# CertiKOS: An Extensible Architecture for Building Certified Concurrent OS Kernels
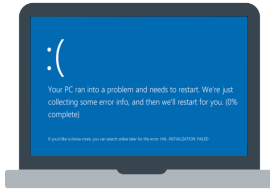
Ronghui Gu, Zhong Shao, Hao Chen, Xiongnan (Newman) Wu, Jieung Kim, Vilhelm Sjöberg, David Costanzo
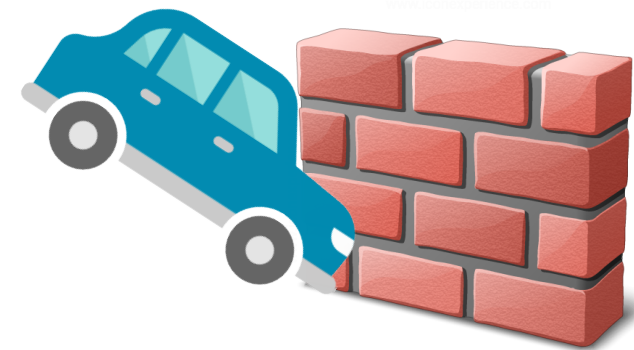
Yale University

OS Kernel

OS Kernel
error

"Complete formal verification is the only known way to guarantee that a system is free of programming errors.

— seL4 [SOSP'09]

seL4
[SOSP'09]

Verve
[PLDI'10]

Ironclad
[OSDI'14]

mCertiKOS
[POPL'15]

FSCQ
[SOSP'15]

CoGENT
[ASPLOS'16]

■ ■ ■

seL4
[SOSP'09]

Verve
[PLDI'10]

mCertiKOS
[POPL'15]

Ironclad
[OSDI'14]

FSCQ
[SOSP'15]

CoGENT
[ASPLOS'16]

. . .

verified sequential kernels

seL4
[SOSP'09]

Verve
[PLDI'10]

Ironclad
[OSDI'14]

mCertiKOS
[POPL'15]

FSCQ
[SOSP'15]

CoGENT
[ASPLOS'16]

■ ■ ■

verified software stacks

seL4
[SOSP'09]

Verve
[PLDI'10]

Ironclad
[OSDI'14]

mCertiKOS
[POPL'15]

FSCQ
[SOSP'15]

CoGENT
[ASPLOS'16]

■ ■ ■

verified sequential file systems

shared-memory concurrency?

seL4
[SOSP'09]

Verve
[PLDI'10]

Ironclad
[OSDI'14]

mCertiKOS
[POPL'15]

FSCQ
[SOSP'15]

CoGENT
[ASPLOS'16]

■ ■ ■

# seL4

" Proofs about concurrent programs are hard, much harder than proofs about sequential programs. "

seL4
[SOSP'09]

Verve
[PLDI'10]

hard

Ironclad
[OSDI'14]

mCertiKOS
[POPL'15]

CoGENT
[ASPLOS'16]

FSCQ
[SOSP'15]

■  ■  ■

# FSCQ

[SOSP'15]

hard

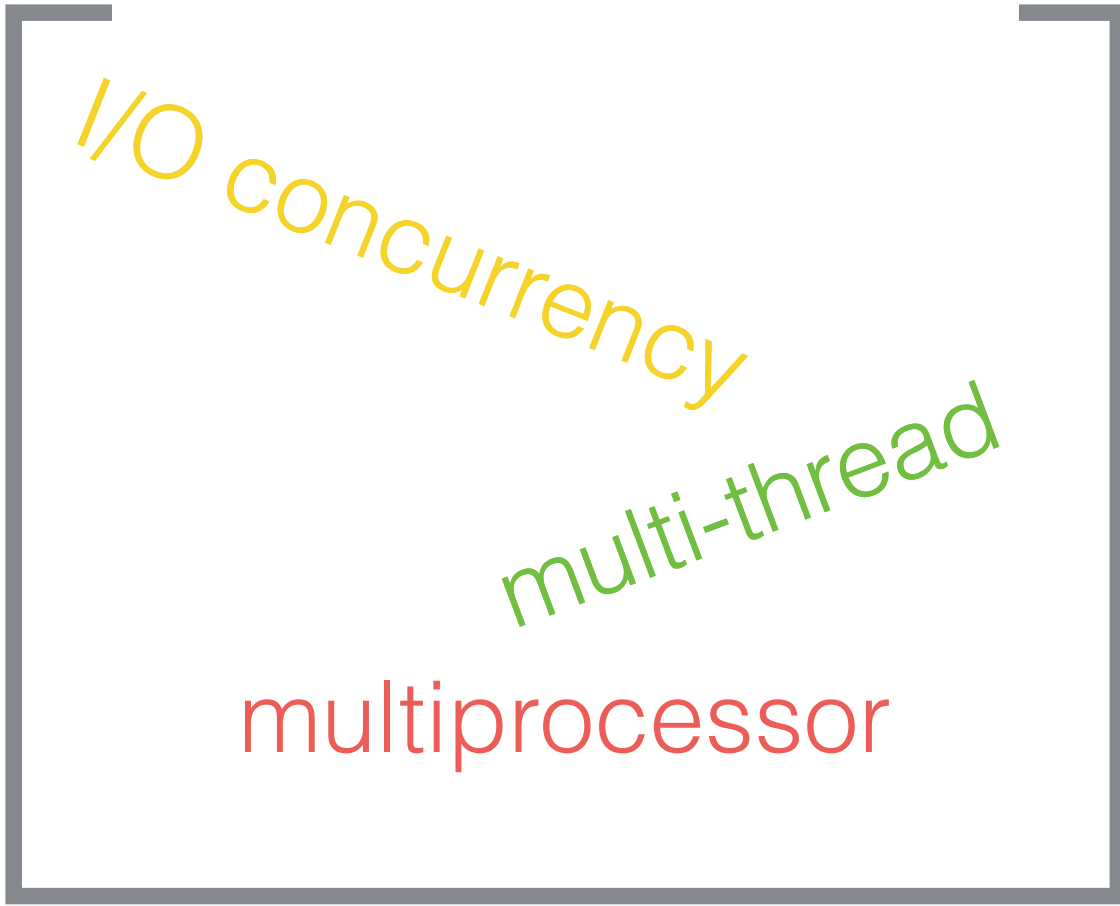" […]multiprocessor support, which may require global changes […] "

# FSCQ

hard
└ global changes

" […]multiprocessor support, which may require global changes […] "

hard

glol

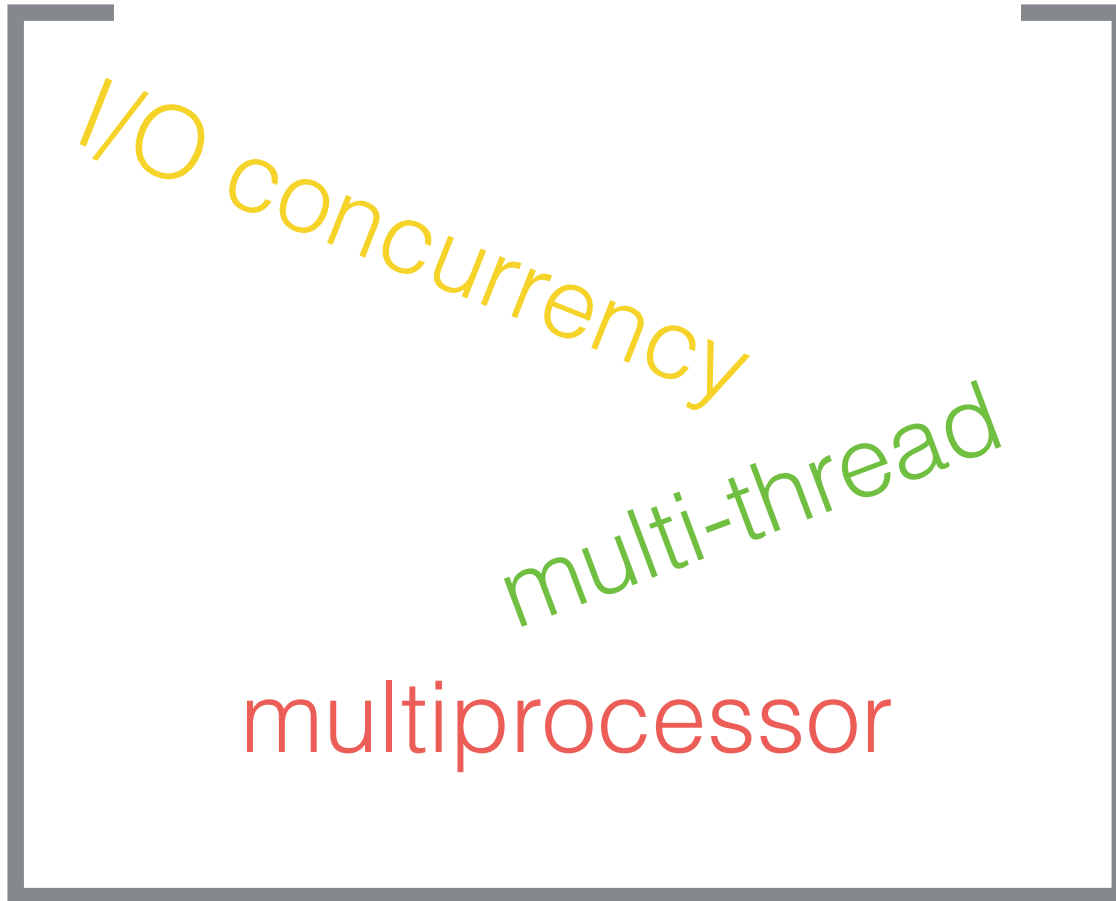I/O concurrency

multi-thread

multiprocessor

I/O concurrency

multi-thread

multiprocessor

hard

global change

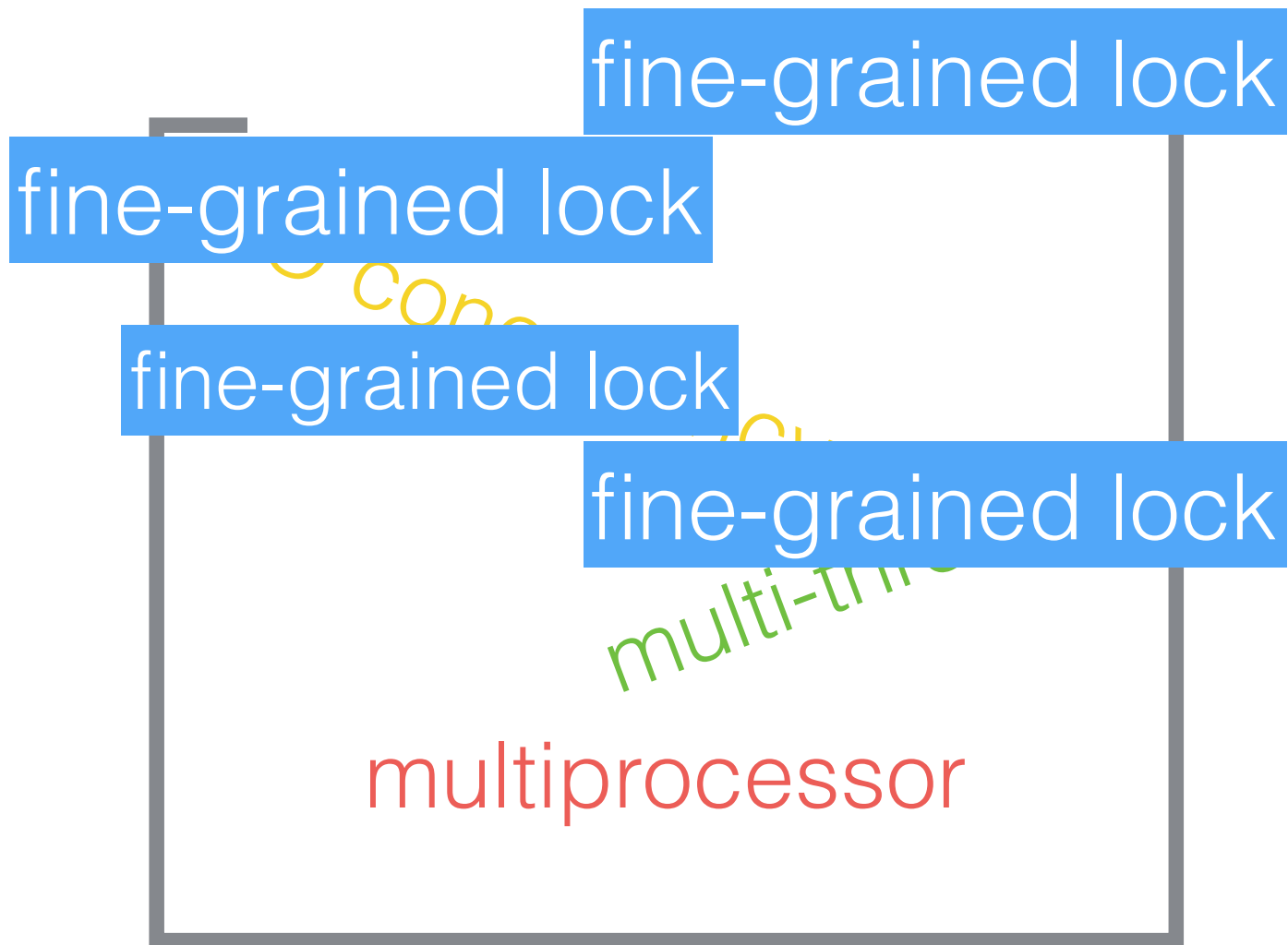I/O concurren

multi-thread

multiprocesso

# Big Lock

multi-threaded

multiprocessor

hard...
glob...
I/O
mu...
mu...

fine-grained lock

fine-grained lock

fine-grained lock

fine-grained lock

multiprocessor

multi-thr

hard

glo

I/O

mu

mu

# S.Peters et al.

hard
├ glol
├ I/O
mu
mu

" the verification to a kernel version with fine-grained locking will far exceed the cost already paid for verifying the single core version. "

# S.Peters et al.

"
the verification to a kernel version with fine-grained locking will far exceed the cost already paid for verifying the single core version.
"

hard
global cha
I/O concu
multi-threa
multiproce
fine-grain

# What to prove?

functional correctness

liveness system calls will eventually return

hard
glob
I/O
mul
mul
fine

concurrent OS kernel

hard
├─ global change
├─ I/O concurren
│  multi-thread
│  multiprocesso
├─ fine-grained l
└─ liveness

concurrent OS kernel

hard

glob

I/O

mul

mu

fine

live

C

Compiler

Asm

⋈

Asm

hard

global changes

I/O concurrency

multithread

**cost**

asm&C

compiler

hard

- global changes
- I/O concurrency
  multi-thread
  multiprocessor
- fine-grained lock
- liveness
- asm&C
- compiler
- cost

# CertiKOS

mC2, the first formally verified concurrent OS kernel with fine-grained locks.

glob
I/O
mul
mu
fine
live
asr
cor
cos

# CertiKOS

- glob
- I/O
  mul
  mu
- live
- asr
- cor
- cos

mC2, the first formally verified concurrent OS kernel with fine-grained locks.

contributions **CertiKOS** hard

⊢mC2
⊢fine-grained lock

⊢glob
⊢I/O
mu
mu

both functional correctness
and liveness

⊢live
⊢asr
⊢cor

⊢cos

# CertiKOS

- mC2
- fine-grained lock
- liveness

both functional correctness
and liveness

hard
- glob
- I/O
- mul
- mu

- asr
- con

- cos

contributions

## CertiKOS

hard

- mC2
- fine-grained lock
- liveness

- glob
- I/O
  mu
  mu

certified concurrent layers

- asn
- cor

- cos

# CertiKOS

hard

├ mC2
├ fine-grained lock
└ liveness

├ glob
├ I/O
mul
mu

reuses sequential verification techniques.

├ asr
├ con

└ cos

certified concurrent layers

contributions

## CertiKOS

- mC2
- fine-grained lock
- liveness
- global changes

reuses sequential verification techniques.

hard

- I/O
- mul
- mu

- asr
- cor

- cos

certified concurrent layers

## CertiKOS

contributions

- mC2
- fine-grained lock
- liveness
- reuse of techs

hard

- I/O
  mu
  mu

- asr
- cor
- cos

handles all 3 kinds of concurrency

certified concurrent layers

# CertiKOS

contributions

- mC2
- fine-grained lock
- liveness
- reuse of techs

- I/O concurrency

multi-thread
multiprocessor

hard

- asr
- cor

- cos

handles all 3 kinds of concurrency

certified concurrent layers

# CertiKOS

- mC2
- fine-grained lock
- liveness
- reuse of techs
- mix of 3

hard

C  6100 LOC

Asm  400 LOC

- asr
- cor
- cos

contributions

# CertiKOS

hard

- mC2
- fine-grained lock
- liveness
- reuse of techs
- mix of 3
- asm&C

C 6100 LOC

Asm 400 LOC

- cor
- cos

# CertiKOS

contributions

- mC2
- fine-grained lock
- liveness
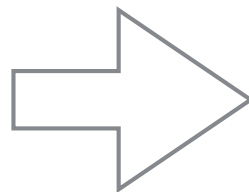- reuse of techs
- mix of 3
- asm&C

C 6100

CompCertX

Asm ⋈ 400 Asm

hard

cor

cos

# CertiKOS

contributions

- mC2
- fine-grained lock
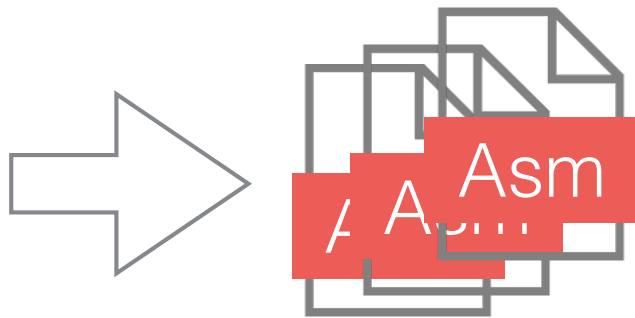- liveness
- reuse of techs
- mix of 3
- asm&C
- compiler

hard

C  6100

CompCertX

Asm  ⋈  Asm  400

cos

# CertiKOS

contributions

- mC2
- fine-grained lock
- liveness
- reuse of techs
- mix of 3
- asm&C
- CompCertX

hard

6100    400

C    Asm

verified

~~model~~ Coq    ~~SMT~~
~~checking~~    ~~solver~~
machine-checkable proof

cos

# CertiKOS

contributions

- mC2
- fine-grained lock
- liveness
- reuse of techs
- mix of 3
- asm&C
- CompCertX

machine-checkable
proof

hard

6100    400

C    Asm

verified

Asm

cos

# CertiKOS

contributions
- mC2
- fine-grained lock
- liveness
- reuse of techs
- mix of 3
- asm&C
- CompCertX

hard

6100    400

C    Asm

verified

le

Asm Asm Asm ⟹ executable

cos

# CertiKOS

contributions
- mC2
- fine-grained lock
- liveness
- reuse of techs
- mix of 3
- asm&C
- CompCertX

hard

6100 400

C Asm

~~verified~~

certified

le ⇒ Asm Asm ⇒ executable

cos

# CertiKOS

hard

mCertiKOS     1 py

[POPL'15]

+ extensions   0.5 py

+ device       0.5 py

[PLDI'16]

+ concurrency  2 py

cos

contributions

**CertiKOS**

- mC2
- fine-grained lock
- liveness
- reuse of techs
- mix of 3
- asm&C
- CompCertX
- extensibility
- certified
- cost

# new technical contributions

certified concurrent layers

logical log + hardware scheduler + environment context

push/pull model

multicore machine lifting

certified sequential layers

certified objects



specification of modules to trust

certified sequential layers

certified **sequential** layers
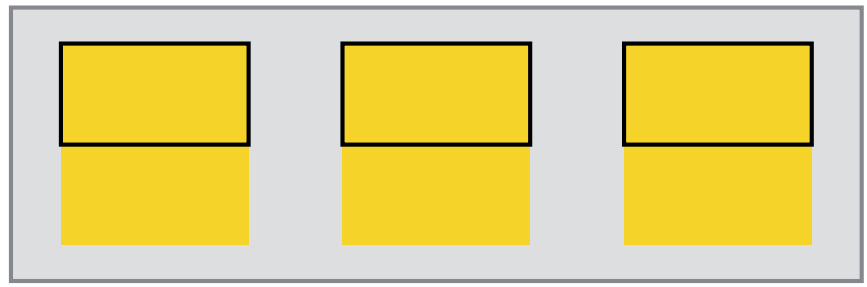
abs-state

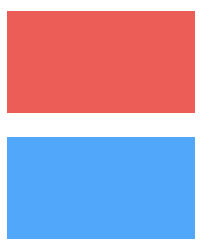certified sequential layers

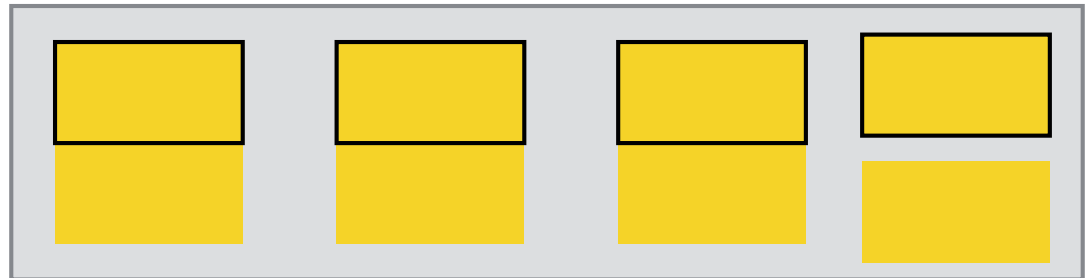abs-state

primitives
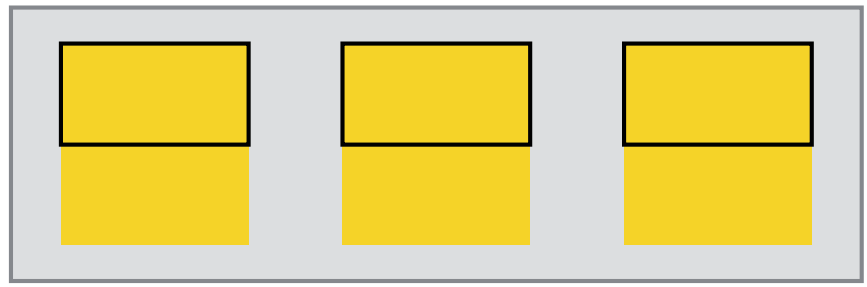
code

memory

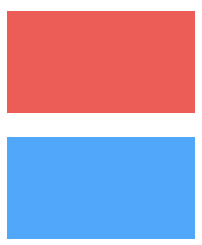implementation

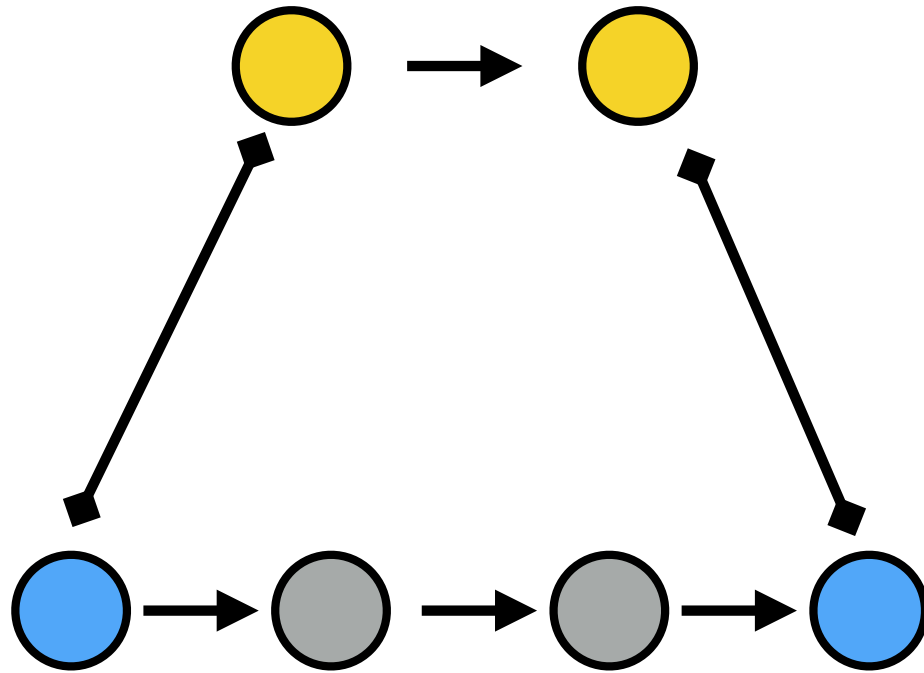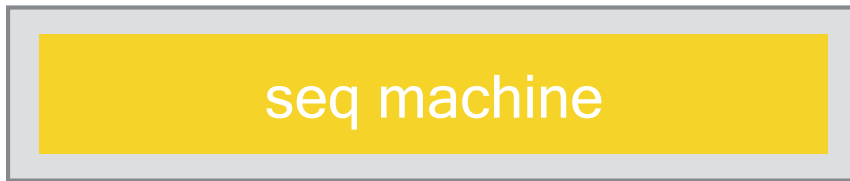specification

implementation

specification

implementation

# simulation proof

specification



⊔⊓

implementation
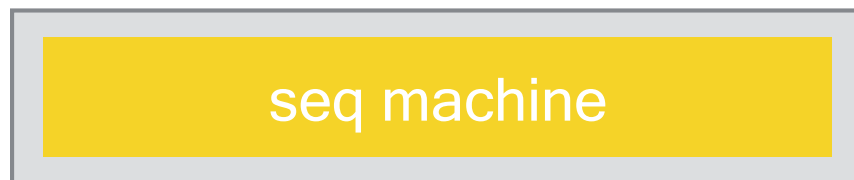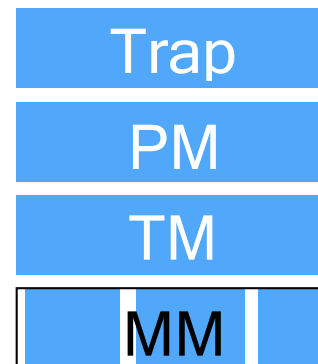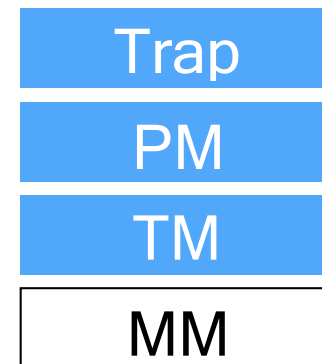
verify a sequential kernel

[POPL'15]

# kernel

code

seq machine

Trap

PM

TM

MM

seq machine

# memory management

Trap
PM
TM
MM

seq machine

trap

Trap

proc

PM

thread

TM

mem

seq machine

Trap
PM
TM
MM

| Trap |
|:---:|
| PM |
| TM |
| MM |

# verified sequential kernel

| trap |
|:---:|
| proc |
| thread |
| mem |
| seq machine |

## TSysCall Layer

(pe, ikern, ihost, ipt, AT, PT, ptp, pbit, kctxp, Htcbp, Htqp, cid, chanp, uctxp, npt, hctx, vmst)

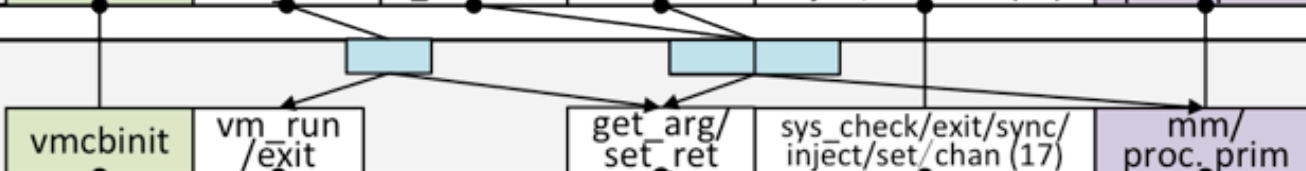| thread_wakeup/kill/sleep/yield | | pt_read | get/set_uctx | | palloc/free | cid_get |
|---|---|---|---|---|---|---|
| sys_chan_send/recv/wait/check | | sys_yield | sys_get_exit_reason | | | sys_get_eip |
| sys_check_shadow/pending_event | | sys_proc_create | | sys_set_seg | | sys_inject |
| sys_get_exit_io_width/port/rep/str/write/eip | | | sys_set_intcept_int | | | sys_npt_instr |
| vmcbinit | pagefault_handler | sys_reg_get/set | | sys_sync | sys_run | vm_exit |

## TSysCall Layer
(mm/proc/virt.abs)

| vmcbinit | sys_run/ vm_exit | PageFault _Handler | sys_yield | sys_check/exit/sync/ inject/set/chan (17) | mm/ proc. prim |
|---|---|---|---|---|---|

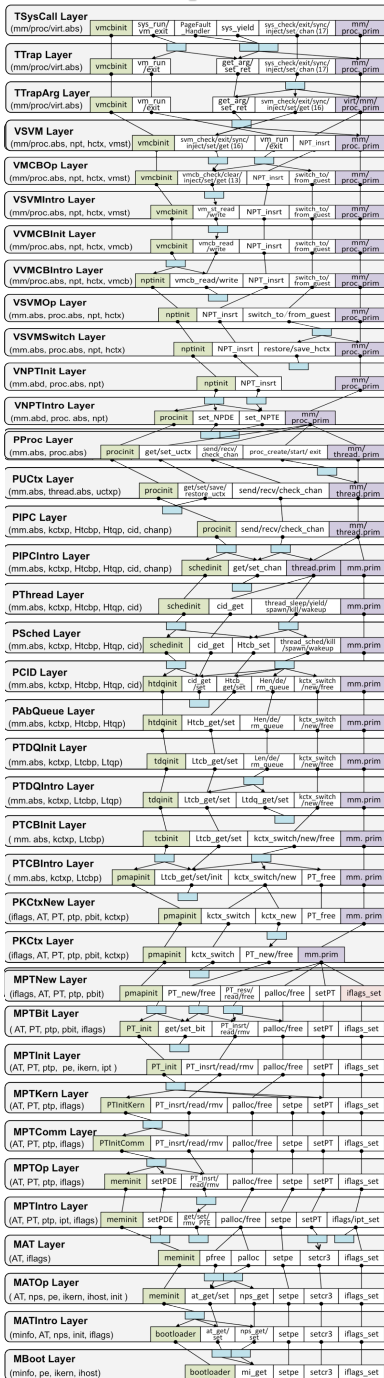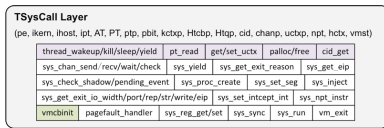## TTrap Layer
(mm/proc/virt.abs)

| vmcbinit | vm_run /exit | | get_arg/ set_ret | sys_check/exit/sync/ inject/set/chan (17) | mm/ proc. prim |
|---|---|---|---|---|---|

1 person year
(cost for tool construction excluded)

trap

proc

thread

mem

seq machine

cost

contributions

extensibility

cost

| | |
|---|---|
| | Trap |
| VM | PM |
| | TM |
| | MM |

trap

proc

thread

mem

seq machine

# contributions

- extensibility
- cost

| | |
|---|---|
| Trap | |
| VM | |
| PM | |
| TM | |
| MM | |

| trap |
|---|

| proc |
|---|

| thread |
|---|

| mem |
|---|

| seq machine | |
|---|---|

contributions

| | |
|---|---|
| | Trap |
| | VM |
| | PM |
| | TM |
| | MM |

trap

virt

VM

proc

thread

extensibility

mem

cost

seq machine

contributions

seq machine

extensibility is the key to support

concurrency

# support concurrency

## contributions

| trap |
| --- |

| virt |
| --- |

| proc |
| --- |

| thread |
| --- |

| mem |
| --- |

seq machine

multicore machine

contributions

reuse

trap

virt

proc

thread

mem

CPU-local machine

multicore machine

contributions

reuse

trap

virt

proc

thread

mem

spin-lock

CPU-local machine

multicore machine

contributions

trap

virt

proc

thread-local machine

reuse

mix of 3

thread

mem

spin-lock

CPU-local machine

multicore machine

contributions

mC2

reuse
mix of 3

trap

virt

proc

thread-local machine

thread

mem

spin-lock

CPU-local machine

multicore machine

# certified concurrent layers

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

# certified concurrent layers

local objects

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

certified concurrent layers

atomic objects

logical log

a sequence of events

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

# certified concurrent layers

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

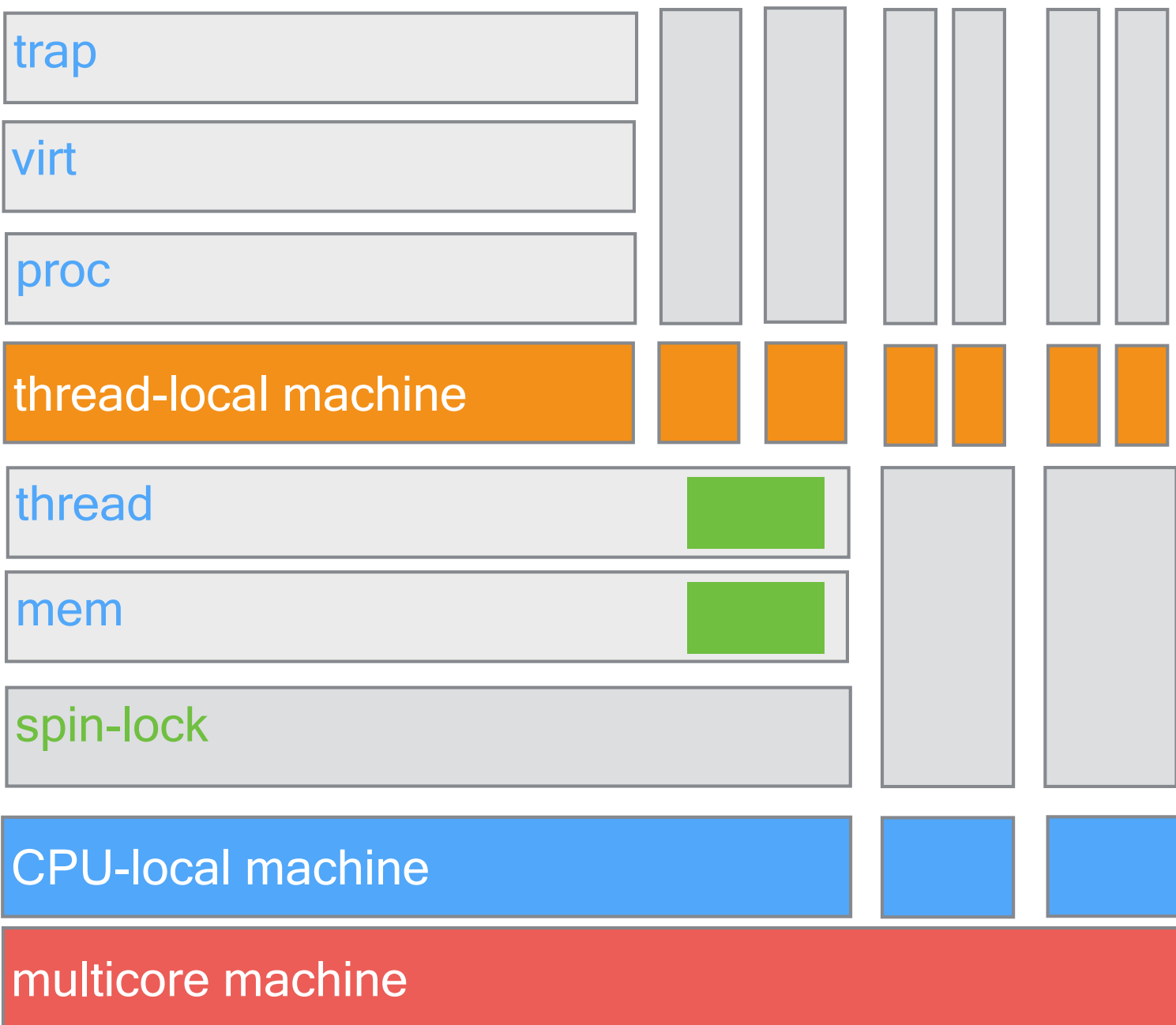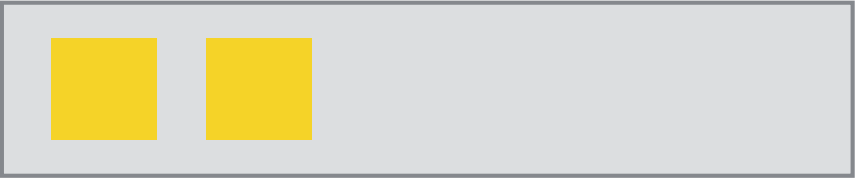# certified concurrent layers

trap

virt

proc

thread

thread
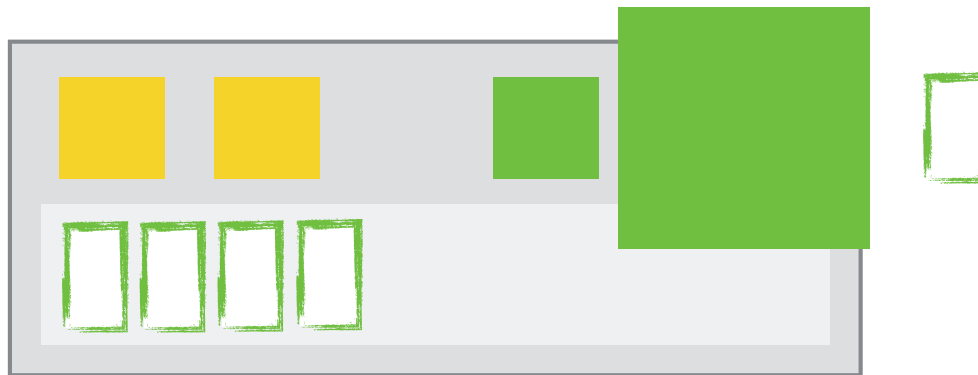
mem

spin-lo

CPU-lo

multico

# certified concurrent layers

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

share

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

fine-grained lock

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

fine-grained lock

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

trap

virt

proc

thread-local machine

thread

mem

spin-lock

CPU-local machine

multicore machine

# step 0: raw x86 multicore model
## assume sequential consistency

0.a

CPU0 — atom — share →

CPU1 — private — atom →

1.a

multicore machine

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

# step 0: raw x86 multicore model

trap

virt

proc

thread

thread

mem

spin-lo

CPU0 — atom ———————— share →

CPU1 ——— private | atom ————————→

logical log  [ 0.a ] [ 1.a ]

CPU-lo

multicore machine

# step 0: raw x86 multicore model
## non-determinism

| | | | |
|---|---|---|---|
| trap | | | |
| virt | | | |
| proc | | | |
| thread | | | |
| thread | | | |
| mem | | | |
| spin-lo | | | |
| CPU-lo | | | |

CPU0 — [atom] ——————— [share] →

CPU1 — [private][atom] —————— →

[ 0.a  1.a ]

**multicore machine**

step 1: hardware scheduler

purely logical

# step 1: hardware scheduler
## purely logical

atom1

pull1

hardware scheduling

step 1: hardware scheduler

private1

atom2

purely logical

# step 1: hardware scheduler

trap

virt

...1 pull1 shared1 shared1 push...

hardware scheduling

private1 atom2

# step 2: push/pull model

trap

virt

CPU0    ▶ share ◀

h1    pull1    shared1    shared1    push

hardware scheduling

private1    atom2

step 2: **push/pull** model

trap

virt

CPU0    pull    share

pull1    shared1    shared1    push

h1

hardware scheduling

private1    atom2

step 2: **push/pull** model

trap

virt

CPU0 ▶ pull ▶ share ▶

n1

hardware scheduling

private1  atom2

pull1  shared1  shared1  push

# step 2: push/pull model

trap

virt

CPU0 ▶ pull ▶ share ▶

n1

hardware scheduling

private1

atom2

pull1

shared1

shared1

push

CPU1 ▶

# step 2: push/pull model

trap

virt

CPU0 ▶ pull ▶ share ▶ push ▶

h1    pull1    shared1    shared1    push

hardware scheduling

private1    atom2

# step 2: push/pull model

trap

virt

CPU0  ▶ pull ▶ share ▶ push ▶

...1   pull1   shared1   shared1   push

hardware scheduling

private1   atom2

trap

virt

h1

hardware scheduling

private1

atom2

pull1

shared1

shared1

push

# step 3: per-CPU machine

hardware scheduling

pull1 | shared1 | shared1 | push1

private1 | atom2 | private2

push1 | pull1 | shared1 | shared1 | push

hardware scheduling

private1 | atom2

# step 3: per-CPU machine

pull1 shared1 shared1 push1

hardware scheduling

0.a

private1 atom2 private2

0

push1

hardware scheduling

pull1 shared1 shared1 push

private1 atom2

# step 3: per-CPU machine



trap

virt

pull1  shared1  shared1  push1

0.a

private1  atom2

1  1  0

push1  pull1  1.a  d1  shared1  push

hardware scheduling

private1  atom2

# step 3: per-CPU machine

0.a

trap

virt

$\mathcal{E}$

1  1  1.a  0

...n1  environment context  pull1  shared1  shared1  push...

hardware scheduling

private1  atom2

trap

virt

h1

hardware scheduling

private1

atom2

pull1

shared1

shared1

push

# step 4: remove unnecessary *interleaving*

trap

virt

...1

pull1

shared1

shared1

push

hardware scheduling

private1

atom2

# step 4: remove unnecessary interleaving

trap

virt

...1 shared1 shared1 push

hardware scheduling pull1

private1 atom2

# step 4: remove unnecessary interleaving

trap

virt

...1

hardware scheduling

private1

atom2

pull1

shared1

shared1

push

contributions

trap

virt

▶ atom

| 0 | 0.a | 1 | 1 | 1.a | 0 |

...1

hardware scheduling

private1

atom2

pull1

shared1

shared1

push

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

trap
virt
proc
thread-local machine
thread
mem
spin-lock
CPU-local machine
multicore machine

# acq-lock specification

trap

virt

proc

thread

thread

mem

safely
pull

logical
copy

CPU-l

multic

spin-lock

# acq-lock specification

safely
pull

logical
copy

pull will
eventually return

trap

virt

proc

thread

thread

mem

CPU-lo

multico

spin-lock

# acq-lock specification

mutual
exclusion

logical
copy

liveness

spin-lock

trap

virt

proc

thread

thread

mem

CPU-lo

multico

# ticket lock

mutual exclusion + liveness

```
void acq_lock (uint i)
{
    uint t = ▶FAI_ticket (i)      FAI
                                  ticket

    while (▶get_now (i) != t)
    { }

    ▶ pull (i);
}
```

trap

virt

proc

thread

thread

mem

CPU-lo

multico

spin-lock

# mutual exclusion + liveness

```
void acq_lock (uint i)
{
  uint t = ▶FAI_ticket (i);

  while (▶get_now (i
  { }

  ▶ pull (i);
}
```

FAI
ticket

get
now

trap

virt

proc

thread

thread

mem

CPU-lo

multico

spin-lock

# mutual exclusion + liveness

```
void acq_lock (uint i)
{
    uint t = ▶ FAI_ticket (i);

    while ( ▶ get_now ( 
    { }

    ▶ pull (i);
}
```

get now

FAI ticket    get now

trap

virt

proc

thread

thread

mem

CPU-lo

multico

spin-lock

# mutual exclusion + liveness

```
void acq_lock (uint i)
{
    uint t = ▶ FAI_ticket (i);

    while ( ▶ get_now (i) != t)
    { }

    ▶ pull (i)  pull
}
```

FAI ticket | get now | get now

trap

virt

proc

thread

thread

mem

CPU-lo

multico

spin-lock

# mutual exclusion + liveness

```
void acq_lock (uint i)
{
    uint t = ▶FAI_ticket (i);

    while (▶get_now (i) != t)
    { }

    ▶ pull  (i);
}
```

FAI ticket    get now    get now    pull

trap

virt
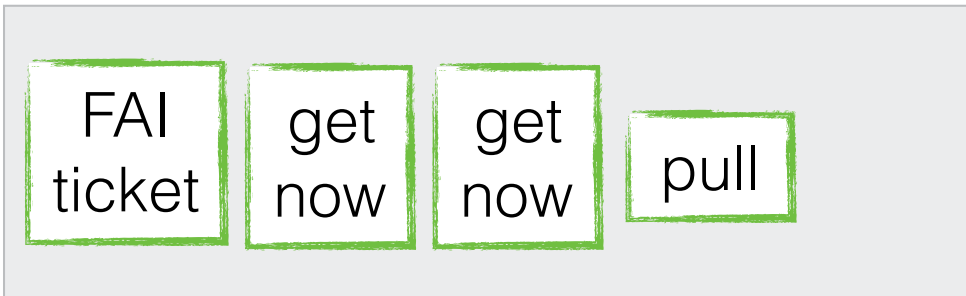
proc

thread

thread

mem

spin-lock

CPU-lo

multico

# mutual exclusion + liveness

```
void acq_lock (uint i)
{
  uint t = ▶FAI_ticket (i);

  while (▶get_now (i) != t)
  { }

  ▶ pull  (i);
}
```
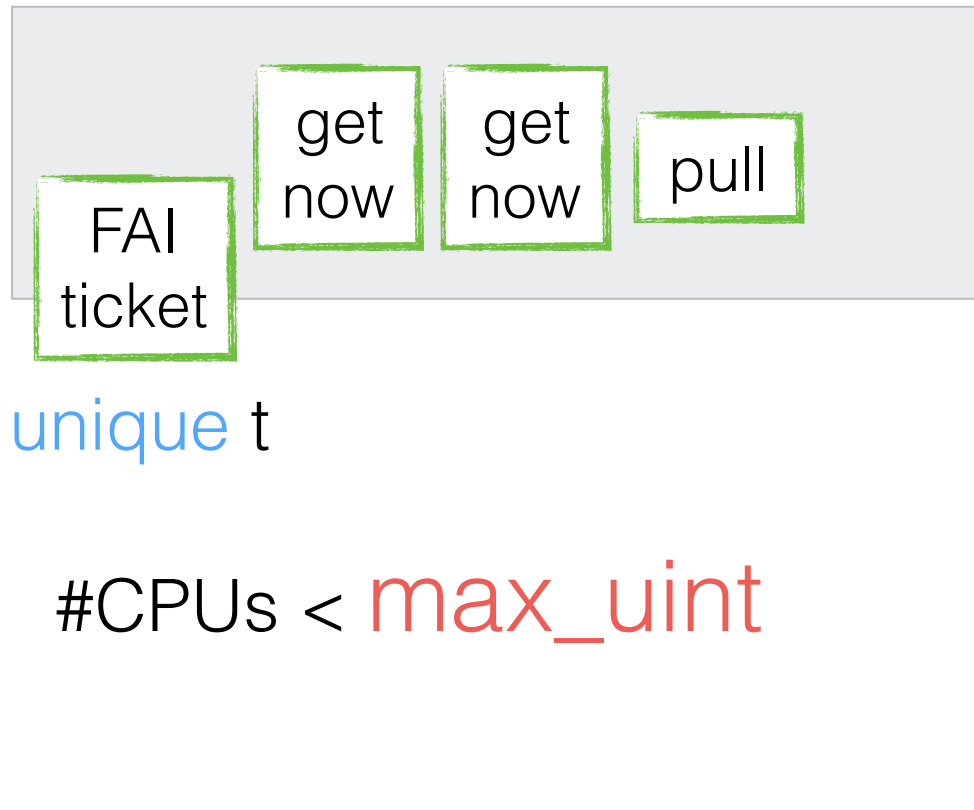
FAI ticket  get now  get now  pull

**unique** t

#CPUs < max_uint

spin-lock

trap

virt
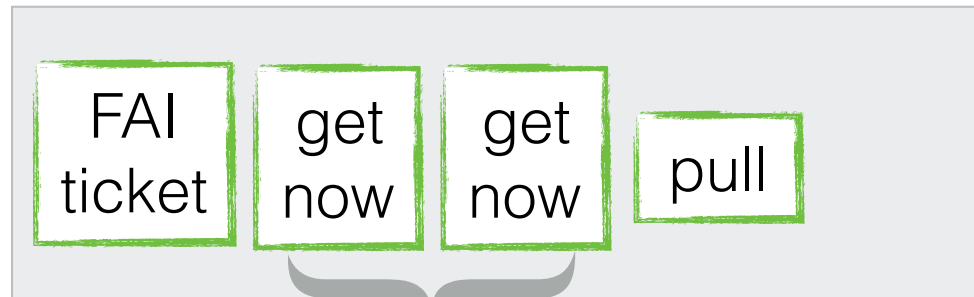
proc

thread

thread

mem

CPU-lo

multico

mutual exclusion + liveness

liveness

```
void acq_lock (uint i)
{
    uint t = ▶FAI_ticket (i);

    while (▶get_now (i) != t)
    { }

    ▶ pull (i);
}
```

FAI ticket    get now    get now    pull

bounded

#CPUs is bounded

a fair scheduler
lock holders will release lock

trap

virt

proc

thread

thread

mem

CPU-lo

multico

spin-lock

trap

virt

proc

thread

thread

mem

acq_lock

acq
lock

acq_lock

FAI
ticket

get
now

get
now

pull

CPU-lo

multico

spin-lock

trap

virt

proc

thread

thread

mem

acq
lock

acq_lock

acq_lock

CPU-lo

multic

FAI
ticket

get
now

get
now

pull

spin-lock

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

trap

virt

proc

thread-local machine

thread

mem

spin-lock

CPU-local machine

multicore machine

enq

local
memory

thread

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

enq

local
memory

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

thread

enq

logical
copy

shared
memory

thread

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

acq lock

enq

logical copy

shared memory

thread

trap

virt

proc

thread

mem

spin-lo

CPU-l

multico

acq
lock

enq

logical
copy

shared
memory

thread
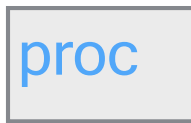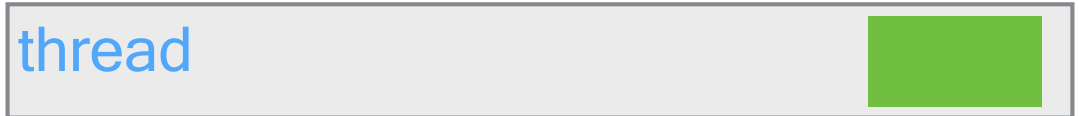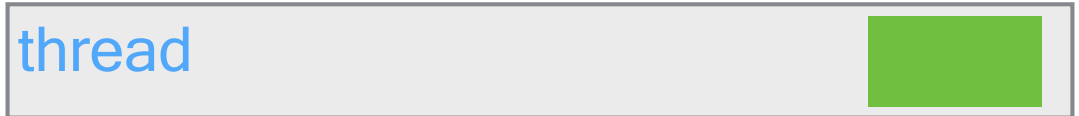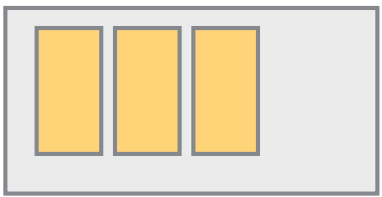
trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

logical
copy

shared
memory

acq
lock

enq

rel
lock

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multic

thread

trap

virt

proc

thread

enq

shared
memory

mem

spin-lo

CPU-l

thread

multico

enq

shared
memory

thread

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multic

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

contributions

trap

virt

proc

```
void yield ()
{
    uint t =  tid();

    ...
    ▶ enq  (t, rdq());


    uint s =  ▶ deq  (rdq());
    ...
    context_switch  (t, s)
}
```

thread

mem

spin-lo

CPU-lo

thread-local machine

multico

contributions



trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

asm&C
CompcertX

```
void yield ()
{
    uint t =  tid();
    …
    ▶enq  (t, rdq());

    uint s =  ▶deq  (rdq());

    context_switch s)

}
```

thread-local machine

contribution

software scheduler

mix of 3

yield

sleep

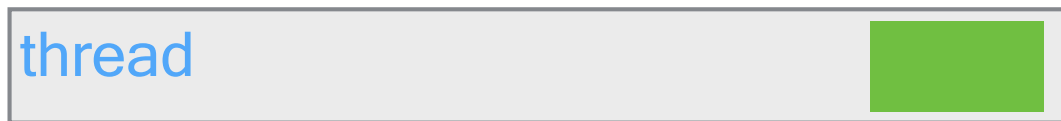wakeup

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

thread-local machine

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

trap

virt

thread

thread

mem

spin-lo

CPU-lo

multico

IPC

CV

proc

## evaluation:
## proof effort for concurrency(LOC)

top spec: 450

machine model: 943

intermediate spec: 40K

proof(concurrency): 50K

Coq &machine checkable
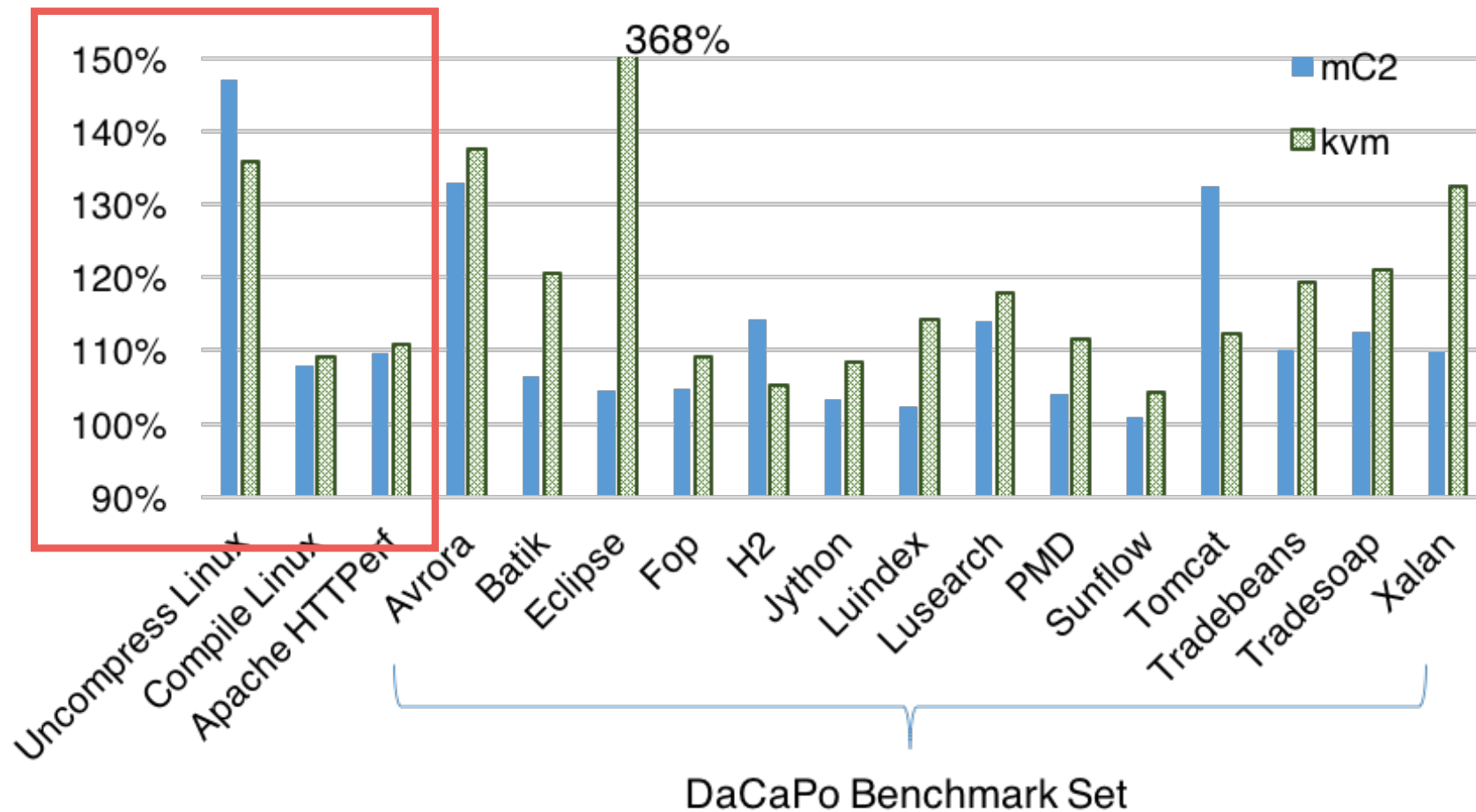
2 person year

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multic

# evaluation:
## performance
## mC2 is comparable with kvm



DaCaPo Benchmark Set
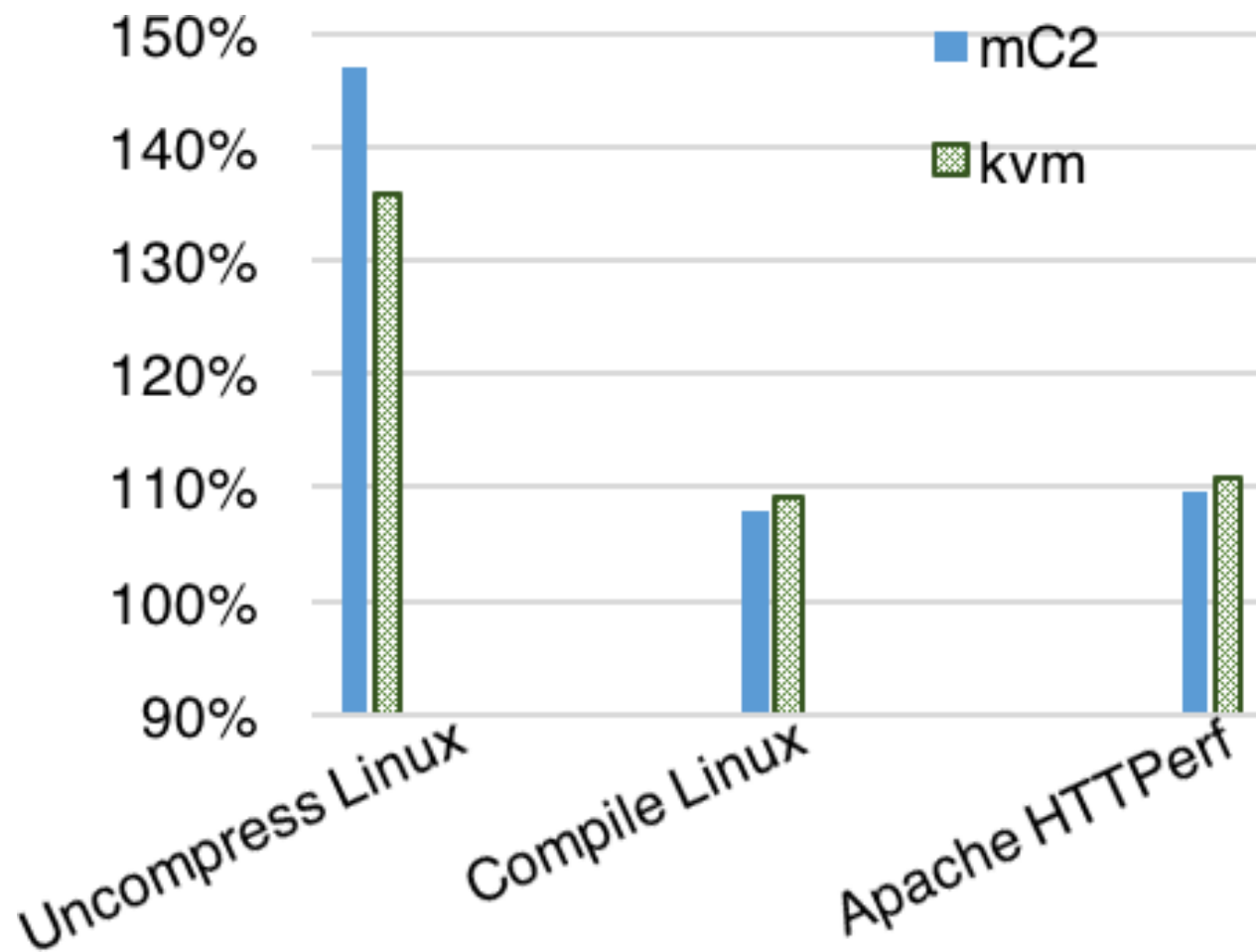
trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

# evaluation:
## performance

## mC2 is comparable with kvm



trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

# limitations & future work

bootloader

assembler of CompCert

machine model is in the TCB

sequential consistency
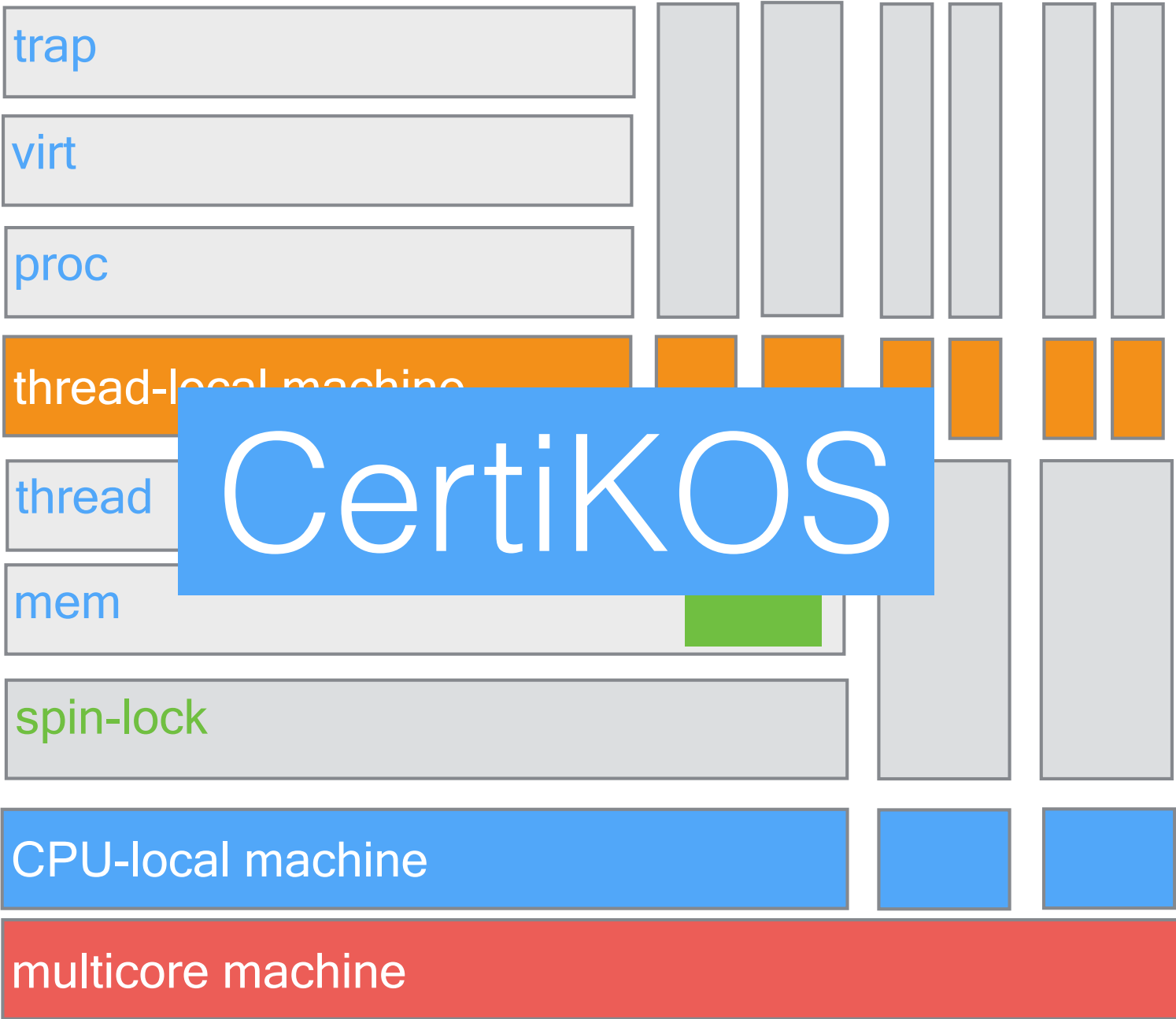
file system & network stack

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

## contributions

- mC2
- fine-grained lock
- liveness
- reuse
- mix of 3
- asm&C
- CompcertX
- extensibility
- Coq &machine checkable
- 2 person year

## CertiKOS

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico

# CertiKOS

new technical contributions

certified concurrent layers

logical log + hardware scheduler + environment context

push/pull model

multicore machine lifting

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico