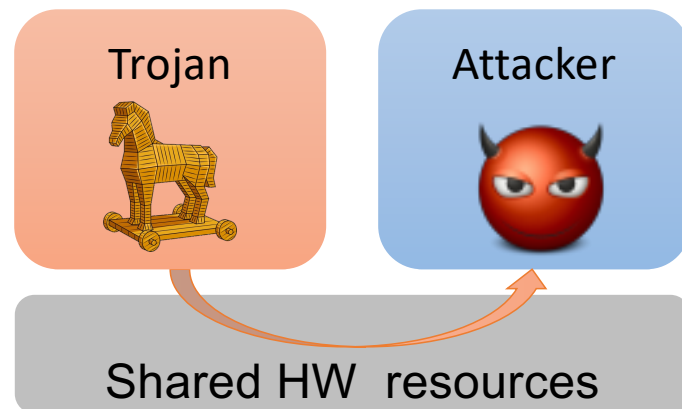
Time Protection

Principled Prevention of Microarchitectural Timing Channels

Refresh: Covert Timing Channels

- Created by contention for shared resource whose effect on timing can be monitored
 - Cache, network bandwidth, CPU load...

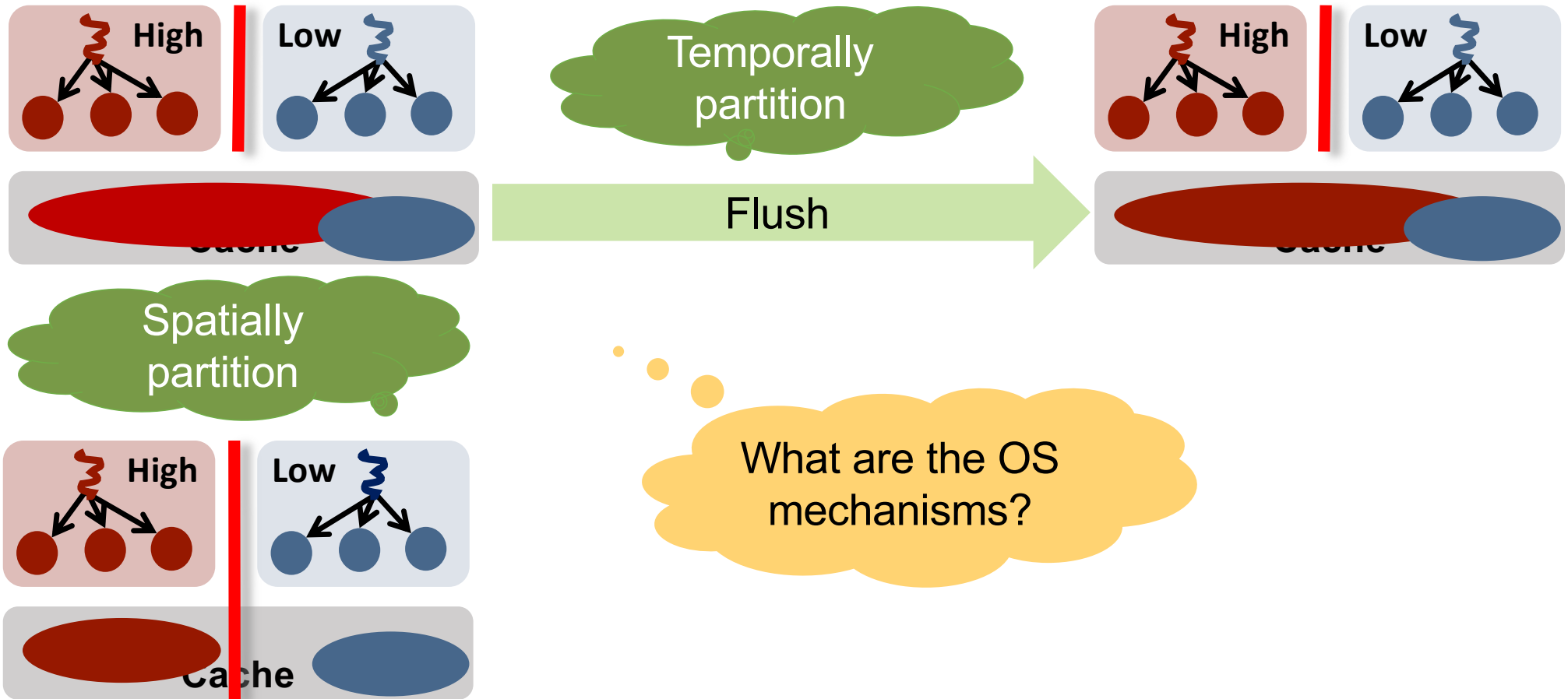
Usually based on some hidden state (e.g. caches)
⇒ timing/storage channel distinction is not deep!



OS protection mechanisms:

- *Memory protection* for spatial isolation
- *Time protection* for temporal isolation

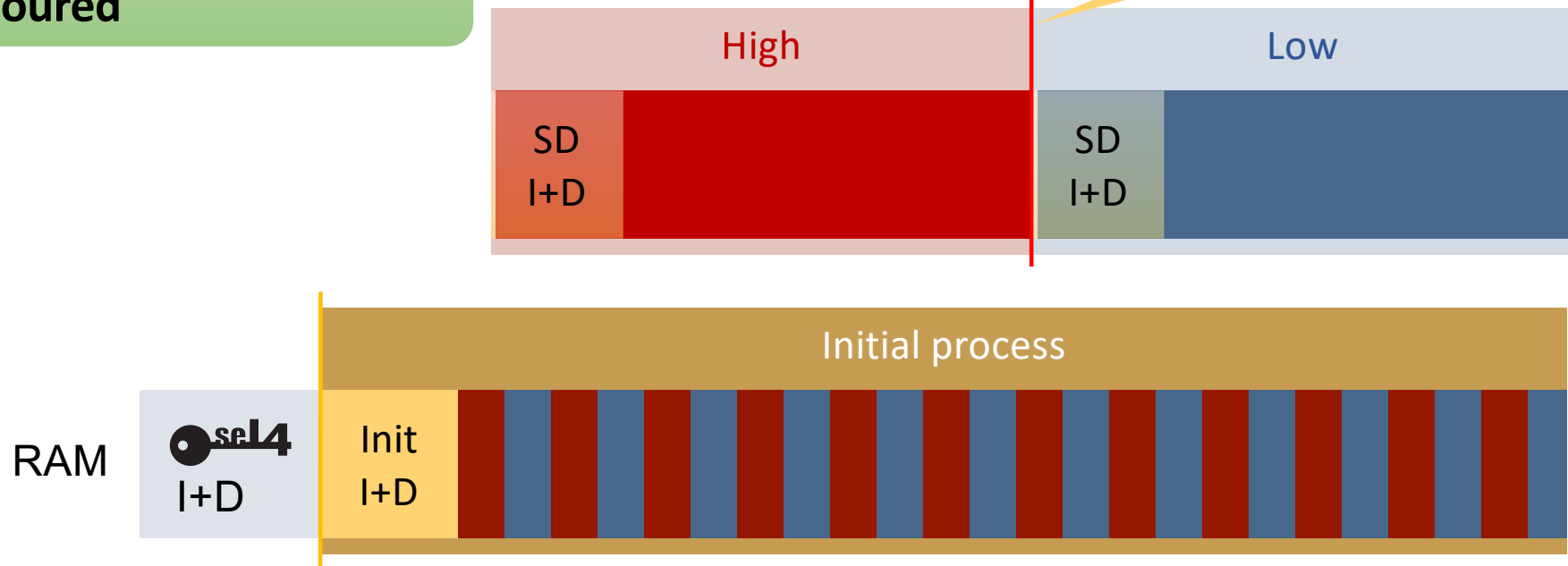
Time Protection: No Sharing of HW State



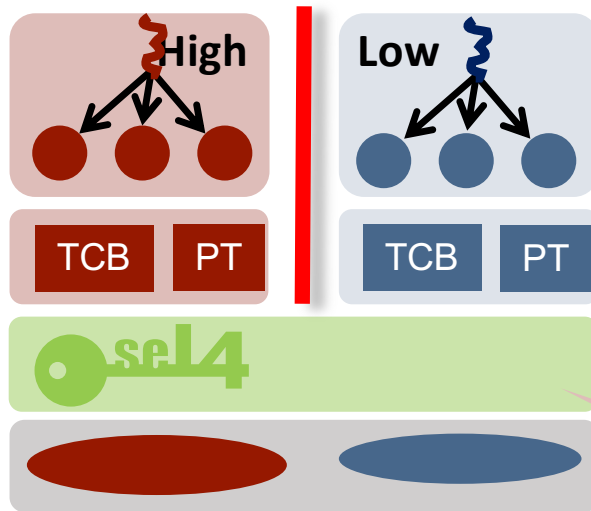
seL4 Spatial Partitioning: Cache Colouring

System permanently coloured

Partitions restricted to coloured memory

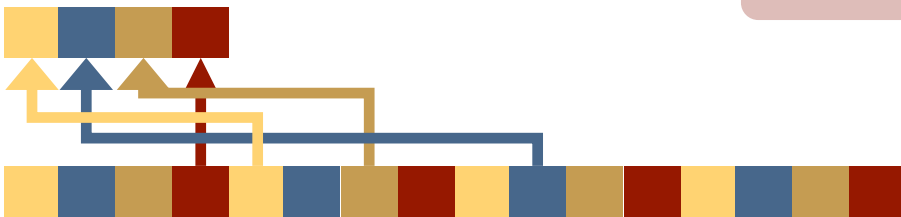


seL4 Spatial Partitioning: Cache Colouring

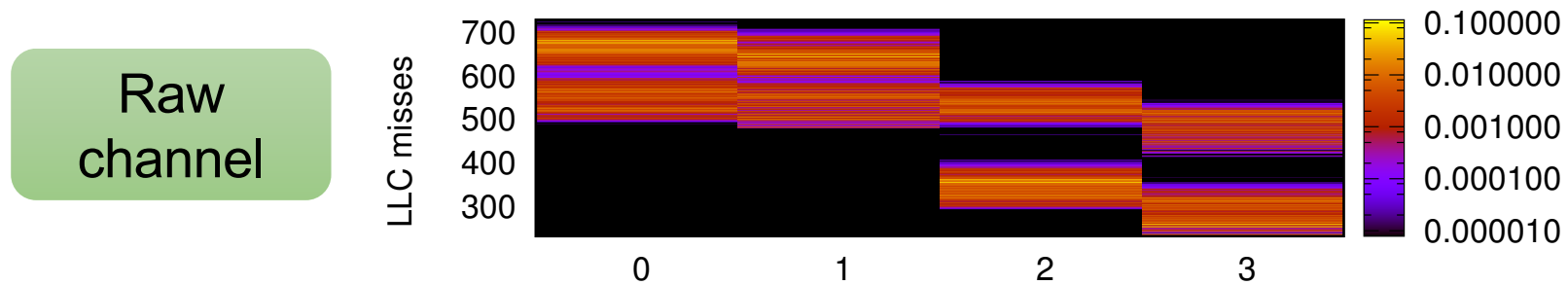


- Partitions get frame pools of disjoint colours
- seL4: userland supplies kernel memory
⇒ colouring userland colours kernel memory

Shared kernel image



seL4 Channel Through Kernel Code



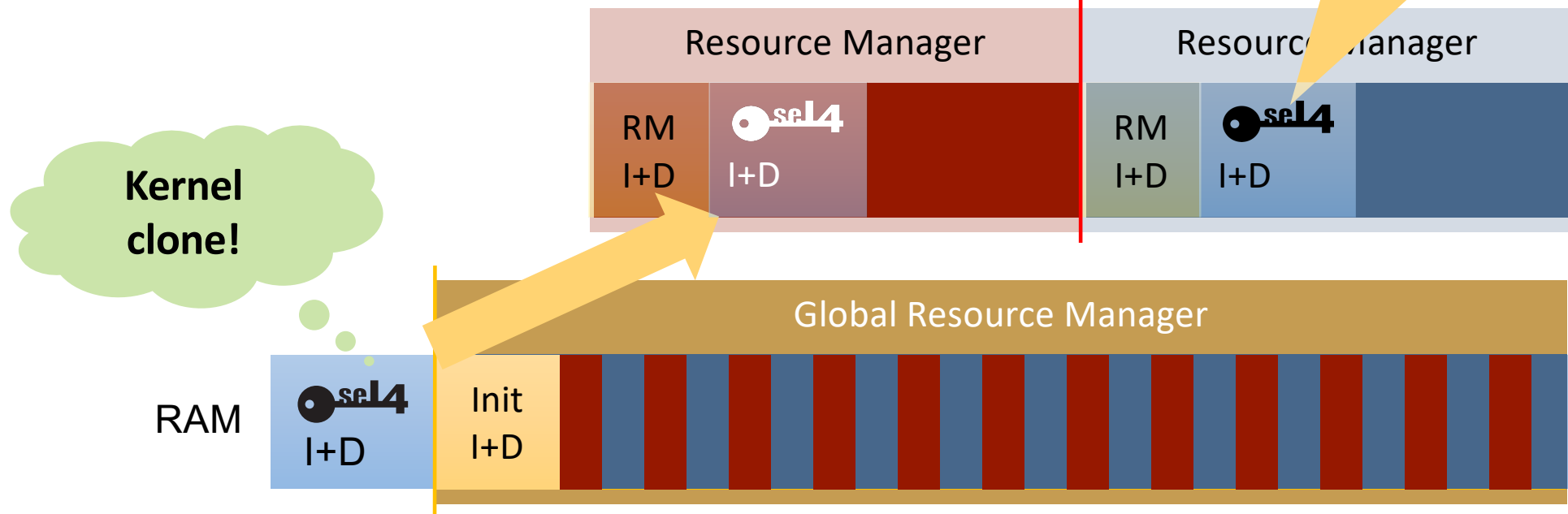
Channel matrix: Conditional probability of observing output signal (time) given input signal (system-call number)

seL4 Colouring the Kernel

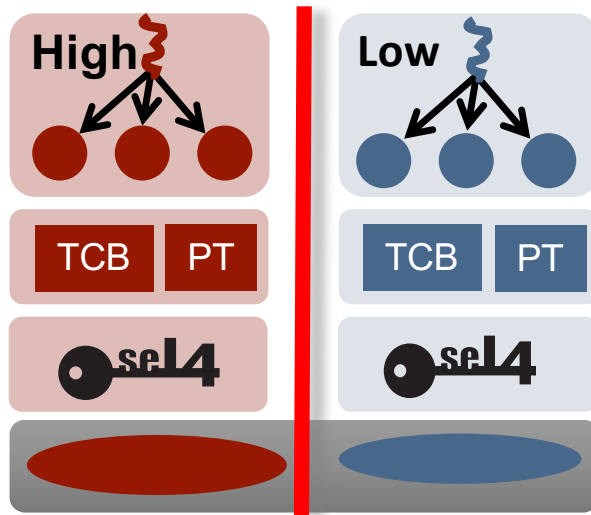
Remaining shared kernel data:

- Scheduler queue array & bitmap
- Few pointers to current thread state

Each partition has own kernel image

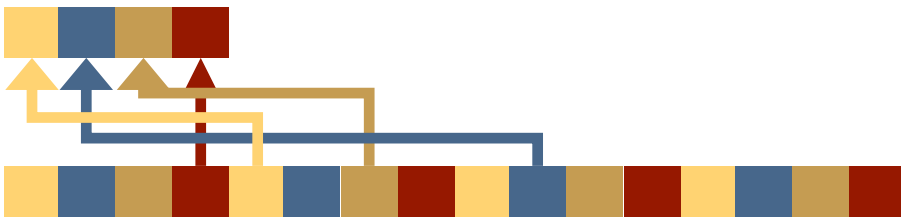


seL4 Spatial Partitioning: Cache Colouring



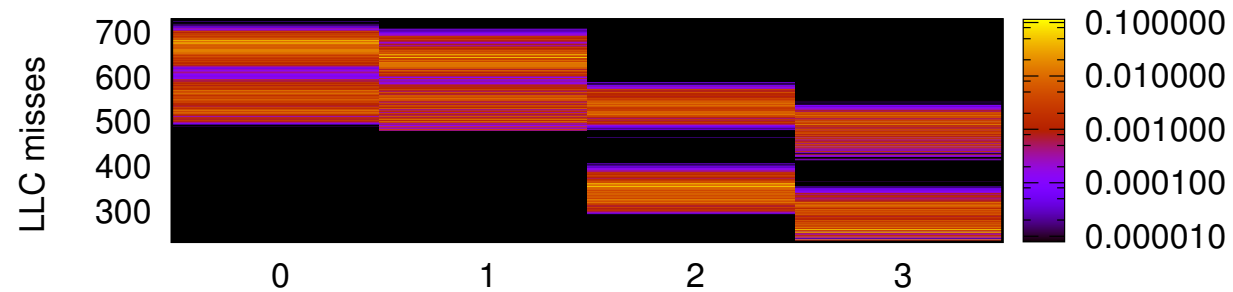
- Partitions get frame pools of disjoint colours
- seL4: userland supplies kernel memory
⇒ colouring userland colours kernel memory
- Per-partition kernel image to colour kernel

Must ensure deterministic access to remaining shared kernel state!

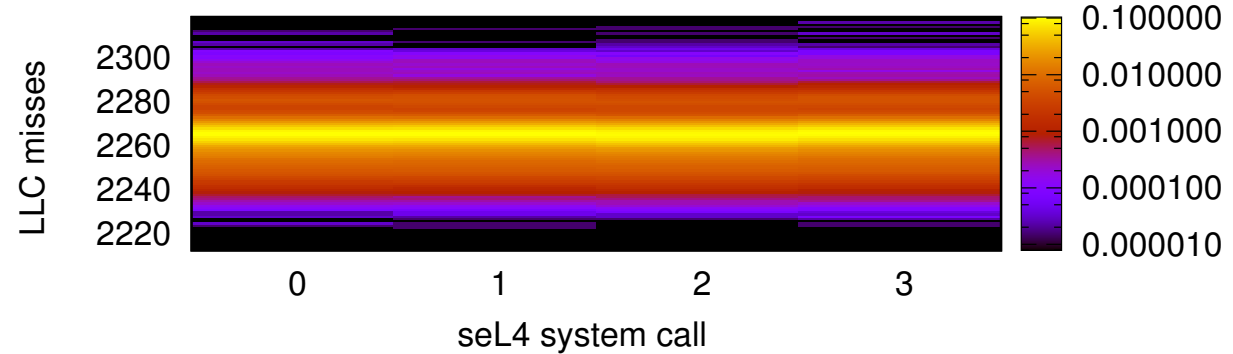


seL4 Channel Through Kernel Code

Raw
channel



Channel with
cloned kernel



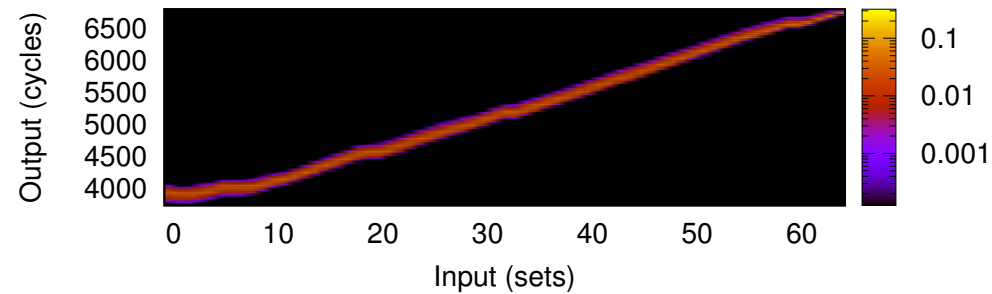
seL4 Temporal Partitioning: Flush on Switch

Must remove any history dependence!

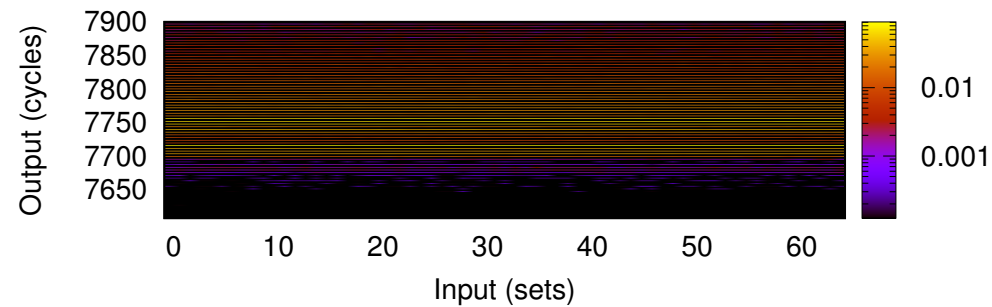
2. Switch user context
3. Flush on-core state
6. Reprogram timer
7. return

seL4 D-Cache Channel

Raw
channel

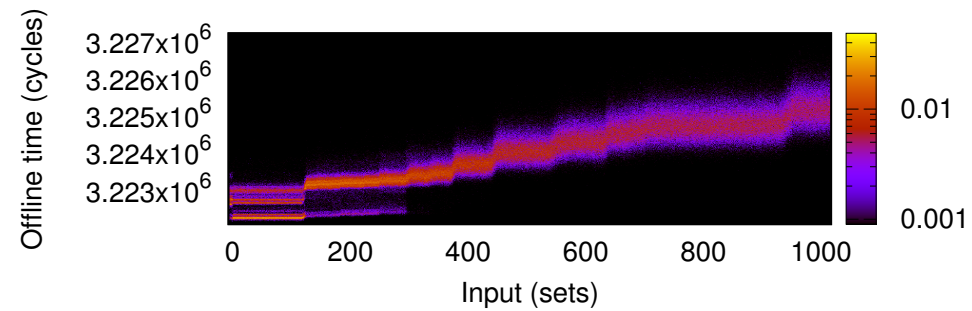


Channel with
flushing



seL4 Flush-Time Channel

Raw
channel



seL4 Temporal Partitioning: Flush on Switch

Must remove any history dependence!

1. $T_0 = \text{current_time}()$
2. Switch user context
3. Flush on-core state
4. Touch all shared data needed for return
5. $\text{while } (T_0 + \text{WCET} < \text{current_time}()) ;$
6. Reprogram timer
7. return

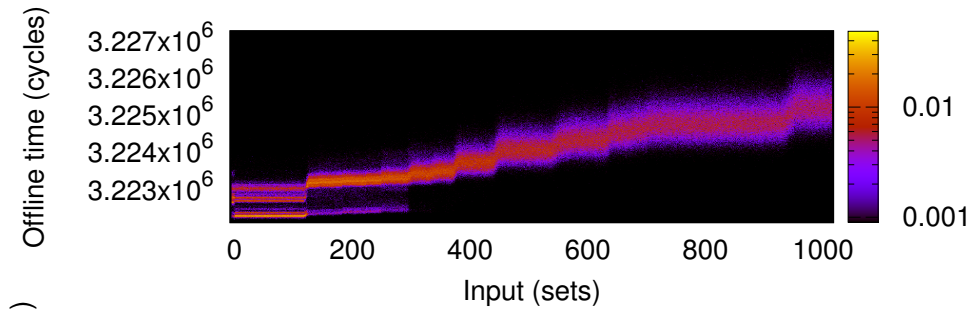
Latency depends on prior execution!

Time padding to remove dependency

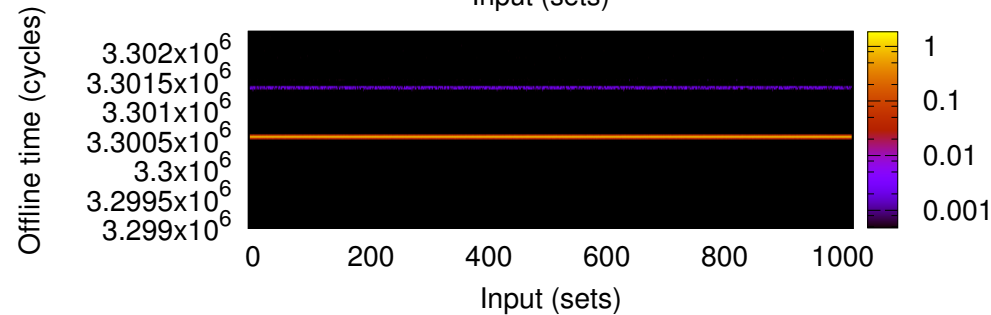
Ensure deterministic execution

seL4 Flush-Time Channel

Raw
channel

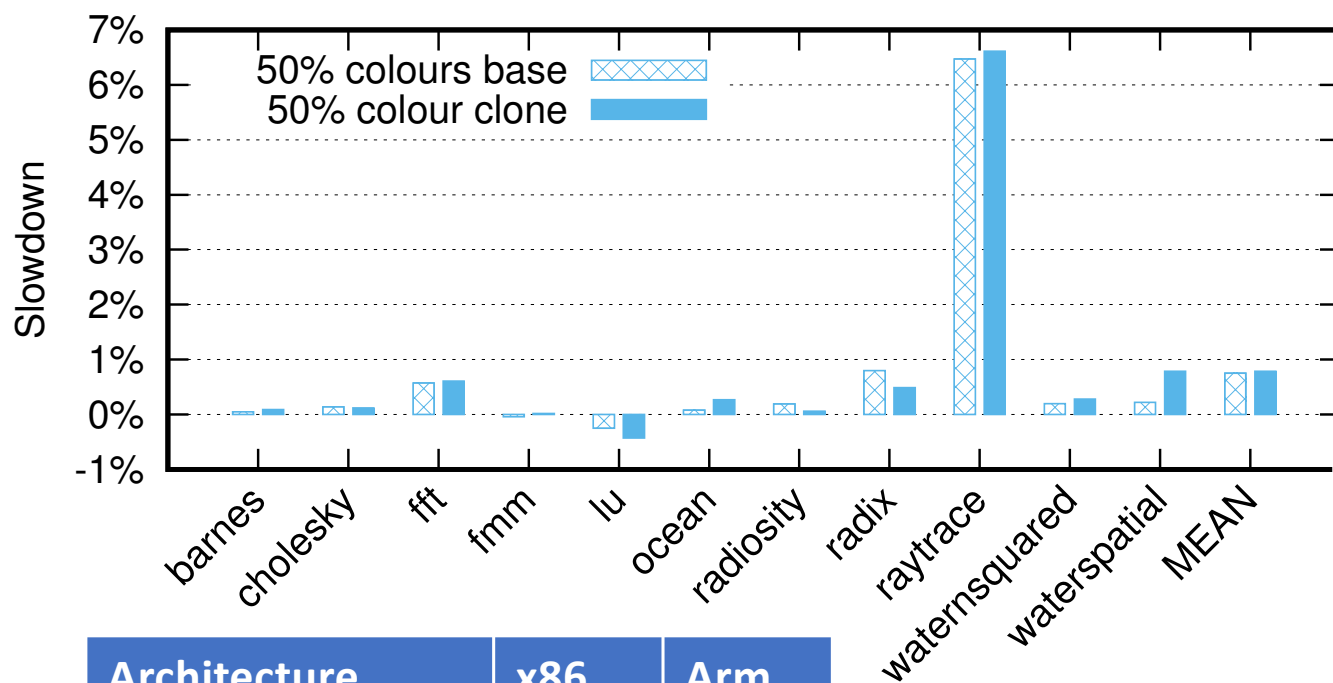


Channel with
deterministic
flushing



seL4 Performance Impact of Colouring

Splash-2 benchmarks on Arm A9



- Overhead mostly low
- Not evaluated is cost of not using super pages [Ge et al., EuroSys'19]

Architecture	x86	Arm
Mean slowdown	3.4%	1.1%

Arch	seL4 clone	Linux fork+exec
x86	79 μ s	257 μ s
Arm	608 μ s	4,300 μ s

seL4 Temporal Partitioning: Flush State

Must remove any history dependence!

1. $T_0 = \text{current_time}()$
2. Switch user context
3. Flush on-core state
4. Touch all shared data needed for return
5. $\text{while } (T_0 + \text{WCET} < \text{current_time}()) ;$
6. Reprogram timer
7. return

Problem: Processors do *not* provide mechanisms for resetting all microarchitectural state!

A New HW/SW Contract

For all shared microarchitectural resources:

aISA: augmented ISA

1. Resource must be spatially partitionable or flushable
2. Concurrently shared resources must be spatially partitioned
3. Resource accessed solely by virtual address must be flushed and not concurrently accessed
4. Mechanisms must be sufficiently specified for OS to partition or reset
5. Mechanisms must be constant time, or of specified, bounded latency
6. Desirable: OS should know if resettable state is derived from data, instructions, data addresses or instruction addresses

Cannot share HW threads across security domains!

[Ge et al., APSys'18]



RISC-V To The Rescue: fence.t

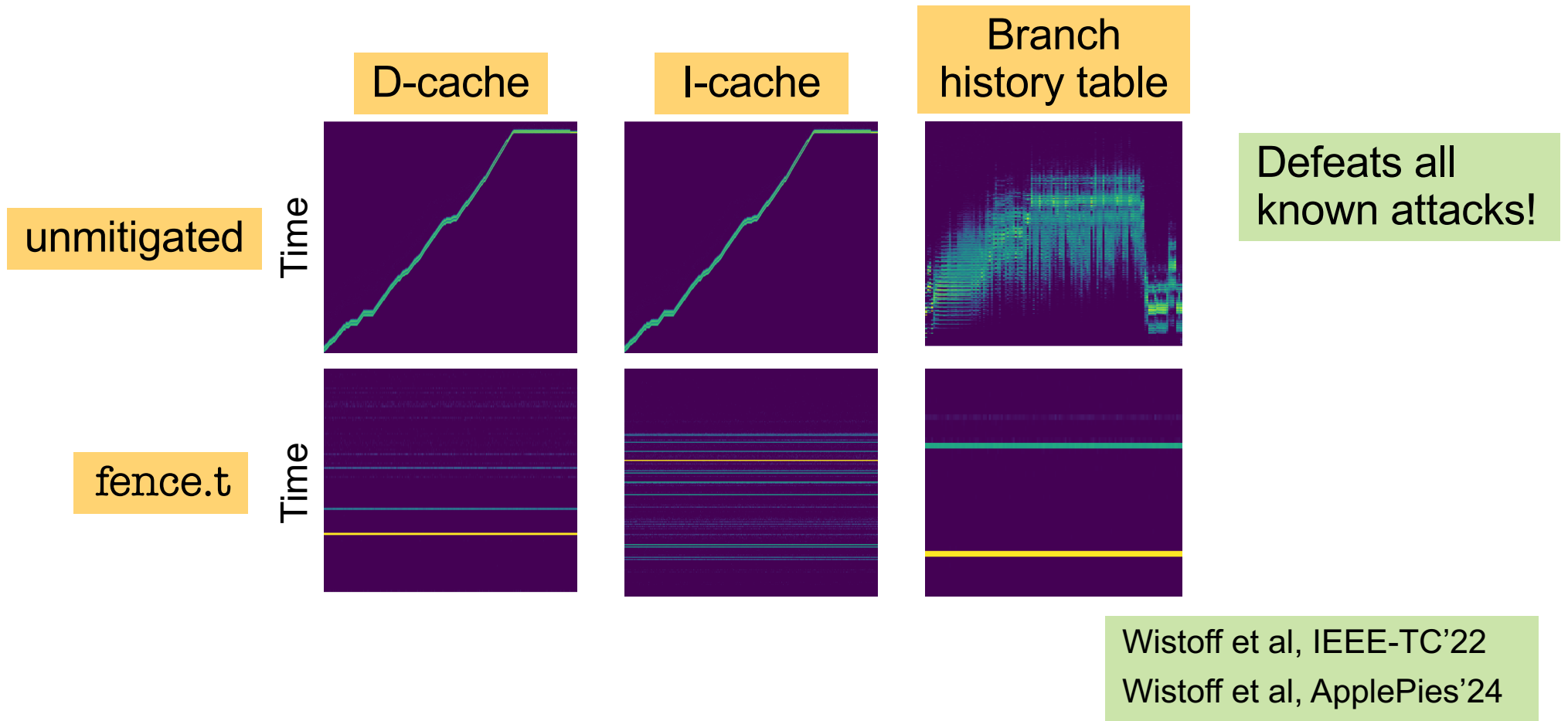
fence.t instruction:

- Flush d-cache
- Reset all flip-flops that are **not** part of architected state



- Prototyped on in-order (CVA6) and OoO (C910) RISC-V processors
- Latency bounded by d-cache flush
- HW cost in the noise

fence.t Instruction on C910





Can Time Protection Be Verified?

1. Correct treatment of spatially partitioned state:

- Need hardware model that identifies all such state (augmented ISA)
- To prove:

No two domains can access the same physical state

Functional property!

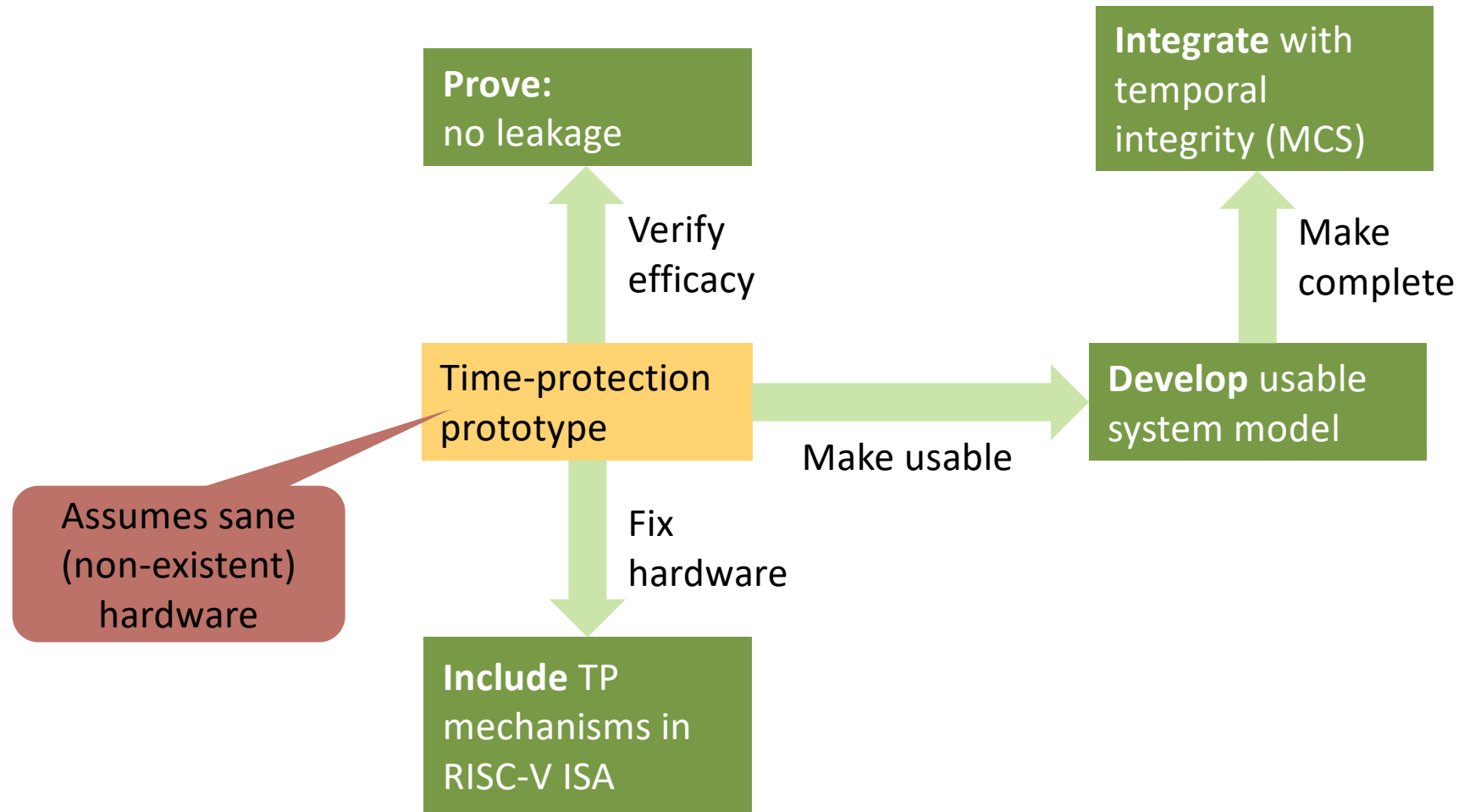
Transforms timing channels
into storage channels!

2. Correct flushing of time-shared state

- Not trivial: eg proving all cleanup code/data are forced into cache after flush
 - Needs an actual cache model
- Even trickier: need to prove padding is correct
 - ... without explicitly reasoning about time!

Functional property!

seL4 Time Protection: On-Going Work





Real-World Use

Courtesy Boeing, DARPA



Thank you!

To the brave AOS students for their interest and dedication

To the world-class Trustworthy Systems team for making all possible

Please remember to do the myExperience survey

There'll also be a more detailed one we'll invite you to fill in