# OFence: Pairing Barriers to Find Concurrency Bugs in the Linux Kernel

Baptiste Lepers
Josselin Giet
Julia Lawall
Willy Zwaenepoel

THE UNIVERSITY OF SYDNEY
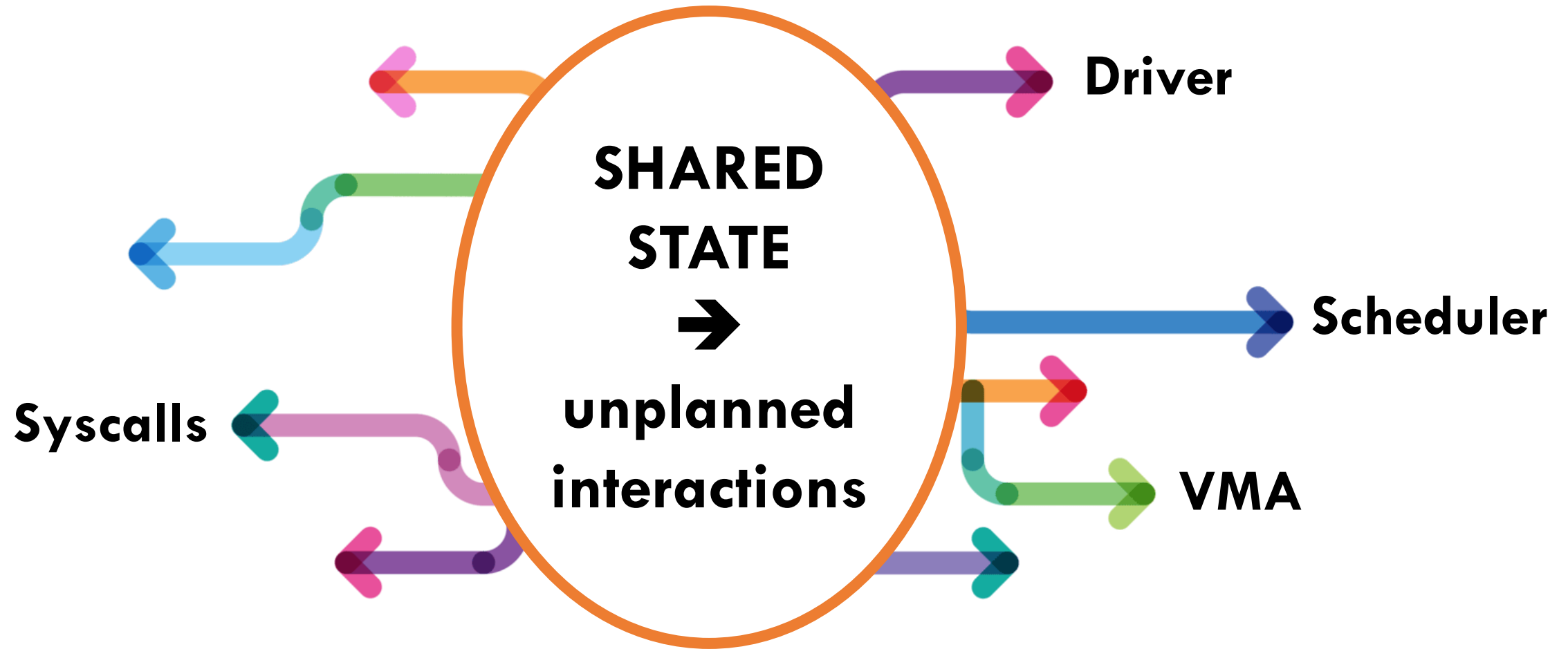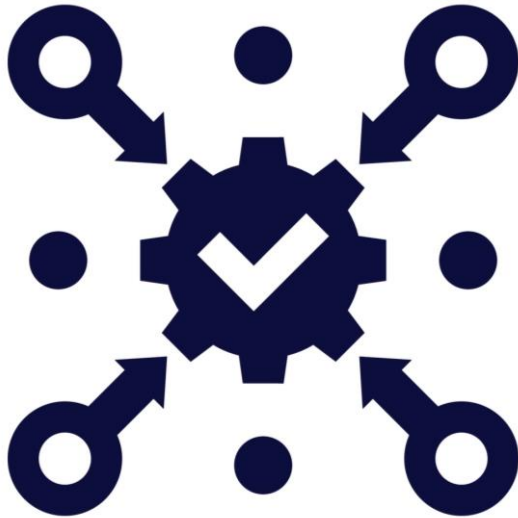
Inria

# Linux is a highly concurrent system



Driver

Scheduler

Syscalls

VMA

# Linux is a highly concurrent system



Driver

Scheduler

VMA

Syscalls

**SHARED STATE → unplanned interactions**

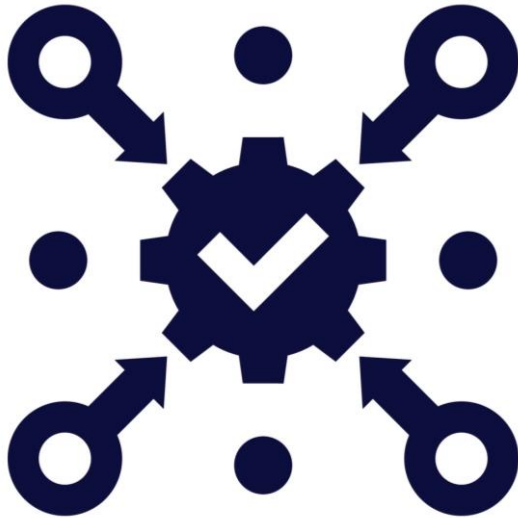# How to find concurrency bugs?

1. **Dynamic analysis (test & observe)**

   *Race reported on correct code*

   *Limited coverage*

# How to find concurrency bugs?

1. Dynamic analysis
(test & observe)

*Race reported on correct code*

*Limited coverage*

2. Static analysis)

*Limited to
lock-based code*

# Limitation of the state-of-the-art

```
1.      lock(&identifier)
2.      shared_var = xxx;
3.      unlock(&identifier)
```

```
1.      lock(&identifier)
2.      …
3.      unlock(&identifier)
4.      if(shared_var) {…}
```

**Traditional approach: locksets.**

*If 2 functions use the same lock*

*Then they likely run concurrently*

*If a variable is modified in a critical section and read outside ➔ bug*

**Lockless code not analyzed.**

# This talk



*Analyze concurrent lockless code.*

*First step: infer concurrency.*

# How to analyze lockless code?

**Observation 1: lockless code (often) uses barriers.**

```
1.    write a
2.    mfence
3.    write b
```

**CPU *must* write a before writing b**
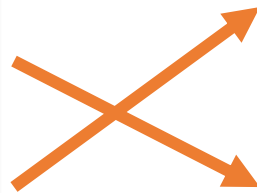
# Idea: use barriers

**Problem: barriers do not have an identifier**

```
1.      write a
2.      mfence
3.      write b
```

# Idea: pair barriers

**Observation 2: barriers run in pairs.**

```
1.    struct type *s = …;
2.    s.field = …;
3.    mfence
4.    s.initialized = 1;
```

```
1.        if(!s.initialized)
2.            return;
3.    mfence
4.    f(s.field);
```

# Idea: pair barriers
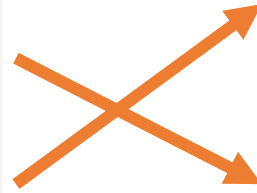
**Observation 2: barriers run in pairs.**

```
1.      struct type *s = …;
2.      s.initialized = 1;
3.      s.field = …;
```

```
1.      if(!s.initialized)
2.              return;
3.      mfence
4.      f(s.field);
```

# Idea: pair barriers using shared objects

## Observation 3: avoid aliasing using types

```
1.    struct type *s = …;
2.    s.field = …;
3.    mfence
4.    s.init = 1;
```

```
1.        if(!x.initialized)
2.            return;
3.    mfence
4.    f(x.field);
```

```
2.    write (type, field)
3.    mfence
4.    write (type, init)
```
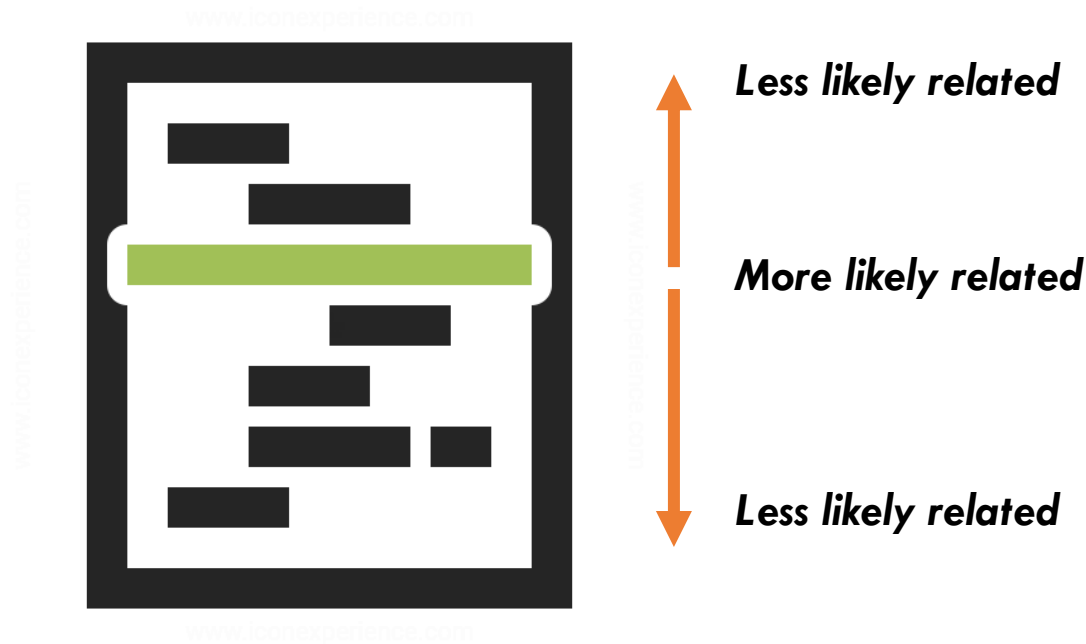
```
1.    read (type, init)
3.    mfence
4.    read (type, field)
```

**We call (type of struct, field name) a « shared object ».**

# Idea: pair barriers using closest shared objects

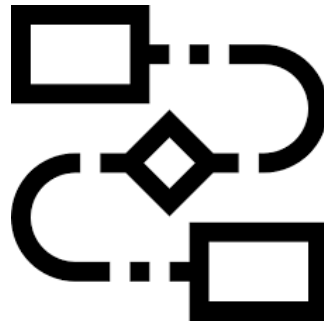**Observation 4: code related to a barrier is close to the barrier.**



Less likely related

More likely related

Less likely related

# Implementation

**Per barrier**
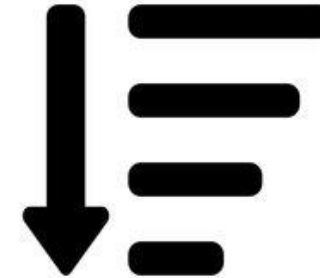
**Per barrier**

**1. Keep files with barriers**

**2. Produce Control Flow Graphs**
*Limit to 10 statements before and 10 statement after each barrier*
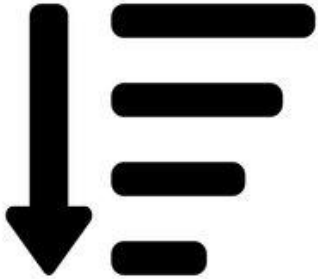
**3. Extract shared objects**

**4. Sort shared objects by distance**
*Distance = number of statements to/from barrier*

# Implementation: pairing

**Barrier 1**



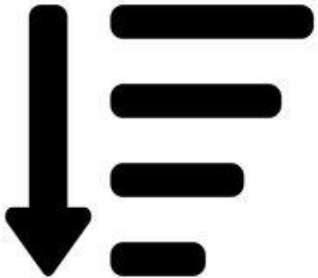**Barrier 2**



1. Have 2 shared objects o1 and o2 in common?
2. Objects are ordered by at least 1 barrier?
3. Pair with weight:

$$weight = o1.distance(barrier1) * o2.distance(barrier1)$$
$$+ o1.distance(barrier2) * o2.distance(barrier2)$$

4. If a barrier is paired with 2 barriers, keep the pairing with lowest weight.

# Use case: check ordering constraints

**A barrier is only useful if:**
**before(barrier 1, a) ⇔ after(barrier 2, a)**

```
1.      write a
2.      mfence
3.      write b
```

```
1.      read b
2.      mfence
3.      read a
```

```
1.      mfence
2.      read b  ⚠
3.      read a
```

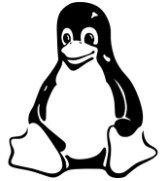|  | New b | Old b |
|---|---|---|
| New a | ✓ | ✓ |
| Old a |  | ✓ |

# Use case: remove unneeded barriers

**If a barrier is unpaired, is it needed?**

```
1.      write a
2.      mfence
3.      compare_and_swap(…)
```

```
1.      write a
2.      mfence
3.      function_with_barrier_semantics(…)
```
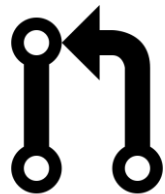
# Evaluation

Entire Linux code base

12 new bugs found & fixed

*Could have led to serious and hard-to-debug crashes*

50 uneeded barriers detected

Patches have been merged.

# Misplaced access in the RPC interface

```
1 void xprt_complete_rqst (...) {
2       req->rq_private_buf.len = ....;
3       smp_wmb();
4       req->rq_reply_bytes_recvd = copied;
5 }
6
7 static void call_decode (...) {
8 +     if (!req->rq_reply_bytes_recvd)
9 +             goto out;
10      smp_rmb();
11 -    if (!req->rq_reply_bytes_recvd)
12 -            goto out;
13      req->rq_rcv_buf.len = req->rq_private_buf.len;
```

# Racy re-read in the socket interface

```
1 int reuseport_add_sock(...) {
2       reuse->socks[reuse->num_socks] = sk ;
3       smp_wmb();
4       reuse->num_socks++;
5 }

469 struct sock *reuseport_select_sock(...) {
470     int socks = READ_ONCE(reuse->num_socks);
471     if (likely(socks)) {
472             smp_rmb();
…
487             reuseport_select_sock_by_hash(reuse);
```

# Racy re-read in the socket interface

```
1 int reuseport_add_sock(...) {
2       reuse->socks[reuse->num_socks] = sk ;
3       smp_wmb();
4       reuse->num_socks++;
5 }

469 struct sock *reuseport_select_sock(...) {
470     int socks = READ_ONCE(reuse->num_socks);
471
472
…
487
```

```
1 static struct sock *reuseport_select_sock_by_hash(…) {
…
5       while (reuse->socks[i]->sk_state == …) {
6               i++;
7               if (i >= reuse->num_socks)
8                       i = 0;
…               …
…       }
… }
```

# Racy re-read in the socket interface

```
1 int reuseport_add_sock(...) {
2       reuse->socks[reuse->num_socks] = sk ;
3       smp_wmb();
4       reuse->num_socks++;
5 }

469 struct sock *reuseport_select_sock(...) {
470     int socks = READ_ONCE(reuse->num_socks);
471
472
…
487
```
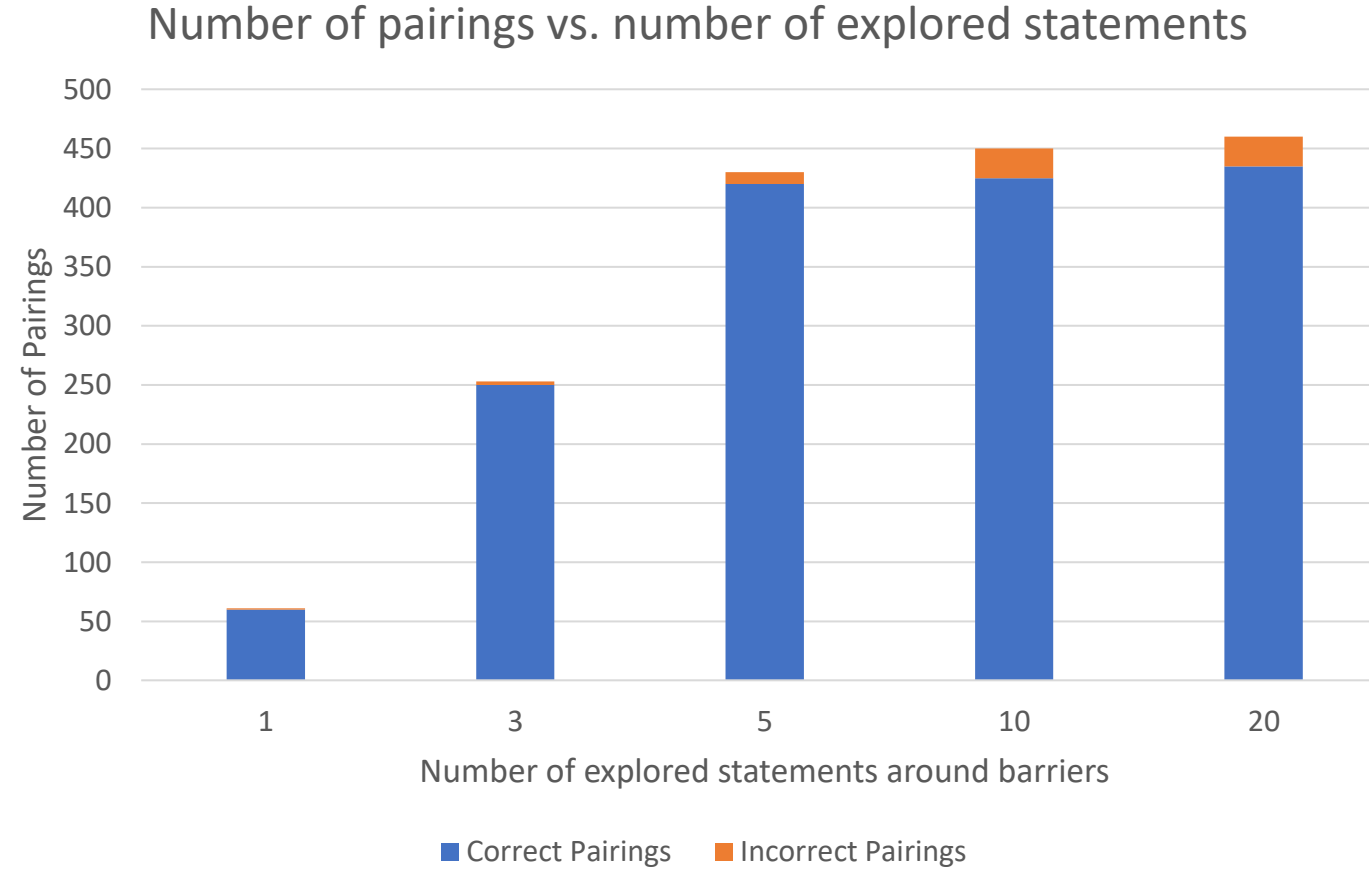
```
1 static struct sock *reuseport_select_sock_by_hash(…, socks){
…
5       while (reuse->socks[i]->sk_state == …) {
6               i++;
7               if (i >= socks)
8                       i = 0;
…               …
…       }
… }
```

# Removing unneeded barrier

```
1 static int rq_qos_wake_function (...) {
2       data->got_token = true ;
3 -     smp_wmb();
4       wake_up_process(data->task);
5 }
```

# Influence of tuning parameters



Number of pairings vs. number of explored statements

# Extension #1: avoid reports on benign races

**Dynamic analysers report races on correct code. We mark the code as safe.**

```
1 static int __pollwake (...) {
2        smp_wmb();
3 -      pwq->triggered = 1;
4 +      WRITE_ONCE(pwq->triggered, 1);
5        return ...;
6 }
7
8 static int poll_schedule_timeout (...) {
9 -      if (!pwq->triggered)
10+      if (!READ_ONCE(pwq->triggered))
11            rc = schedule_hrtimeout_range (...);
12      smp_store_mb(…); // equivalent to smp_mb
13 }
```

# Extension #2: avoid load/store tearing

**Reads/writes on some shared objects have to be atomic.**

```
1 void xprt_complete_rqst (...) {
2      req->rq_private_buf.len = …;
3      smp_wmb();
4 -    req->rq_reply_bytes_recvd = copied;
5 +    WRITE_ONCE(req->rq_reply_bytes_recvd,copied);
6 }
```

**Without WRITE_ONCE, the compiler could do:**

write 1$^{st}$ 32b of copied to rq_reply_bytes

write 2$^{nd}$ 32b of copied to rq_reply_bytes

(Actually happens with clang on arm64 CPUs!)

**… resulting in readers reading possibly partially written values.**

# Conclusion (of the 1$^{st}$ part of the talk)

**It is possible to infer concurrency by pairing barriers**
- Shared objects
- Distance

**Lockless code is error prone**

**Proof of concept: orderings – 12 bugs**

**Use it to check other concurrency bugs! (Use-after-free, …)**
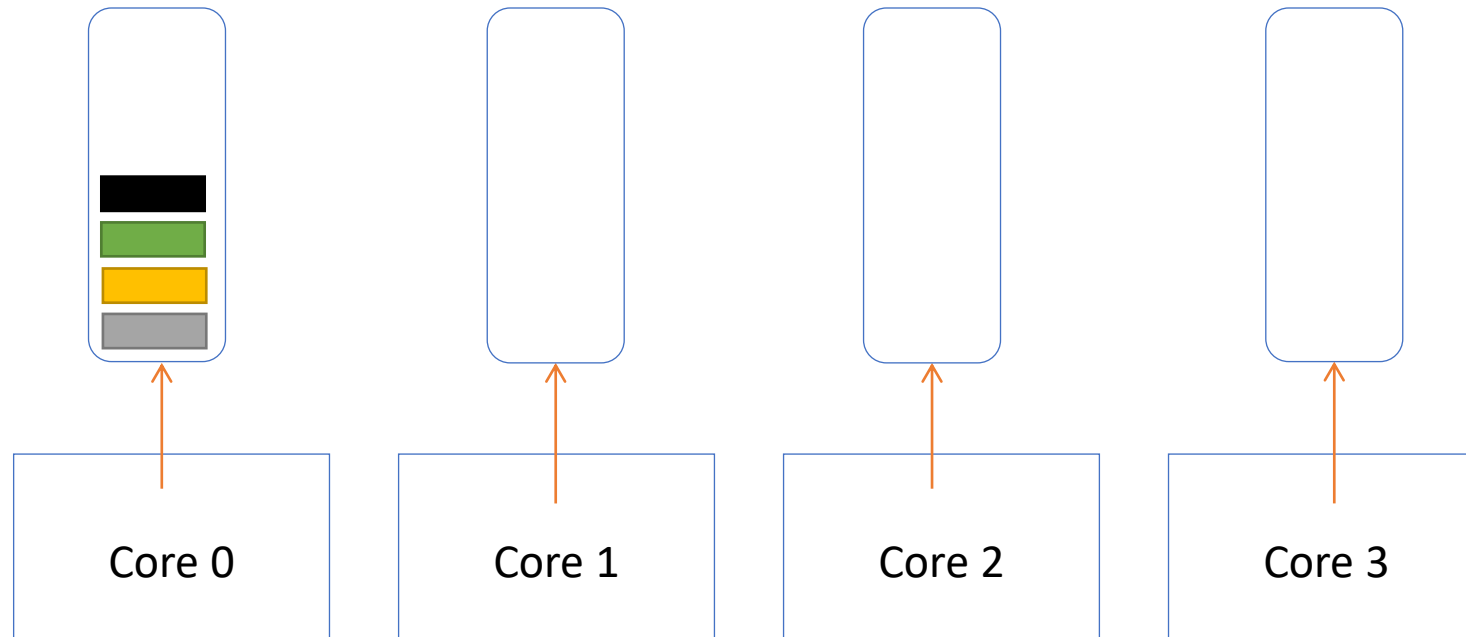
## Questions (before the 2$^{nd}$ part)?
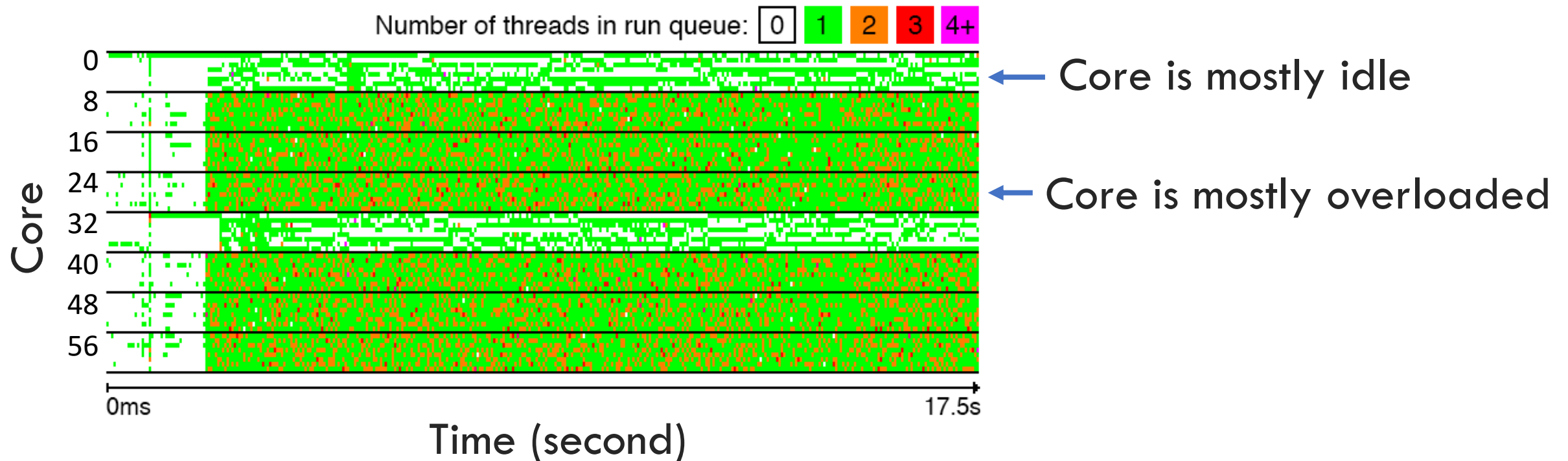
# Work conservation

## "No core should be left idle when a core is overloaded"



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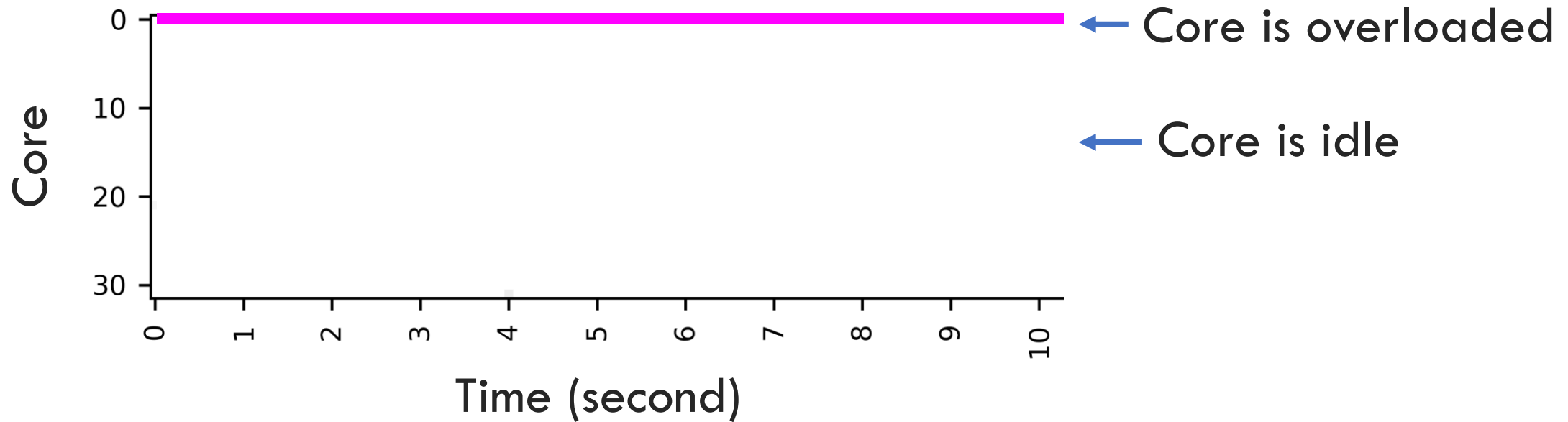core 0 is overloaded, other cores are idle

# Problem

## Linux (CFS) suffers from work conservation issues



Number of threads in run queue: 0 1 2 3 4+

Core is mostly idle

Core is mostly overloaded

Core

Time (second)

0ms    17.5s

*[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.

137x slowdown on HPC applications
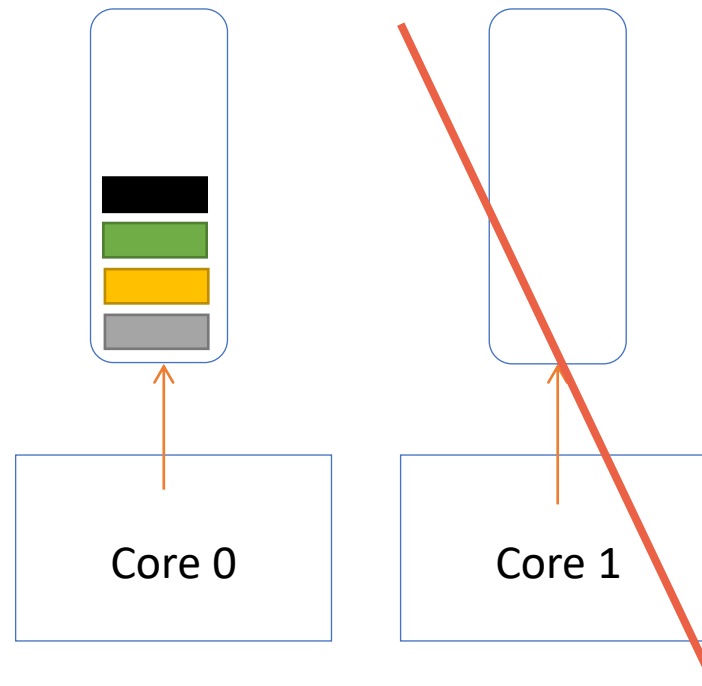23% slowdown on a database.
*[Lozi et al. 2016]*

# This talk

## Formally prove work-conservation

# Work Conservation Formally
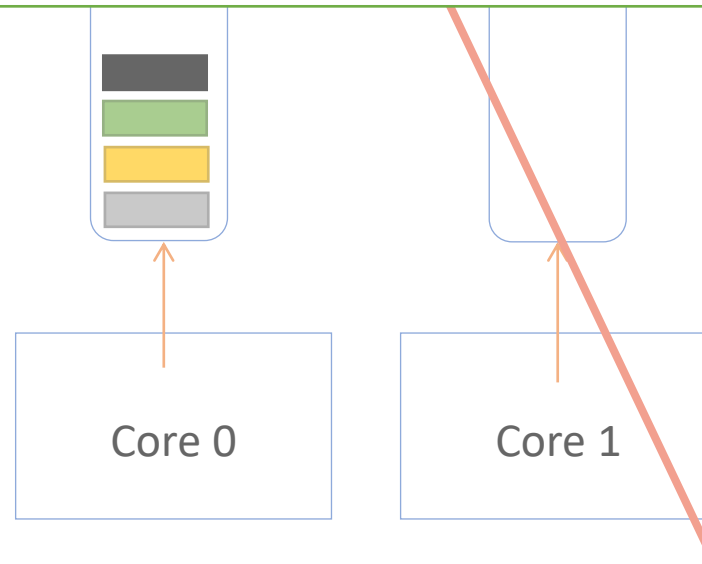
$$(\exists c \,.\, O(c)) \Rightarrow (\forall c' \,.\, \neg I(c'))$$

**If a core is overloaded, no core is idle**

Core 0

Core 1

# Work Conservation Formally

$$(\exists c \, . \, O(c)) \Rightarrow (\forall c' \, . \, \neg I(c'))$$
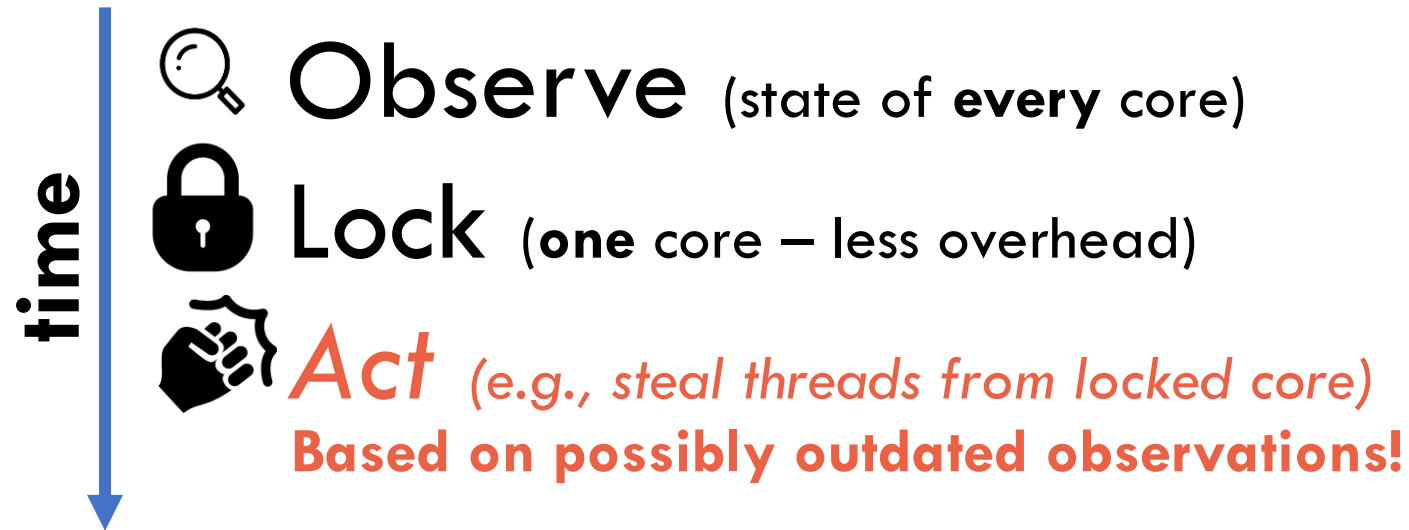
**Does not work for realistic schedulers!**

Core 0

Core 1

# Challenge #1

## Concurrent events & optimistic concurrency

# Challenge #1

## Concurrent events & optimistic concurrency

**time**

Observe (state of **every** core)

Lock (**one** core – less overhead)

*Act* (e.g., steal threads from locked core)
**Based on possibly outdated observations!**

# Challenge #1

## Concurrent events & optimistic concurrency



Runs load balancing

Core 0          Core 1          Core 2          Core 3

# Challenge #1

## Concurrent events & optimistic concurrency



Observes load
(no lock)

# Challenge #1

## Concurrent events & optimistic concurrency



**Ideal scenario**: no change since observations

Core 0

Core 1

Core 2

Core 3

Locks **busiest**

# Challenge #1

**Concurrent events & optimistic concurrency**

Core 0

Core 1

Core 2

Core 3

**(Fail to)
Steal from busiest**

# Challenge #1

**Concurrent events & optimistic concurrency**

time

🔍 Observe

🔒 Lock

✊ *Act*

**Based on possibly outdated observations!**

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

# Concurrent Work Conservation Formally

**Definition of overloaded with « failure cases »:**

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

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

# Concurrent Work Conservation Formally

$$\exists c \,.\, (O(c) \wedge \neg fork(c) \wedge \neg unblock(c) \ldots)$$
$$\Rightarrow$$
$$\forall c' \,.\, \neg(I(c') \wedge \ldots)$$

# 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)**
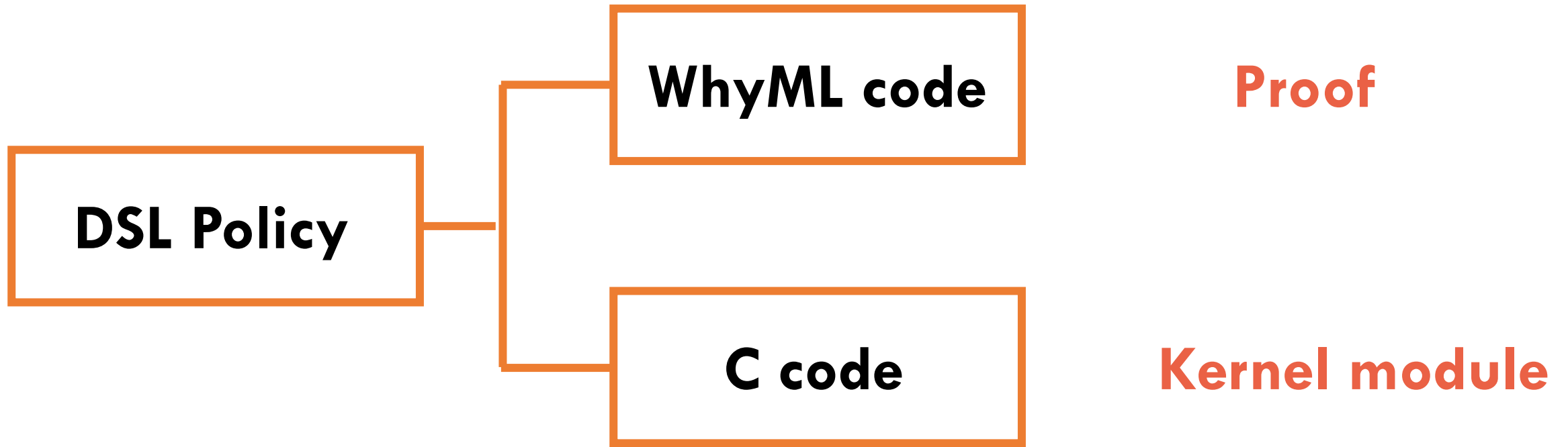
# DSL advantages

**Trade expressiveness for expertise/knowledge:**

**Robustness**: (static) verification of properties

**Explicit concurrency**: explicit shared variables

**Performance**: efficient compilation

# DSL-based proofs

```
                    ┌──────────────────┐
              ┌─────┤   WhyML code     │        Proof
┌───────────┐ │     └──────────────────┘
│ DSL Policy ├─┤
└───────────┘ │     ┌──────────────────┐
              └─────┤    C code        │     Kernel module
                    └──────────────────┘
```

**DSL Policy**

**WhyML code** — Proof

**C code** — Kernel module

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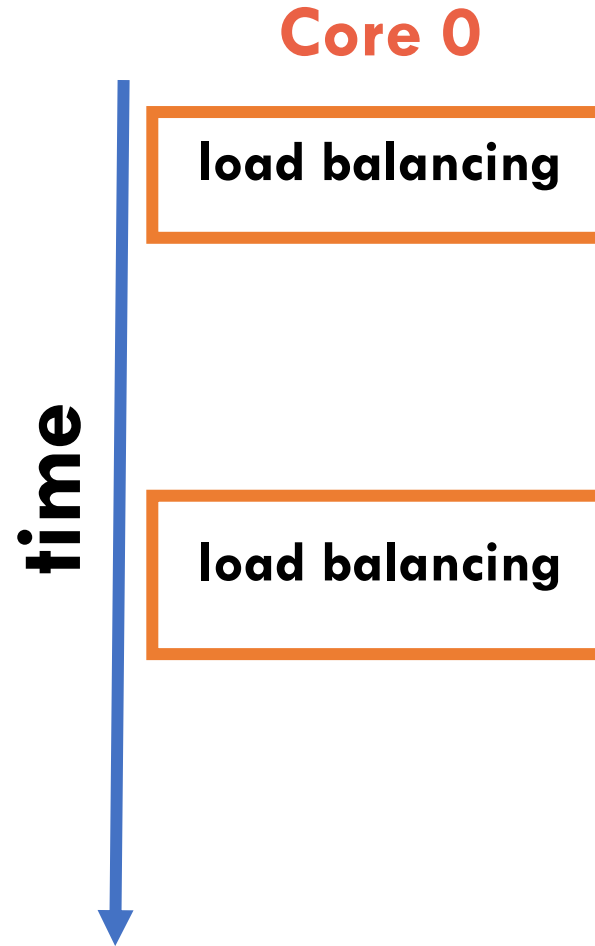Easy learn and to compile to WhyML and C

# DSL-based proofs

**Proof on all possible
interleavings**

# DSL-based proofs

**Proof on all possible interleavings**

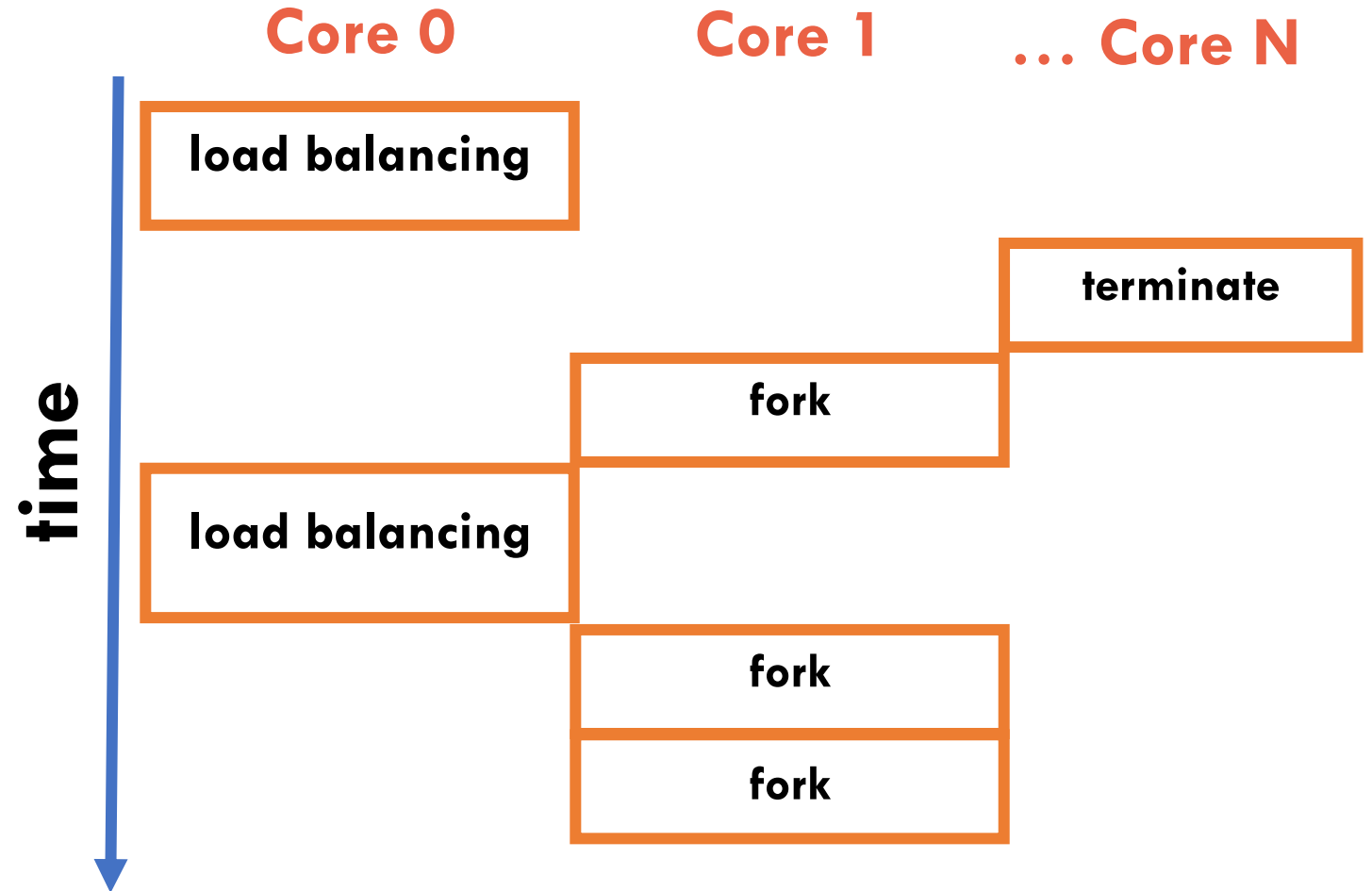**Split code in blocks (1 block = 1 read or write to a shared variable)**

**Core 0**

time

load balancing

load balancing

# DSL-based proofs
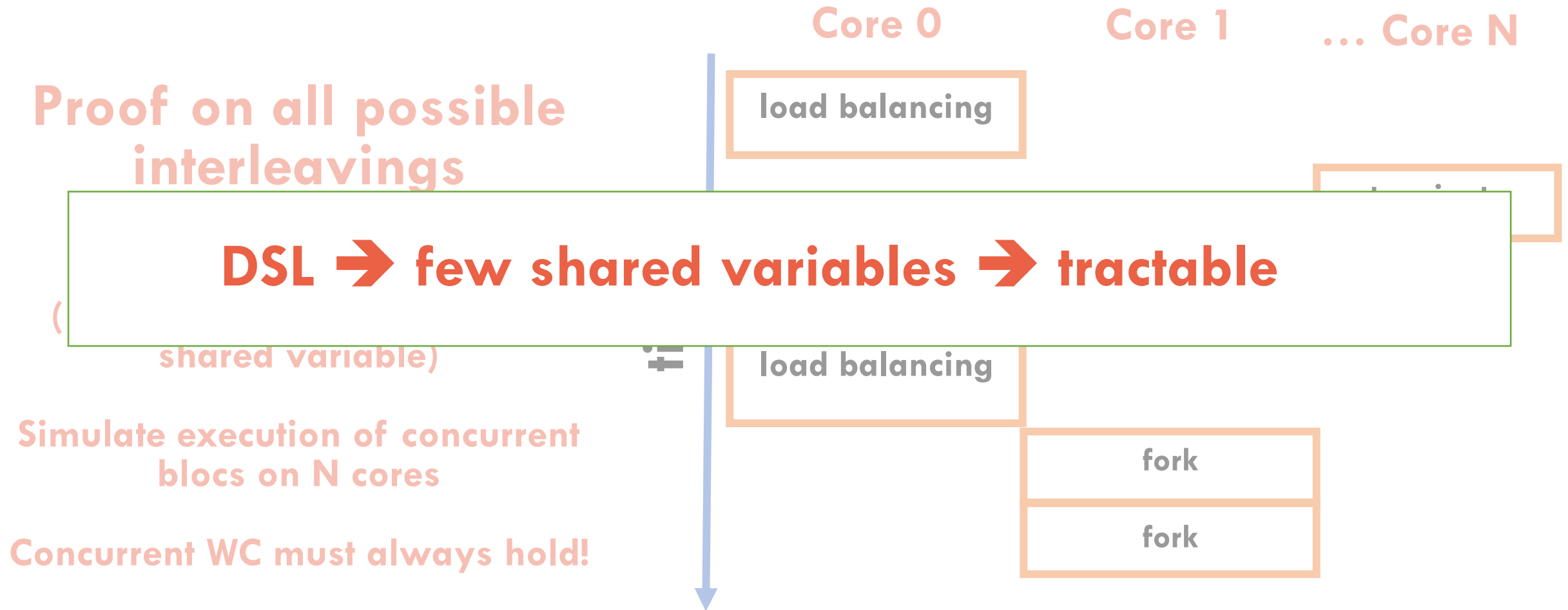
**Proof on all possible interleavings**

**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**

**Core 0**     **Core 1**     **… Core N**

time

| load balancing |
| terminate |
| fork |
| load balancing |
| fork |
| fork |

# DSL-based proofs

**Proof on all possible interleavings**

load balancing

DSL ➜ few shared variables ➜ tractable

shared variable)

load balancing

**Simulate execution of concurrent blocs on N cores**

fork

fork

**Concurrent WC must always hold!**

# Evaluation

## CFS-CWC (365 LOC)
**Hierarchical CFS-like scheduler**

## CFS-CWC-FLAT (222 LOC)
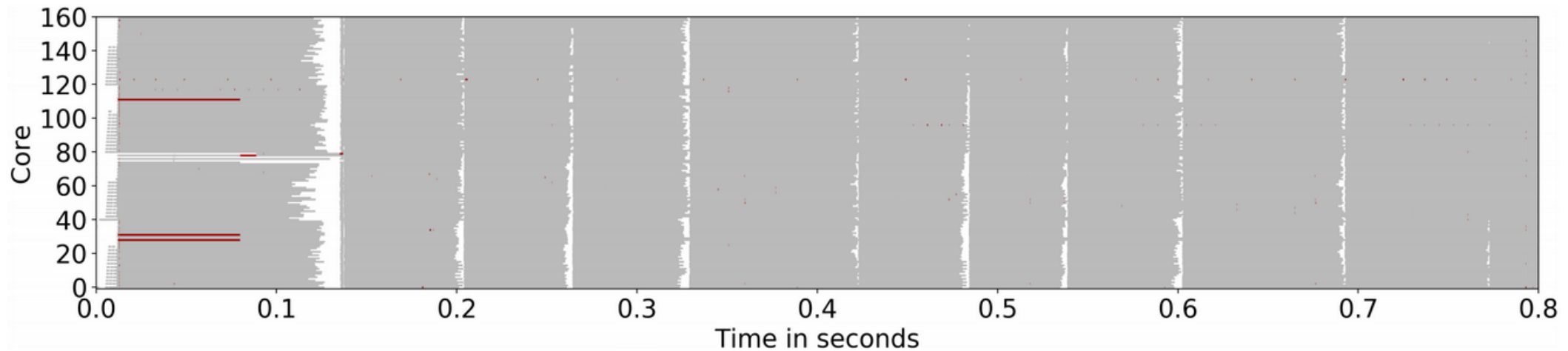**Single level CFS-like scheduler**

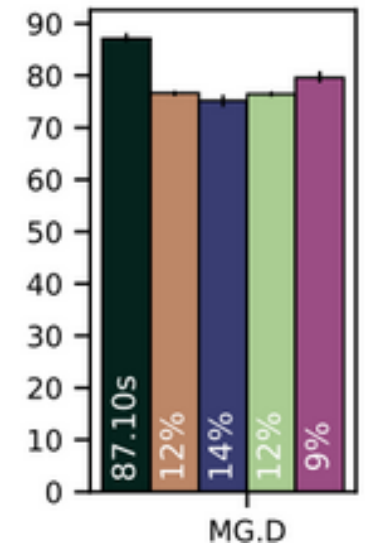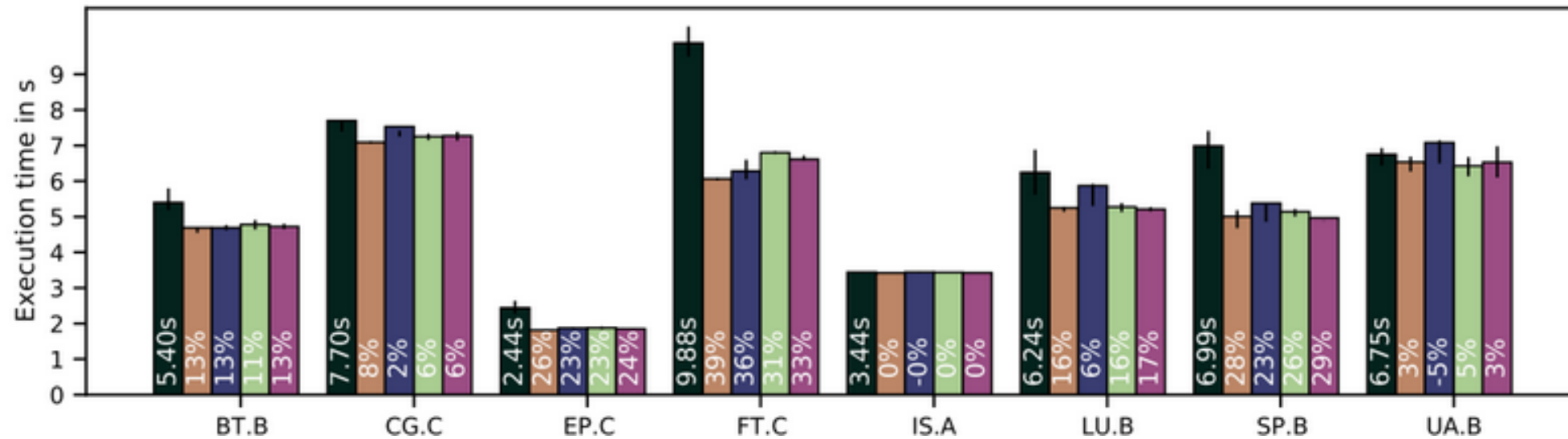## ULE-CWC (244 LOC)
**BSD-like scheduler**

# Less idle time

Execution with vanilla CFS.



Execution with CFS-CWC.

# Comparable or better performance



NAS benchmarks (lower is better)

# Conclusion

Work conservation: not straighforward!
 … new formalism: *concurrent* work conservation!

Complex concurrency scheme
 …proofs made tractable using a DSL.

Performance: similar or better than CFS.