

Provable Multicore Schedulers with **Ipanema**: Application to Work-Conservation

Baptiste Lepers

Redha Gouicem

Damien Carver

Jean-Pierre Lozi

Nicolas Palix

Virginia Aponte

Willy Zwaenepoel

Julien Sopena

Julia Lawall

Gilles Muller



THE UNIVERSITY OF
SYDNEY



SORBONNE
UNIVERSITÉ



UNIVERSITÉ
Grenoble
Alpes

Inria

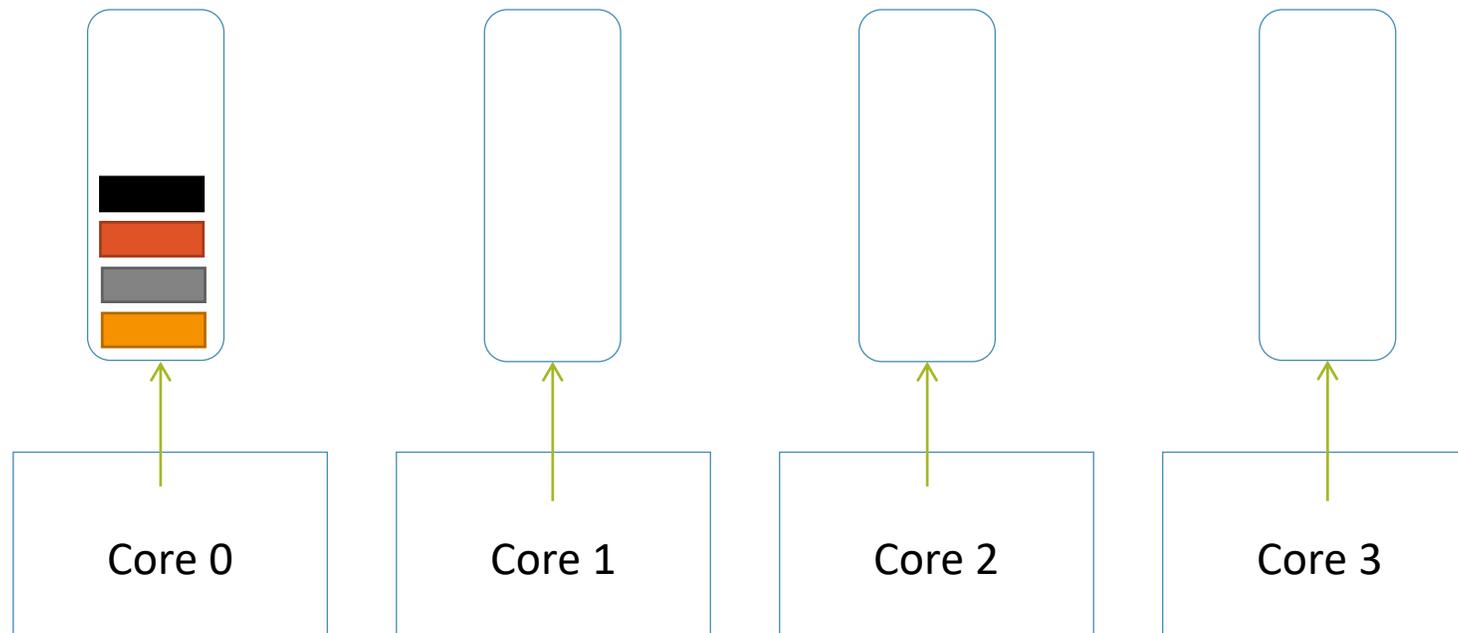
cham

ORACLE®



Work conservation

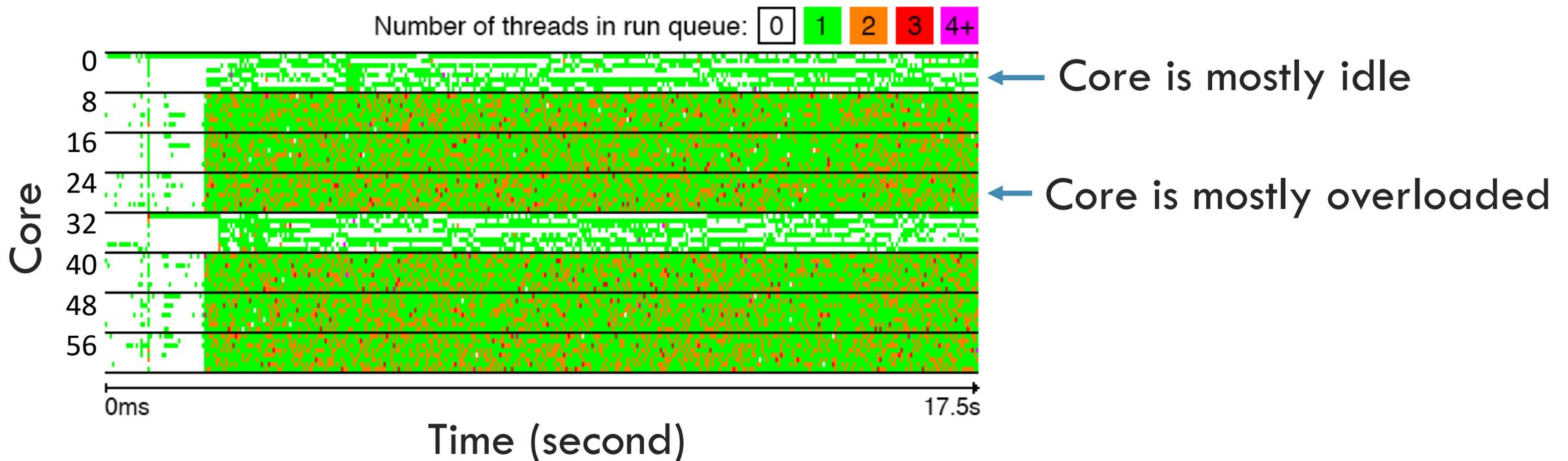
“No core should be left idle when a core is overloaded”



Non work-conserving situation: core 0 is overloaded, other cores are idle

Problem

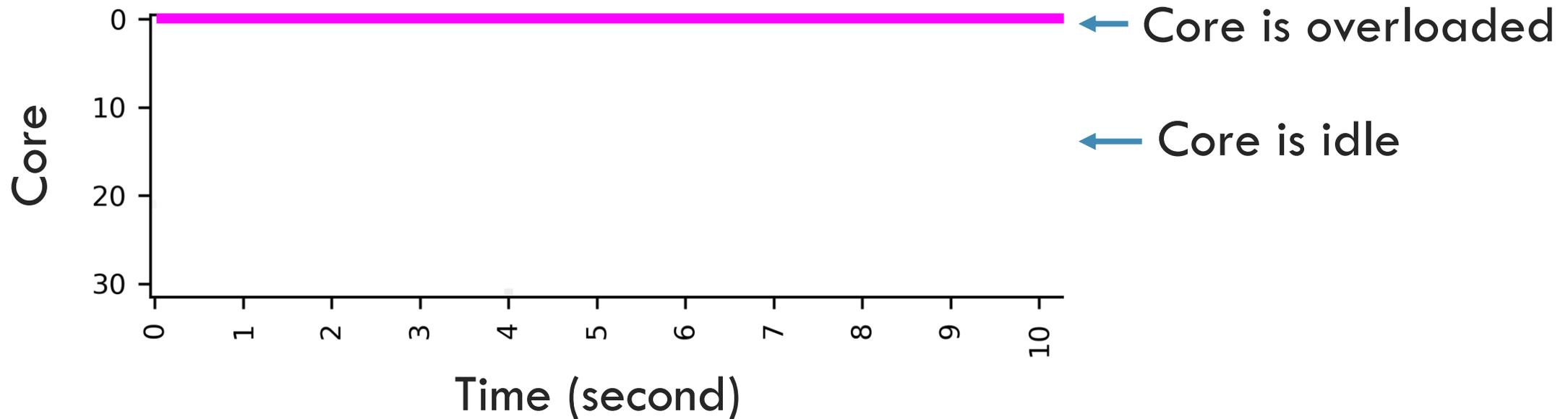
Linux (CFS) suffers from work conservation issues



[Lozi et al. 2016]

Problem

FreeBSD (ULE) suffers from work conservation issues



[Bouron et al. 2018]

Problem

Work conservation bugs are hard to detect

No crash, no deadlock. No obvious symptom.

1.37x slowdown on HPC applications

23% slowdown on a database.

[Lozi et al. 2016]

This talk

Formally prove work-conservation

$a = \frac{\Delta v}{\Delta t} = \frac{500 \text{ m/s}}{0,02 \text{ s}} = 500 \text{ m/s}^2$

$F = 500 \text{ m/s}^2 \cdot 0,5 \text{ kg} = 250 \text{ N}$

$\frac{2 \cdot 10^{17}}{(7 \cdot 10^5)^2} \cdot 6,67 \cdot 10^{-11}$

$286,42 \text{ m/s}^2 < 1,9 \cdot 10^6$

$\Delta v = mg$

$\vec{v}_2 = \frac{1}{\Delta t} \cdot \dots$

$T = 2\pi\sqrt{\frac{p}{g}}$

Einstein

$e = -1,6 \cdot 10^{-19}$

$\vec{p} = 1,6 \cdot 10^{-19}$

$S = \frac{v^2 - v_0^2}{2a}$

$E_1 + E_2 = E_p + E_{p2}$

$\vec{v}_2 = (m_1 + m_2) \dots$

$A = Fh = mgd$

$m = \rho SP = 390$

$S = \sqrt{E} \quad p = m \dots$

$\vec{F}_1, \vec{F}_2, \vec{F}_3$

$E = h\nu$

$w = 2\pi f$

$\vec{D}, \vec{d}, \vec{S}, \vec{F}$

$\vec{v} = mg$

\vec{m}, \vec{g}

$T = 2\pi\sqrt{\frac{m}{k}}$

$T = \frac{1}{w}$

$S \cdot \vec{v}_0 t \pm \frac{at^2}{2}$

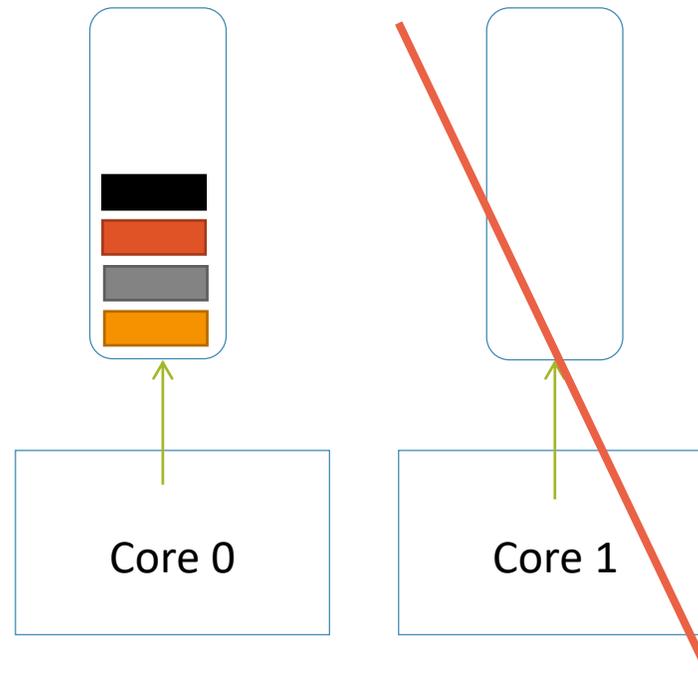
$9/3e + 4 \rightarrow 239 \text{ Pu} \rightarrow 239 \text{ Pu} + 0,1e$

$239 \text{ Pu} \rightarrow 239 \text{ Pu} + 0,1e$

Work Conservation Formally

$$(\exists c . O(c)) \Rightarrow (\forall c' . \neg I(c'))$$

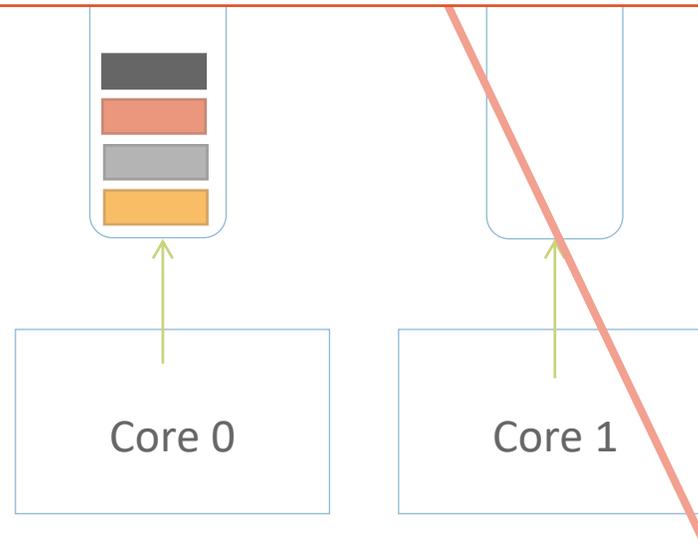
If a core is overloaded, no core is idle



Work Conservation Formally

$$(\exists c . O(c)) \Rightarrow (\forall c' . \neg I(c'))$$

Does not work for realistic schedulers!



Challenge #1

Concurrent events & optimistic concurrency

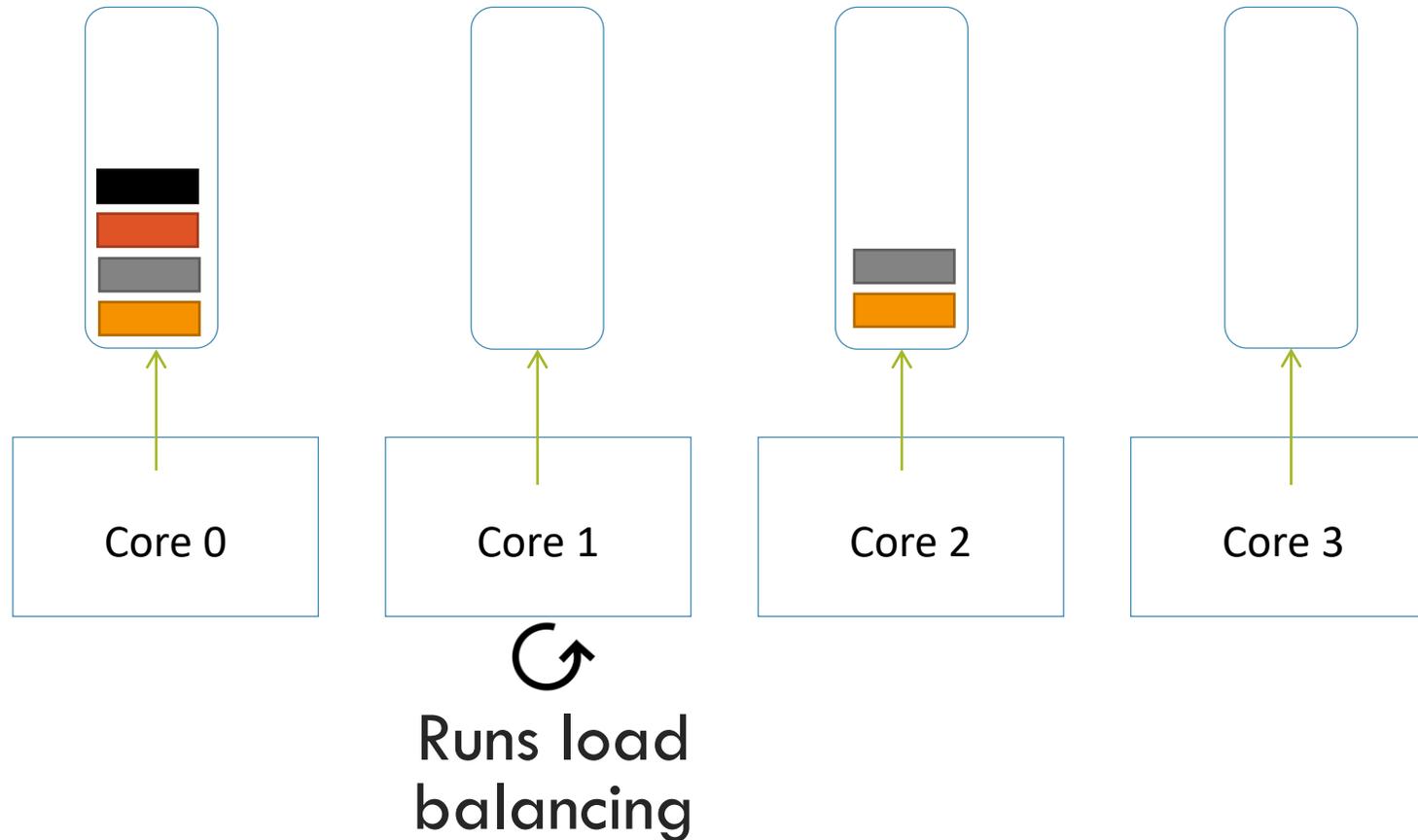
Challenge #1

Concurrent events & optimistic concurrency



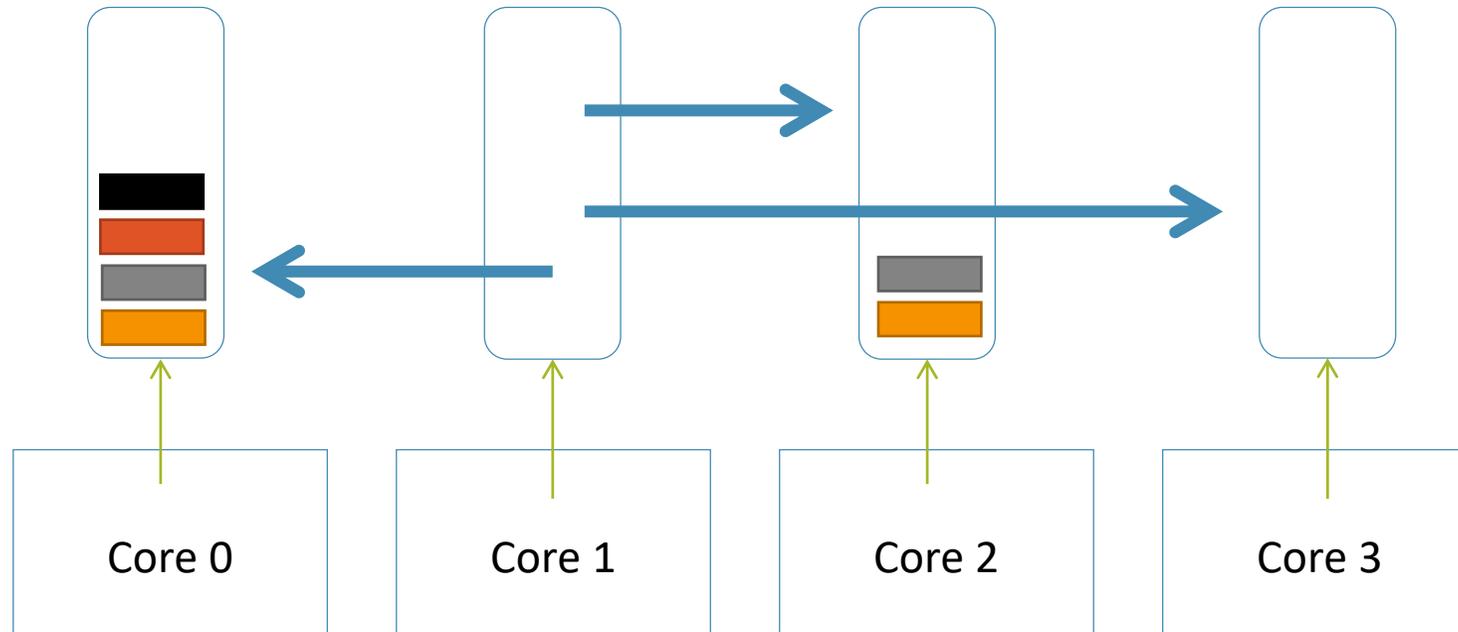
Challenge #1

Concurrent events & optimistic concurrency



Challenge #1

Concurrent events & optimistic concurrency

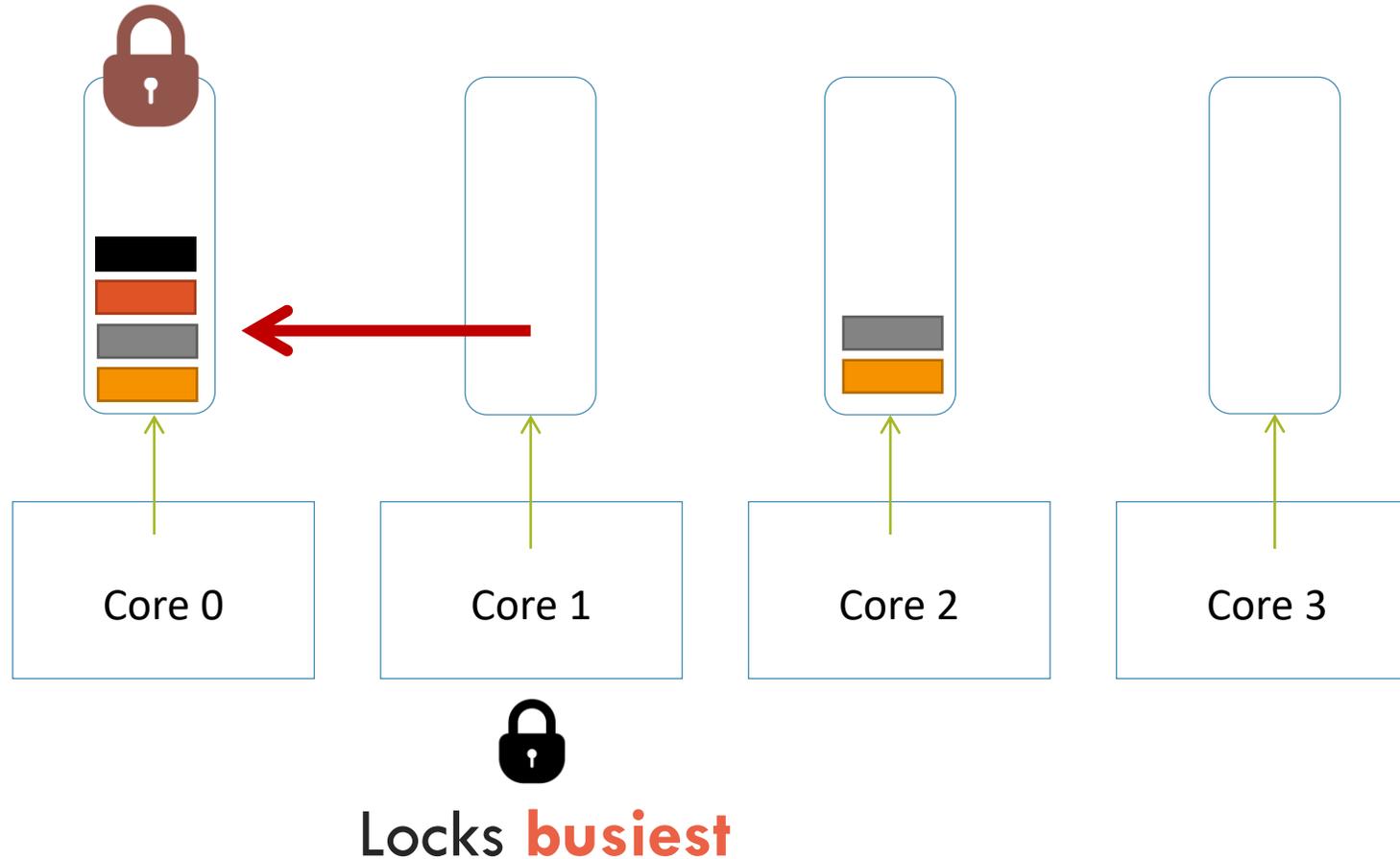


Observes load
(no lock)

Challenge #1

Concurrent events & optimistic concurrency

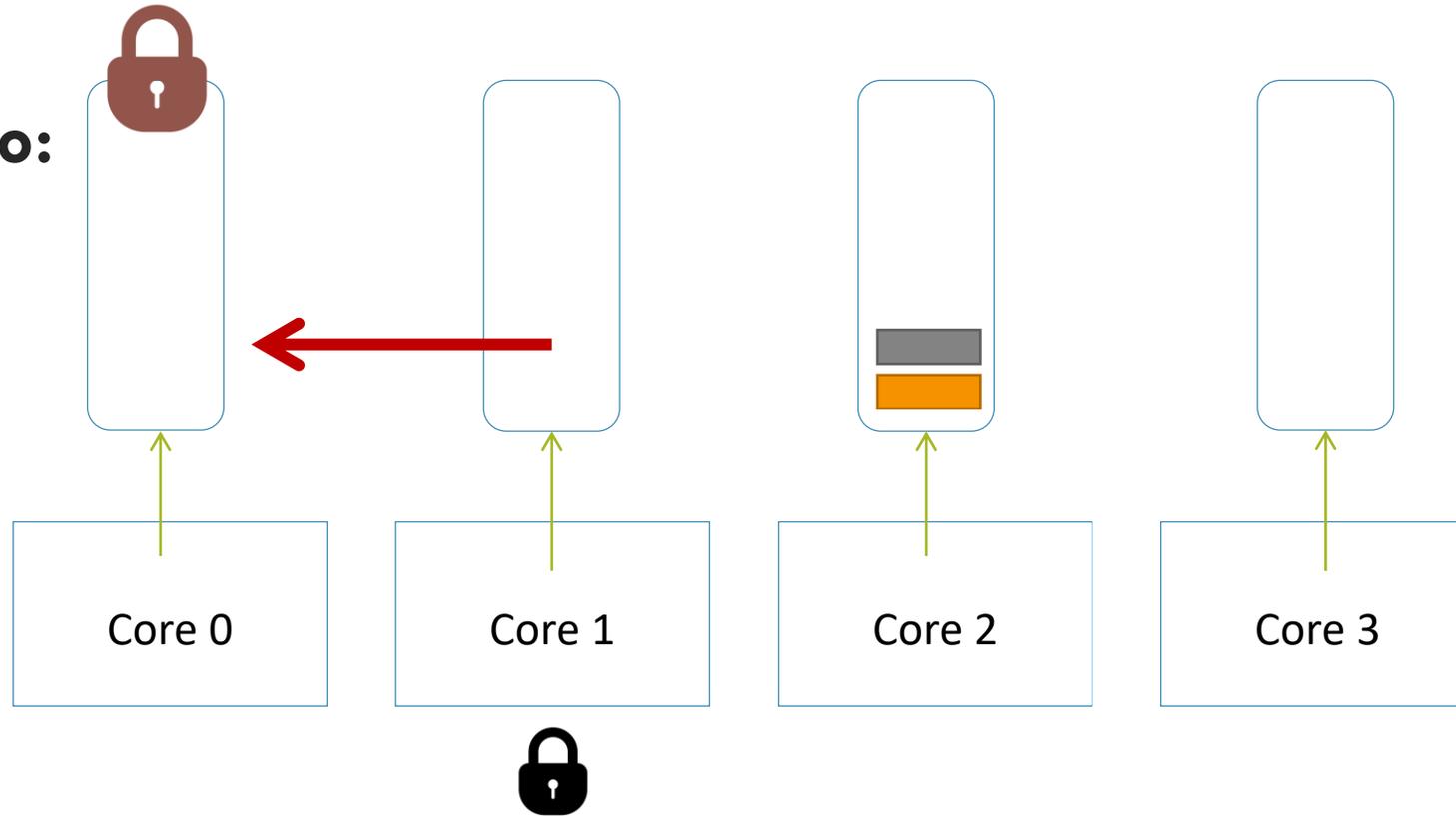
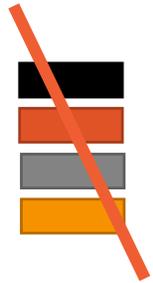
Ideal scenario: no change since observations



Challenge #1

Concurrent events & optimistic concurrency

Possible scenario:

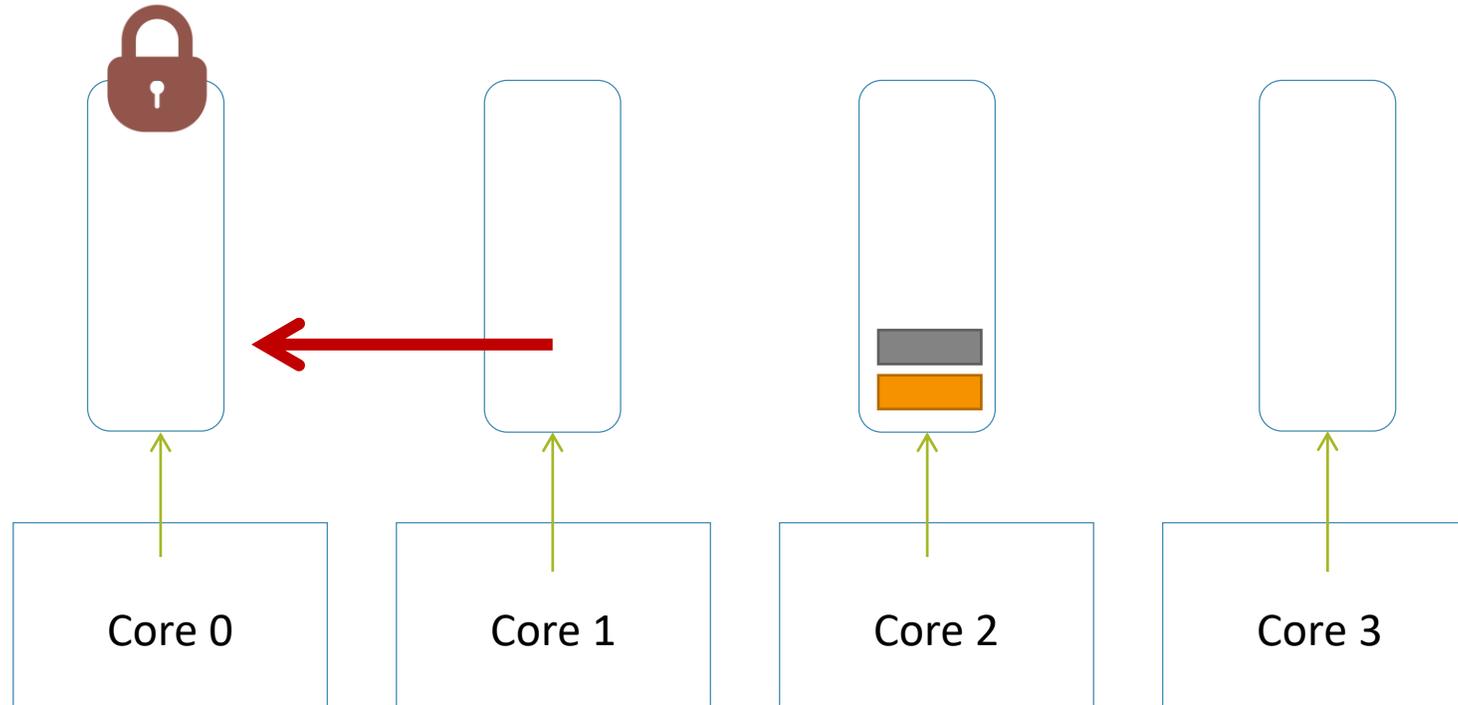


Locks **“busiest”**

Busiest might have no thread left! (Concurrent blocks/terminations.)

Challenge #1

Concurrent events & optimistic concurrency



**(Fail to)
Steal from busiest**

Challenge #1

Concurrent events & optimistic concurrency



**Definition of Work Conservation must take
concurrency into account!**

Concurrent Work Conservation Formally

Definition of overloaded with « failure cases »:

$\exists c . (O(c) \wedge \neg \text{fork}(c) \wedge \neg \text{unblock}(c) \dots)$

If a core is overloaded
(but not because a thread was concurrently created)



Concurrent Work Conservation Formally

$$\begin{aligned} \exists c . (O(c) \wedge \neg \text{fork}(c) \wedge \neg \text{unblock}(c) \dots) \\ \Rightarrow \\ \forall c' . \neg (I(c') \wedge \dots) \end{aligned}$$

Challenge #2

Existing scheduler code is hard to prove

- ⌚ Schedulers handle millions of events per second
Historically: low level C code.

Challenge #2

Existing scheduler code is hard to prove

- ⌚ Schedulers handle millions of events per second
Historically: low level C code.

Code should be easy to prove AND efficient!

Challenge #2

Existing scheduler code is hard to prove

- ⌚ Schedulers handle millions of events per second
Historically: low level C code.

Code should be easy to prove AND efficient!



Domain Specific Language (DSL)

DSL advantages

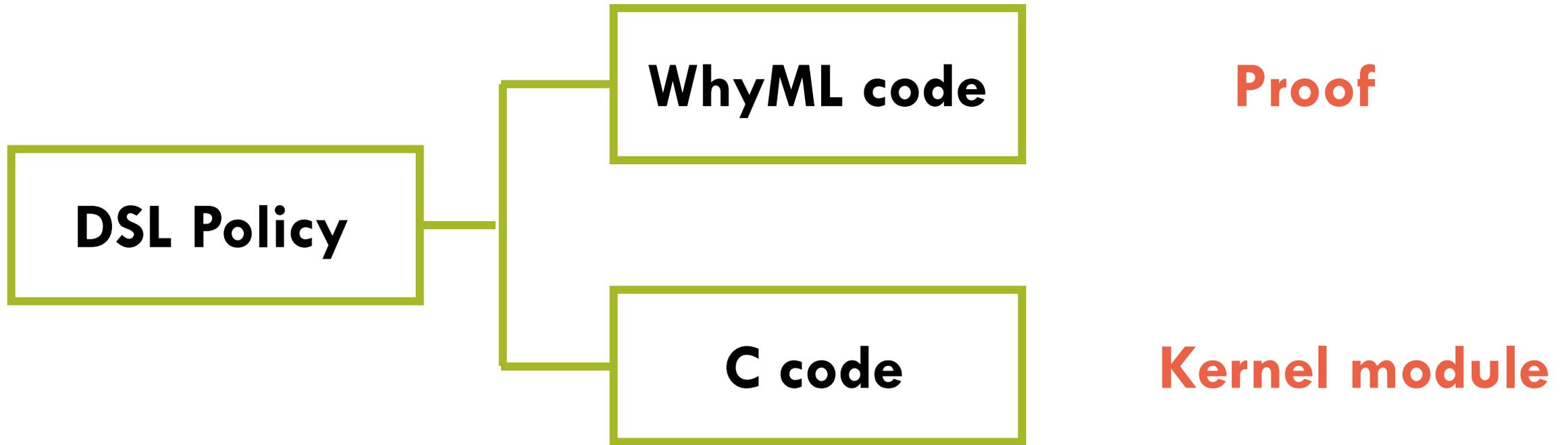
Trade expressiveness for expertise/knowledge:

Robustness: (static) verification of properties

Explicit concurrency: explicit shared variables

Performance: efficient compilation

DSL-based proofs



DSL: close to C
Easy learn and to compile to WhyML and C

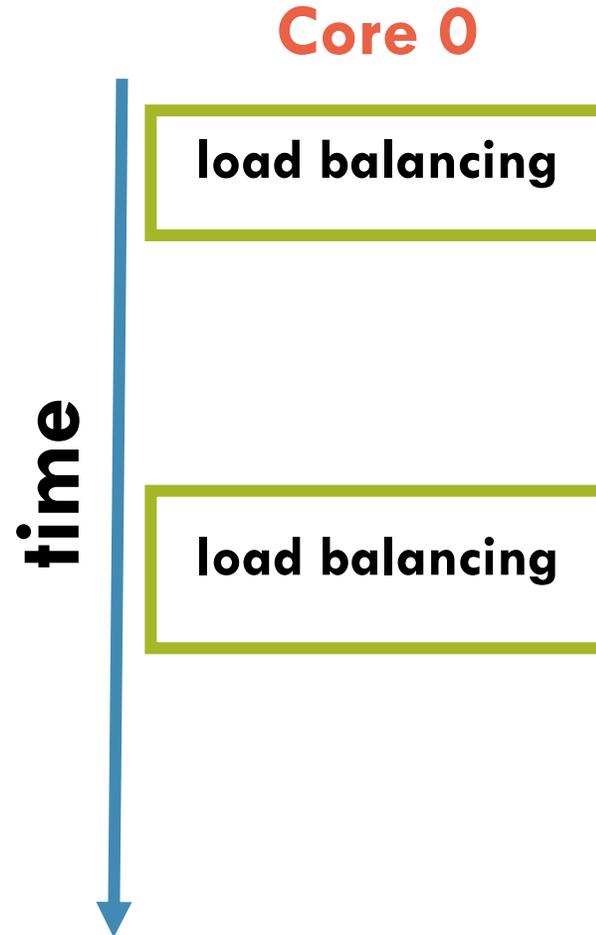
DSL-based proofs

**Proof on all possible
interleavings**

DSL-based proofs

Proof on all possible interleavings

**Split code in blocks
(1 block = 1 read or write to a shared variable)**



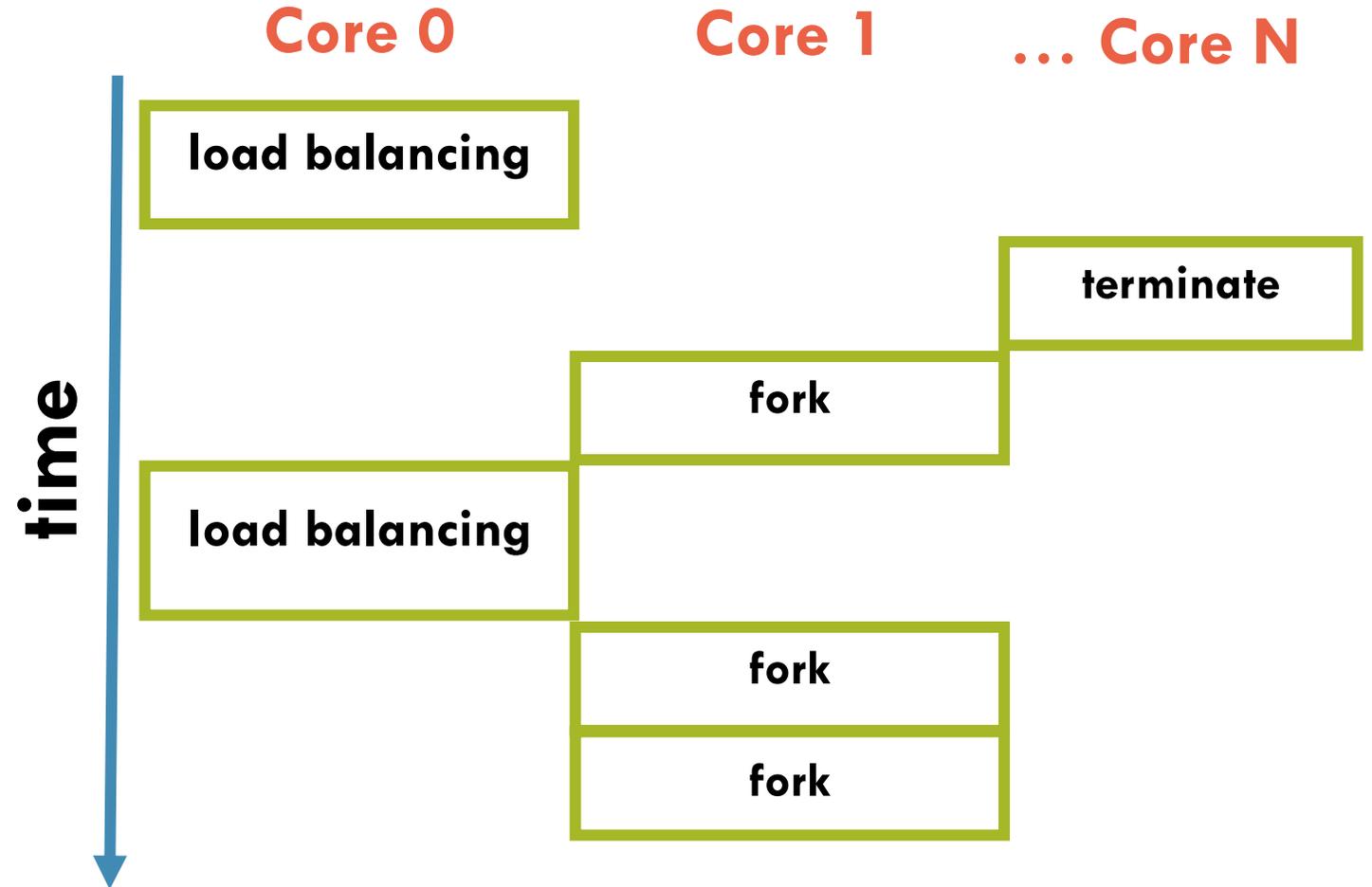
DSL-based proofs

Proof on all possible interleavings

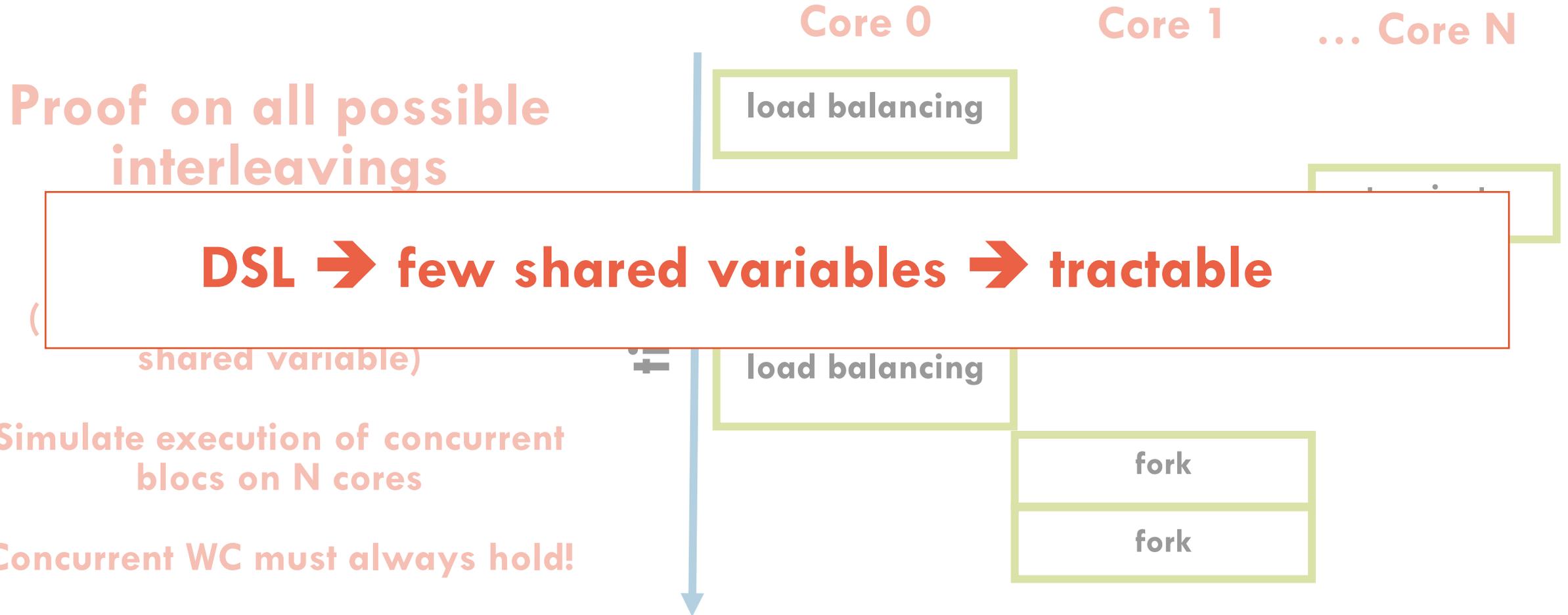
Split code in blocks
(1 block = 1 read or write to a shared variable)

Simulate execution of concurrent blocs on N cores

Concurrent WC must hold at the end of the load balancing



DSL-based proofs



Evaluation

CFS-CWC (365 LOC)

Hierarchical CFS-like scheduler

CFS-CWC-FLAT (222 LOC)

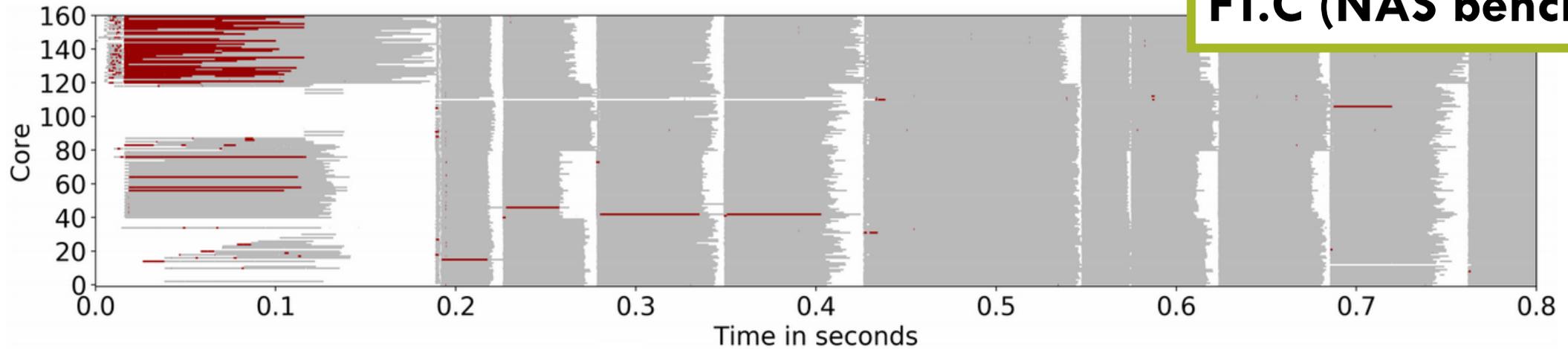
Single level CFS-like scheduler

ULE-CWC (244 LOC)

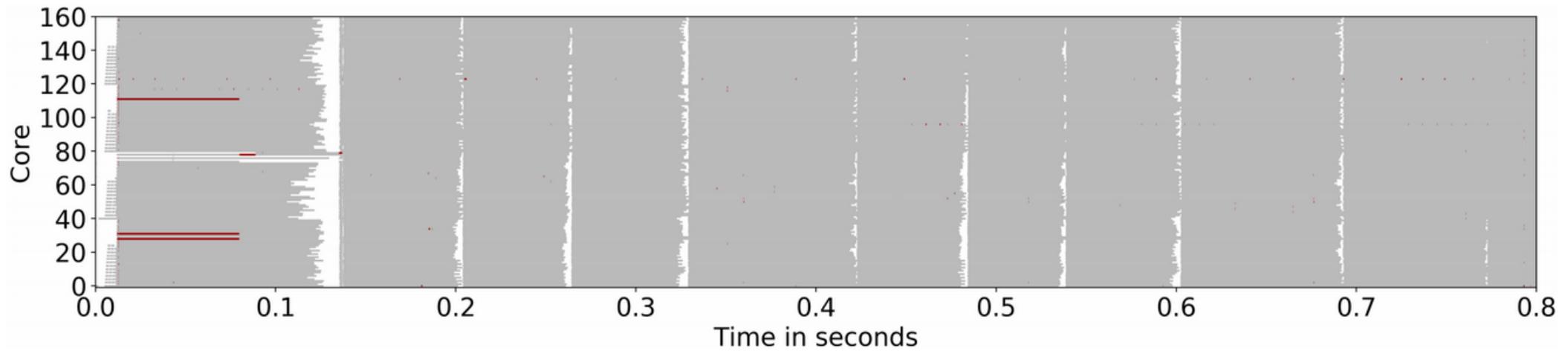
BSD-like scheduler

Less idle time

FT.C (NAS benchmark)



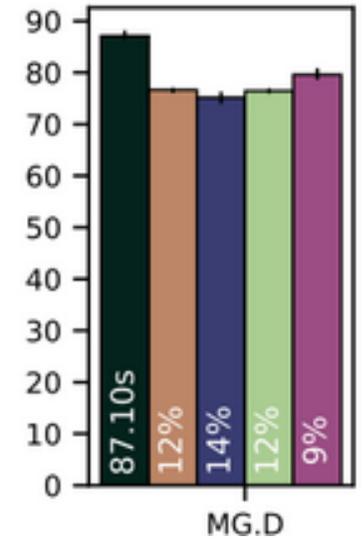
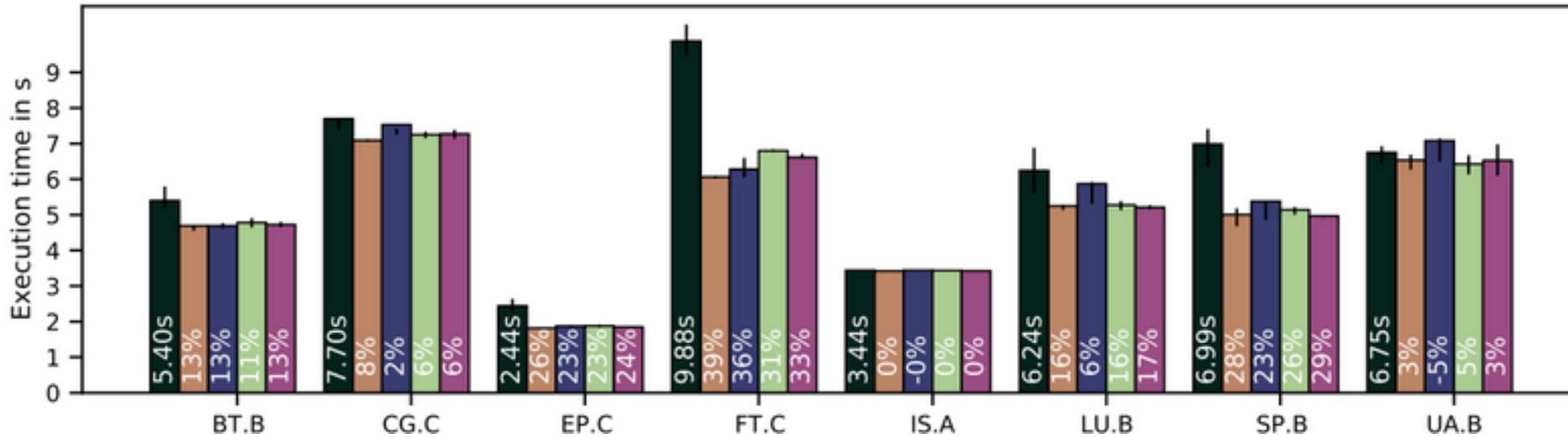
Execution with vanilla CFS.



Execution with CFS-CWC.

Comparable or better performance

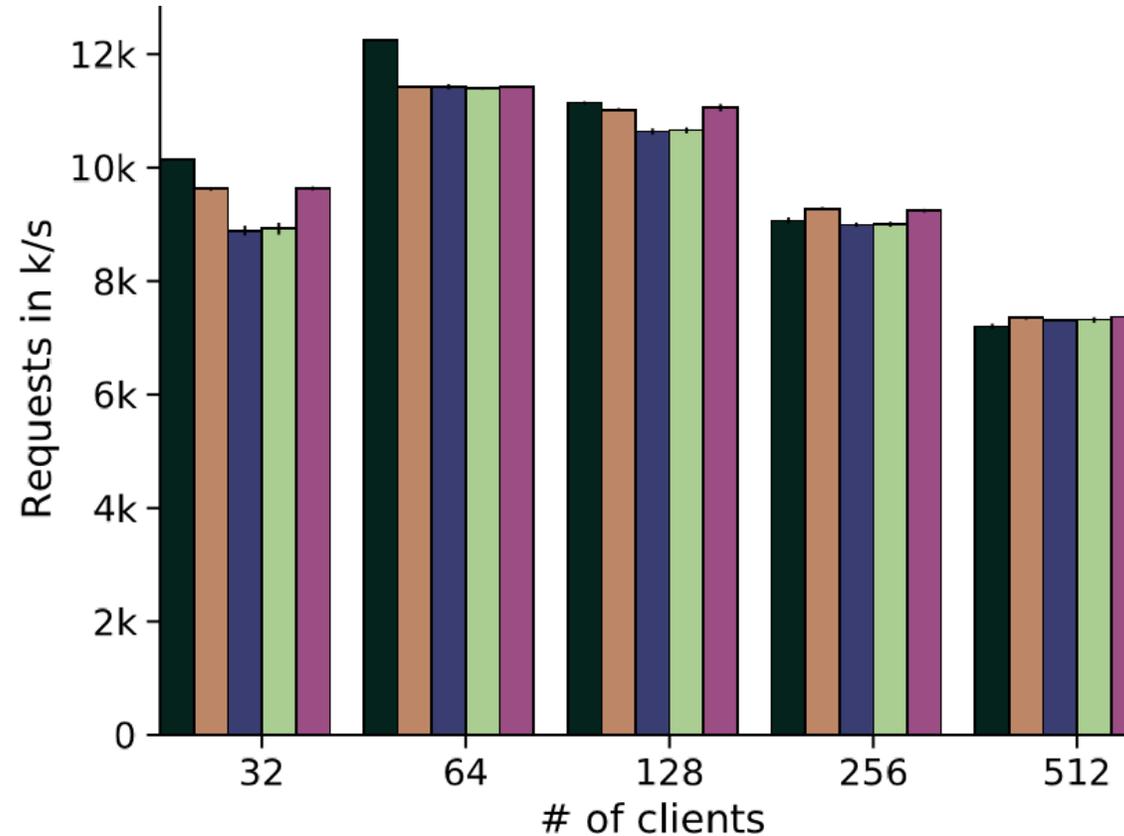
CFS ULE CFS-CWC CFS-CWC-FLAT ULE-CWC



NAS benchmarks (lower is better)

Comparable or better performance

■ CFS ■ ULE ■ CFS-CWC ■ CFS-CWC-FLAT ■ ULE-CWC



Sysbench on MySQL (higher is better)

Conclusion

Work conservation: not straightforward!

... new formalism: *concurrent* work conservation!

Complex concurrency scheme

...proofs made tractable using a DSL.

Performance: similar or better than CFS.