

CertiKOS: An Extensible Architecture for Building Certified Concurrent OS Kernels

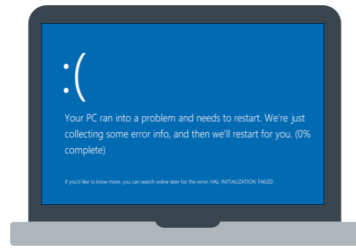
Ronghui Gu, Zhong Shao, Hao Chen, Xiongnan (Newman) Wu, Jieung Kim,
Vilhelm Sjöberg, David Costanzo

Yale University

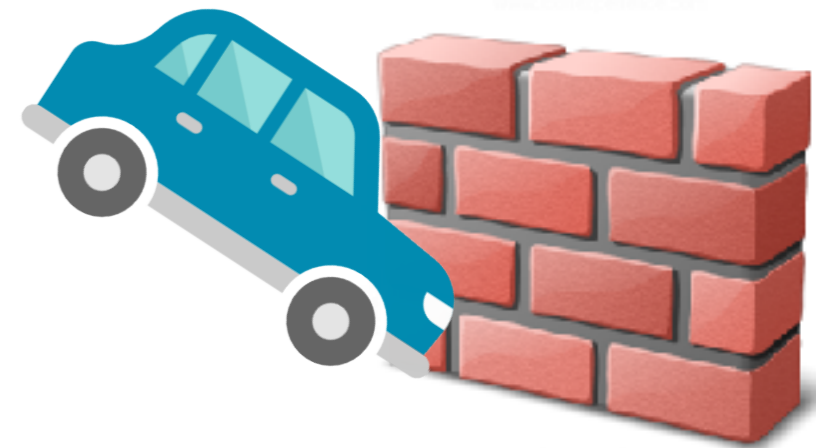


OS Kernel





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]

Ironclad
[OSDI'14]

mCertikOS
[POPL'15]

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?



You shall not pass!

shared-memory concurrency?

seL4

[SOSP'09]

Verve

[PLDI'10]

Ironclad

[OSDI'14]

mCertikOS

[POPL'15]

FSCQ

[SOSP'15]

CoGENT

[ASPLOS'16]



seL4

[SOSP'09]

“ 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

[SOSP'15]

hard

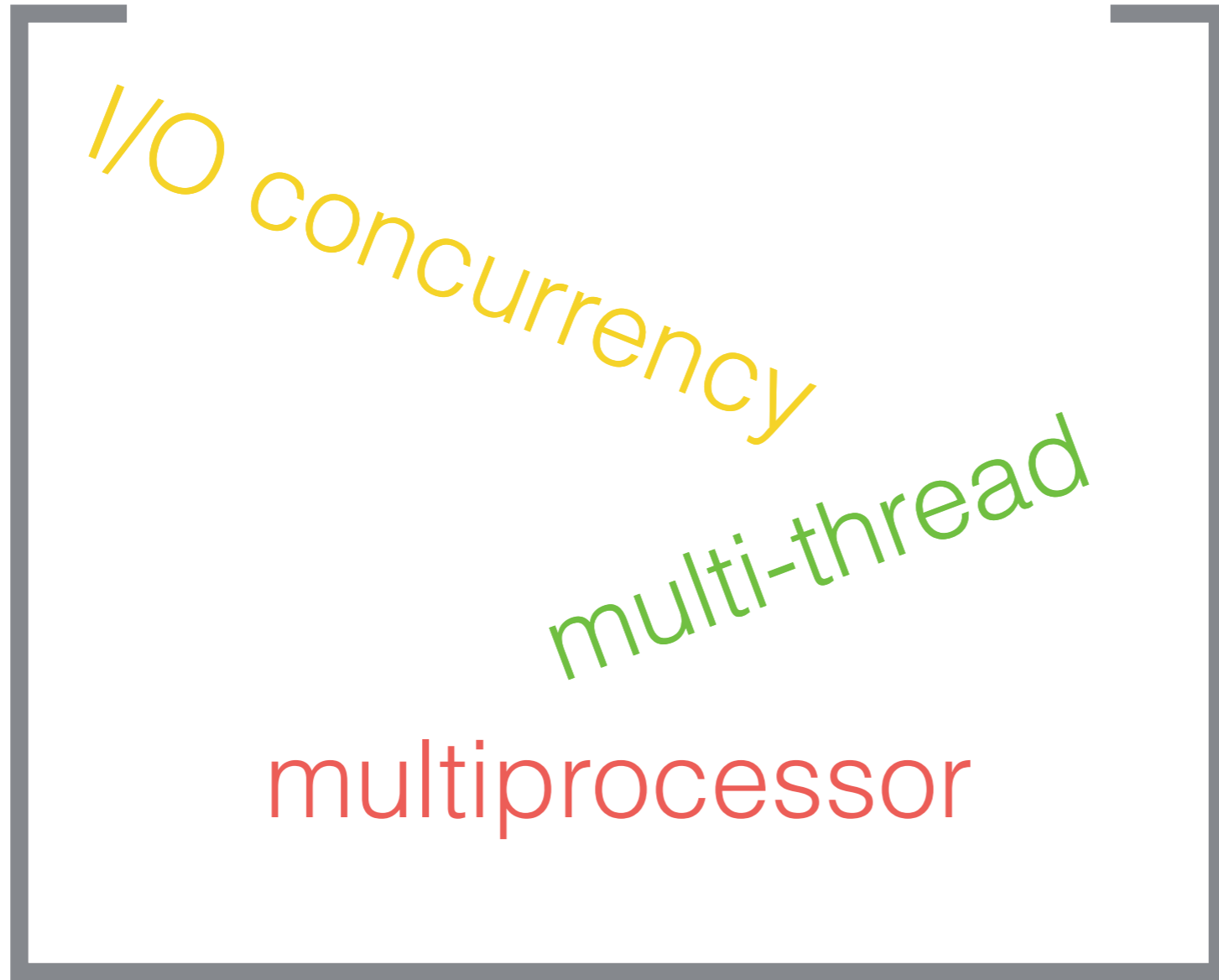
└ global changes

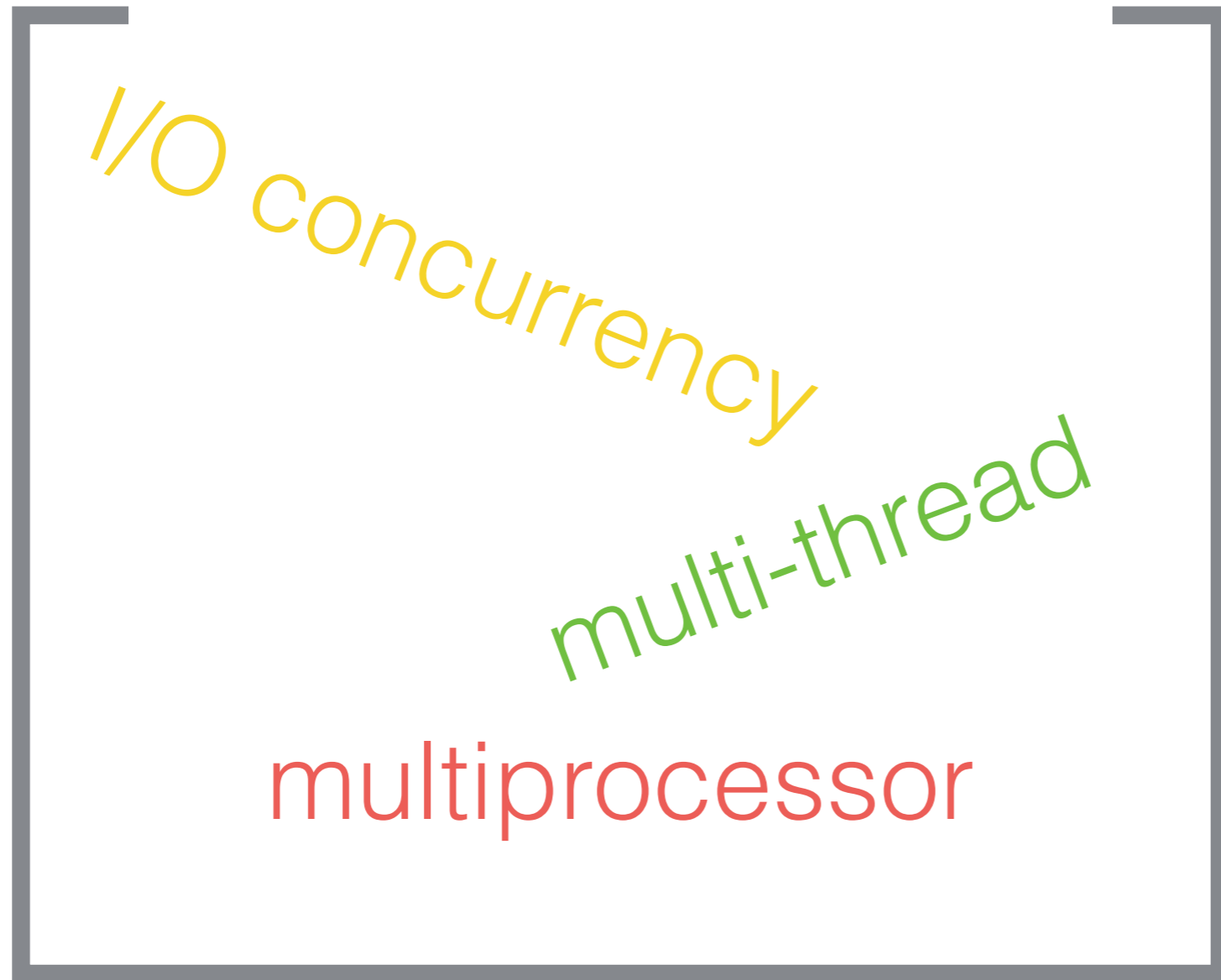
“

[...] **multiprocessor** support, which
may require global changes [...]

”

hard
└─ glo





hard

global change

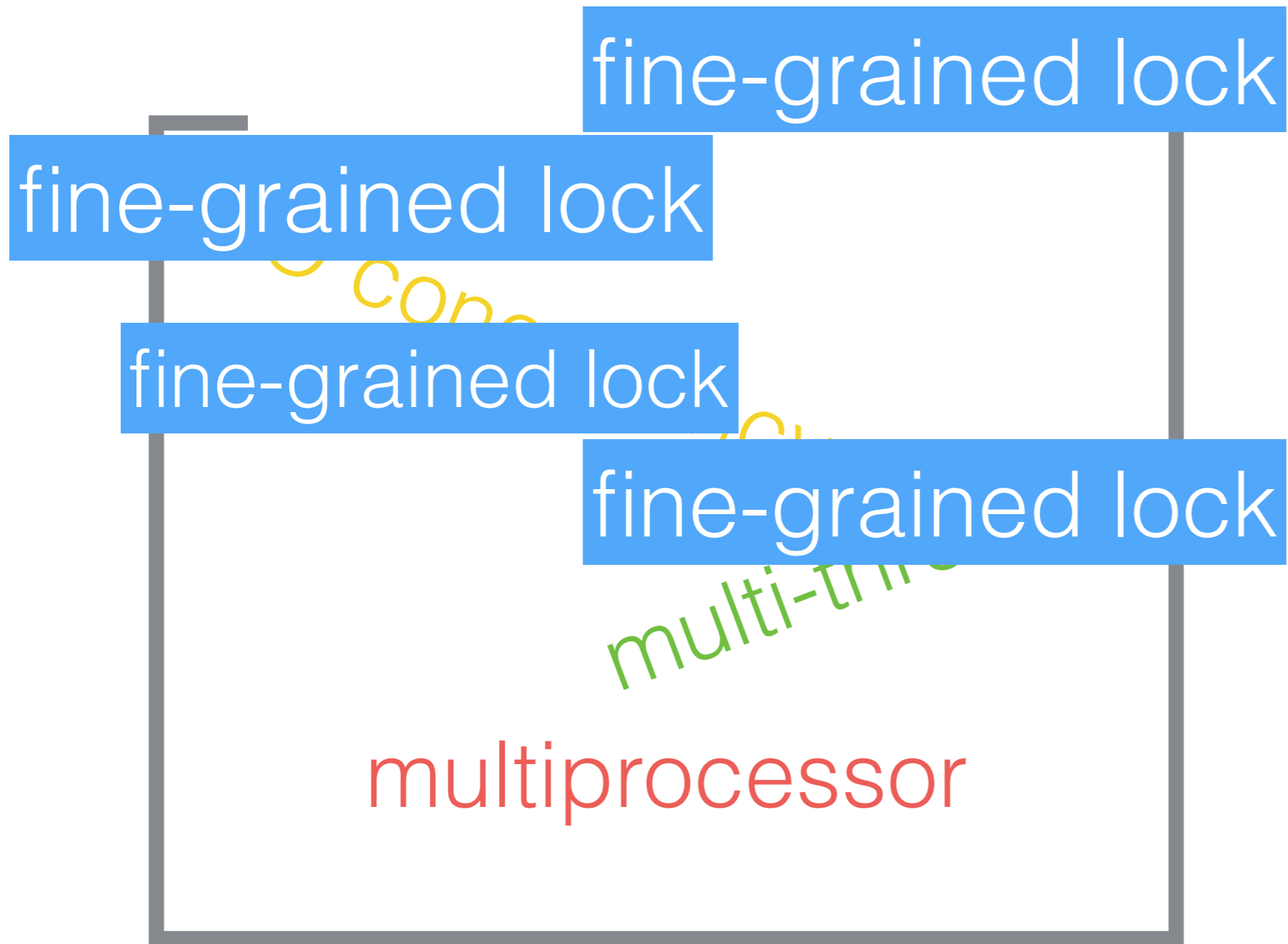
I/O concurrency
multi-thread
multiprocessor

Big Lock

multi-processor

multiprocessor

hard
|
•-glob
|
•-I/O
mu
mu



hard
|
• glo
|
• I/O
mu
mu

S.Peters et al.

[APSys'15]

hard
|
- glo
- 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.

[APSys'15]

“

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 change
- I/O concurrency
- multi-threaded
- multiprocessor
- fine-grained

What to prove?

functional correctness

liveness system calls will eventually return

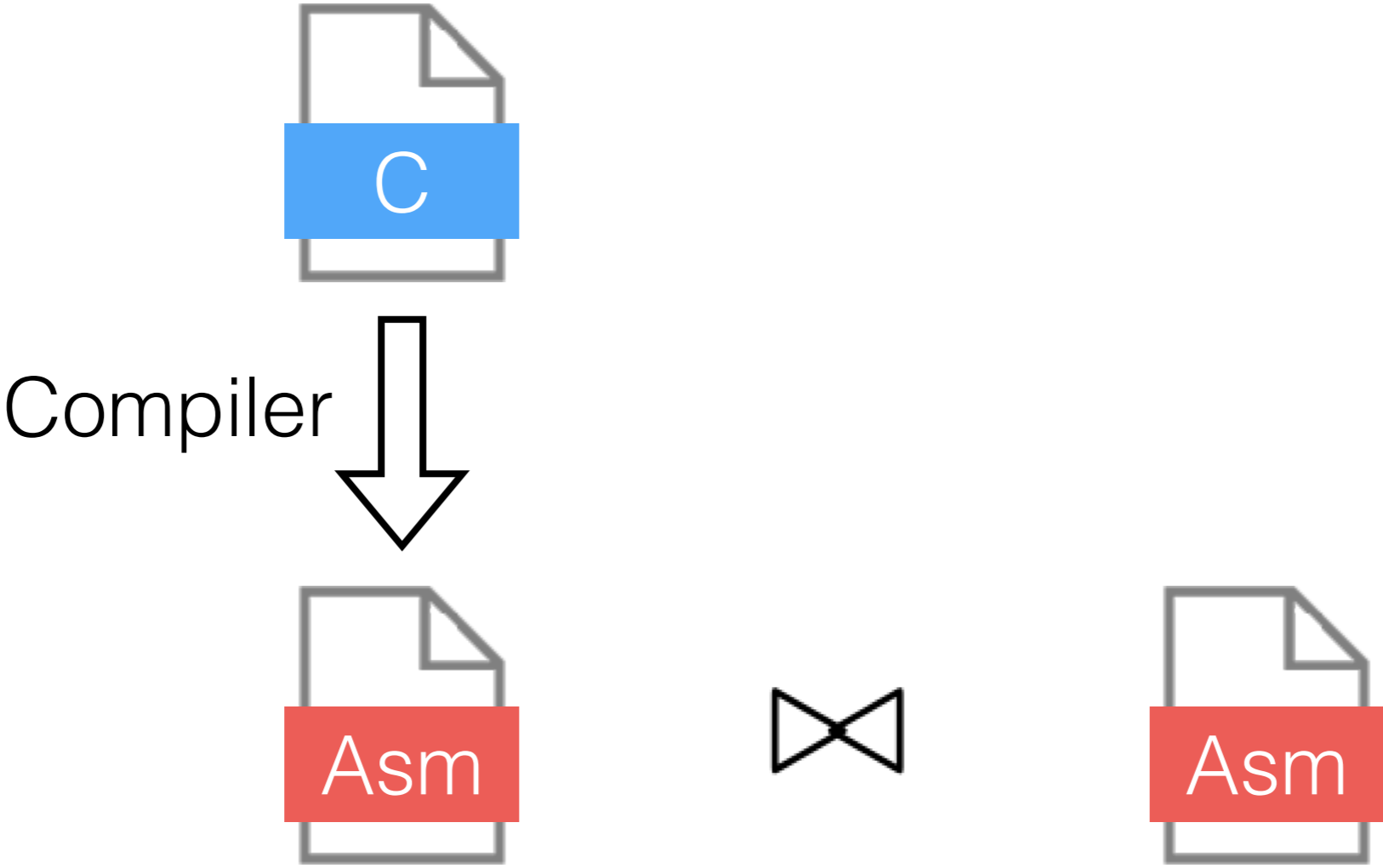


concurrent OS kernel

hard

- global change
- I/O concurrent
multi-thread
multiprocessor
- fine-grained I/O
- liveness

concurrent OS kernel



- hard
- glok
- I/O
- mul
- mu
- fine
- live

hard

- global changes

- I/O concurrency

cost

- asm&C

- compiler

hard

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

CertiKOS

solves all these challenges

- hard
- glok
- I/O
- mul
- mu
- fine
- live
- asr
- con
- cos

contributions

Certikos

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

hard

glok

I/O

mul

mu

fine

live

asr

cor

cos

CertIKOS

contributions

- mC2
- fine-grained lock

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

hard

• glob

• I/O

mul

mu

• live

• asr

• cor

• COS

CertIKOS

contributions

- mC2
- fine-grained lock

both functional correctness
and liveness

hard

• glob

• I/O

mul

mu

• live

• asr

• cor

• COS

Certikos

contributions

- mC2
- fine-grained lock
- liveness

both functional correctness
and liveness

hard

• glob

• I/O

mul

mu

• asr

• cor

• COS

Certikos

contributions

- mC2
- fine-grained lock
- liveness

certified concurrent layers

hard

• glob

• I/O

mul

mu

• asr

• cor

• COS

Certikos

contributions

- mC2
- fine-grained lock
- liveness

reuses sequential verification techniques.

certified concurrent layers

hard

• glob

• I/O

mul
mu

• asr

• cor

• COS

Certikos

contributions

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

reuses sequential verification techniques.

certified concurrent layers

hard

I/O
mul
mu

asr

cor

cos

Certikos

contributions

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

handles all 3 kinds of
concurrency

certified concurrent layers

hard

- I/O
- mul
- mu

• asr

• cor

• COS

Certikos

contributions

- mC2
- fine-grained lock
- liveness
- reuse of techs
- I/O concurrency
 - multi-thread
 - multiprocessor

handles all 3 kinds of
concurrency

certified concurrent layers

hard

asr

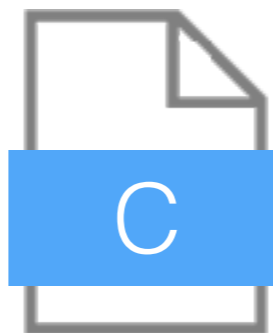
con

cos

Certikos

contributions

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



6100 LOC



400 LOC

hard

asm

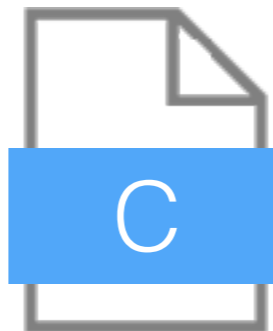
cor

COS

Certikos

contributions

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



6100 LOC



400 LOC

hard

cor

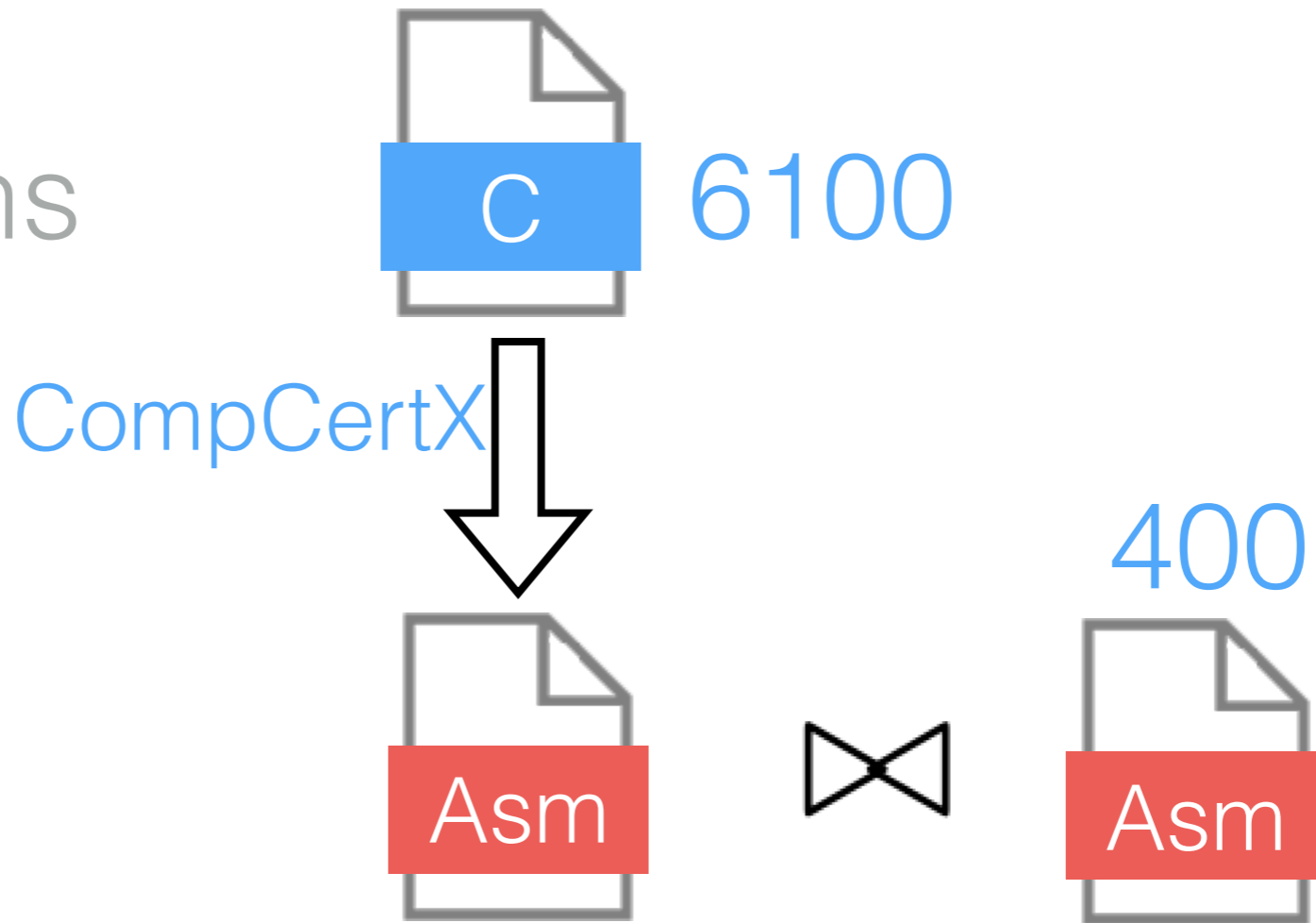
COS

Certikos

contributions

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

hard



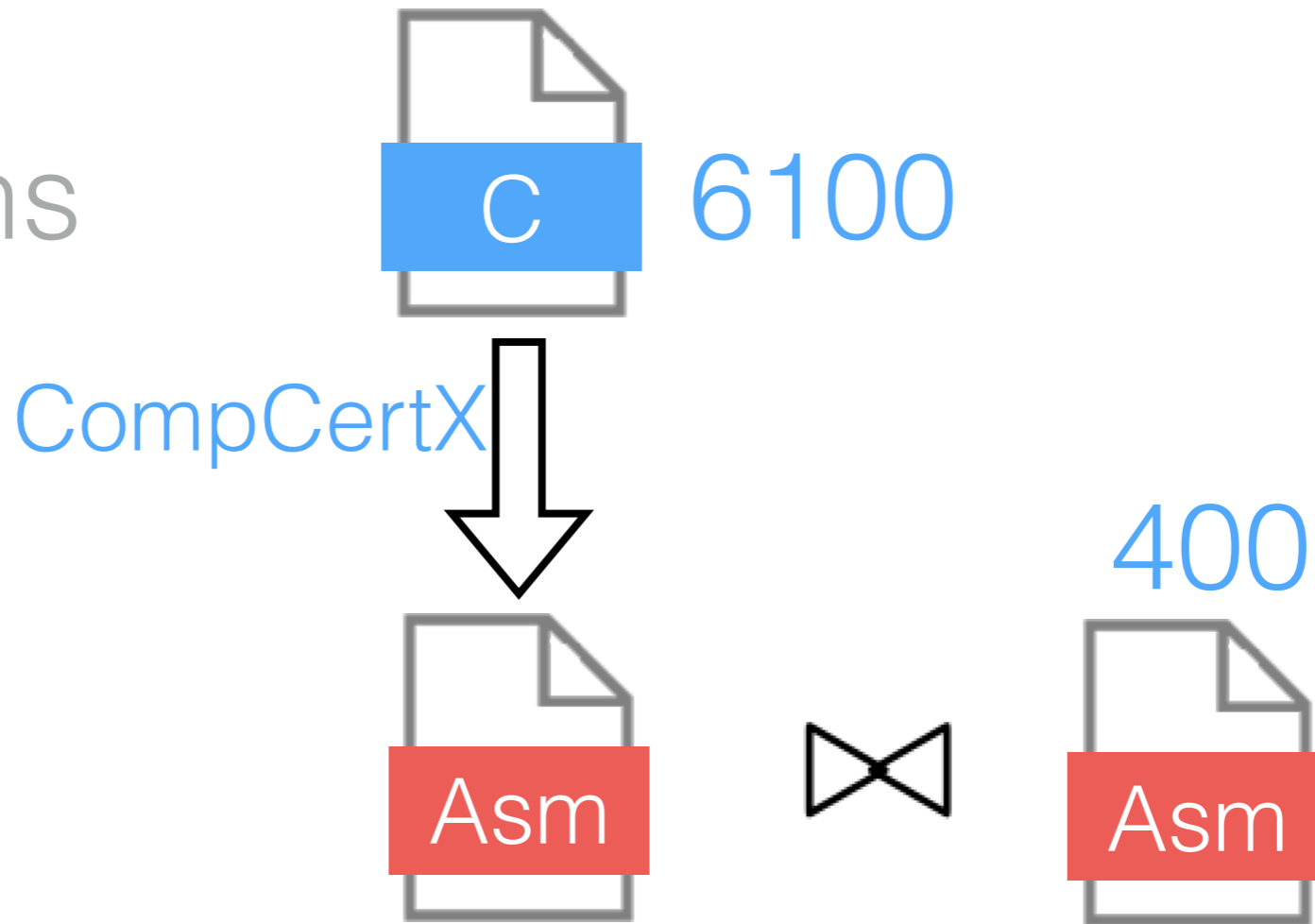
cor

cos

CertIKOS

contributions

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



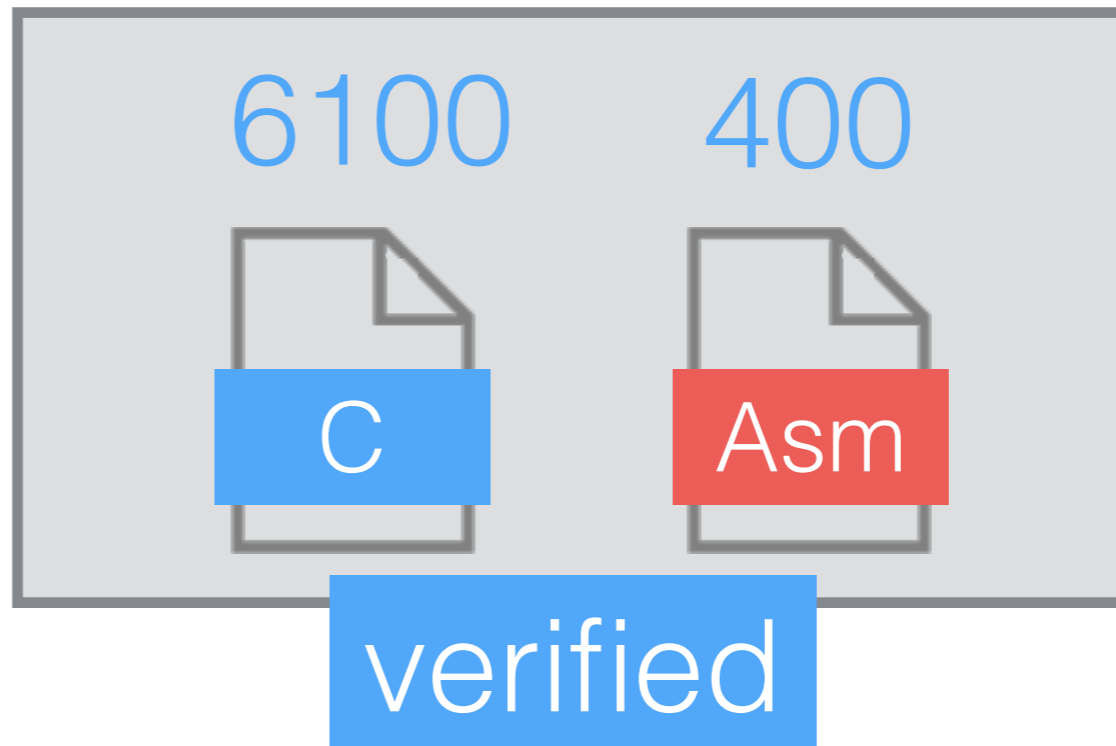
hard

COS

Certikos

contributions

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



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

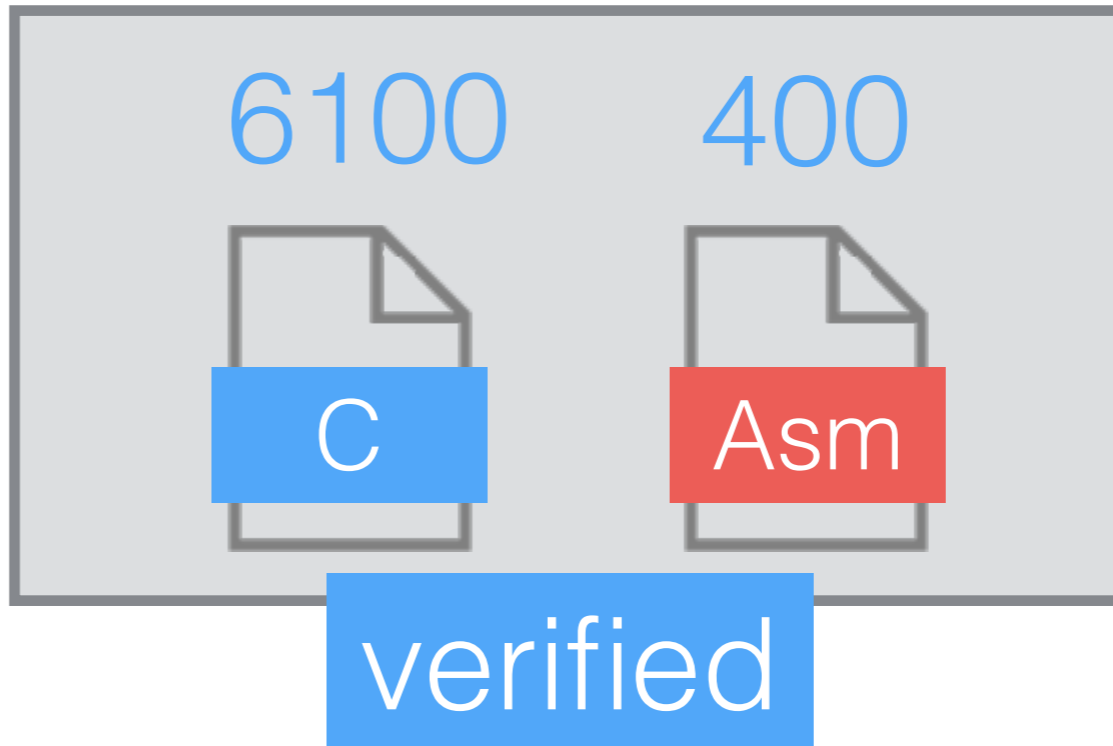
hard

cos

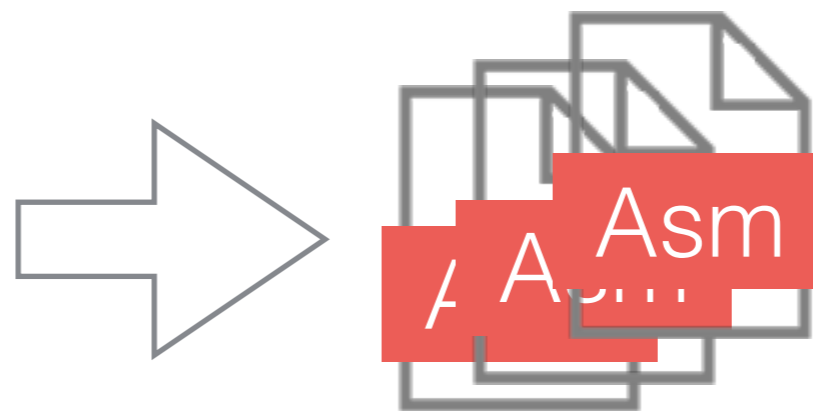
CertiKOS

contributions

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



machine-checkable
proof



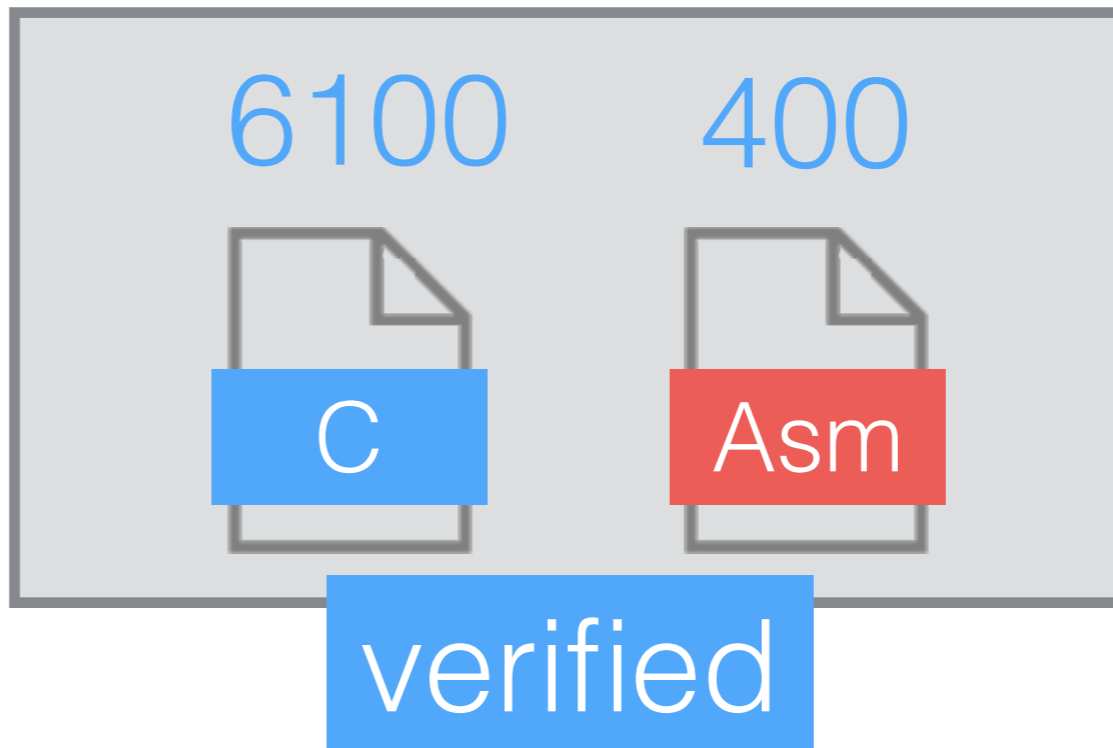
hard

COS

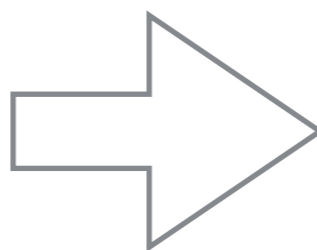
CertiKOS

contributions

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



le



executable

hard

cos

Certikos

contributions

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

hard

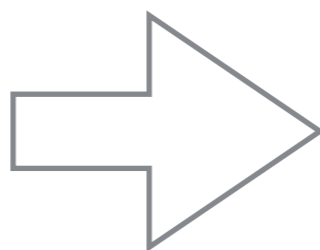
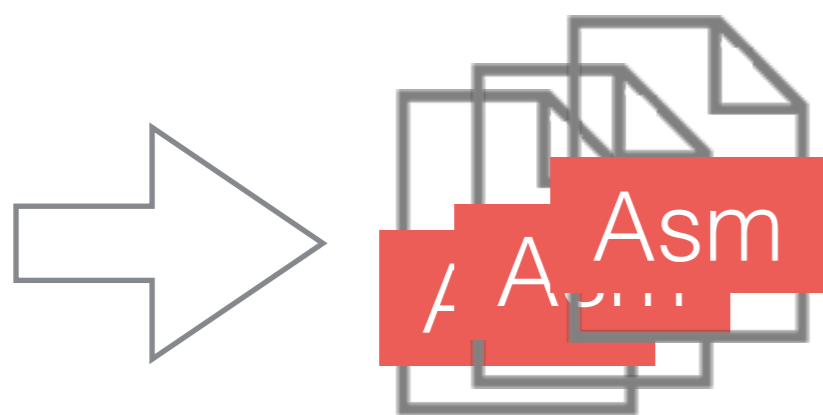


~~verified~~

certified

executable

le



cos

CertiKOS

contributions

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

mCertikOS

[POPL'15]

1 py

+ extensions 0.5 py

+ device 0.5 py

[PLDI'16]

+ concurrency 2 py

hard

cos

Certikos

contributions

- 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

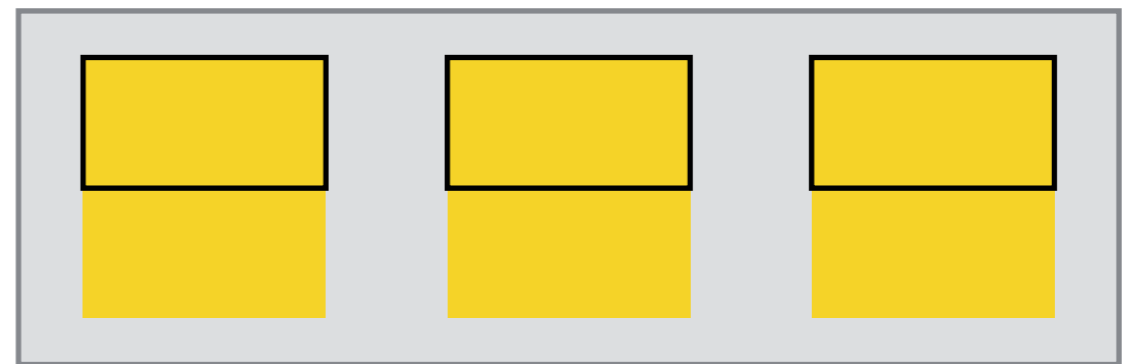
push/pull model

multicore machine **lifting**

certified sequential layers

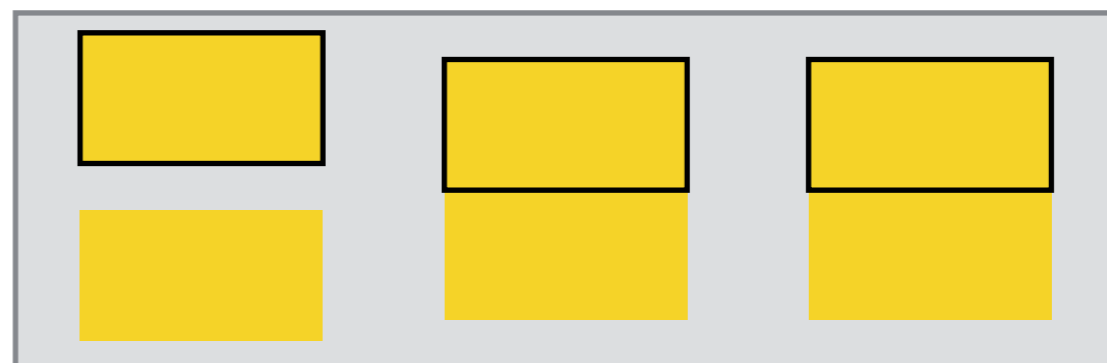


certified objects



specification of modules to trust

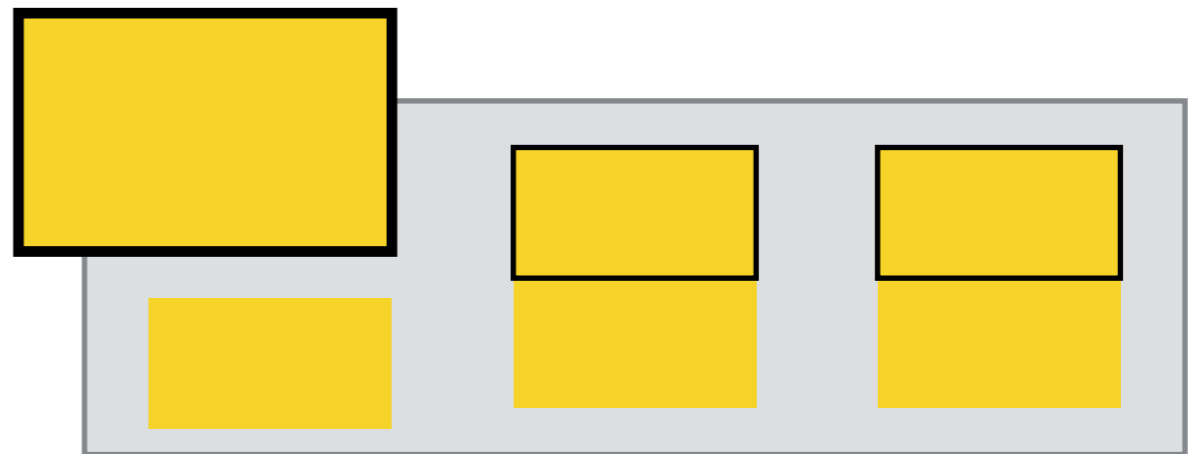
certified sequential layers



certified sequential layers



abs-state



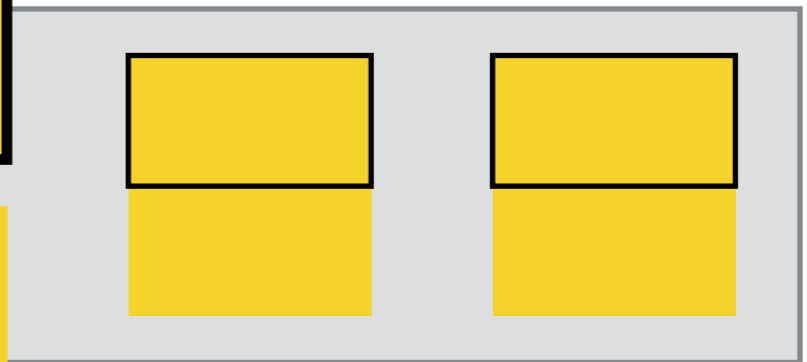
certified sequential layers



abs-state

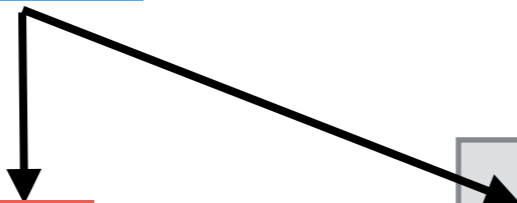


primitives

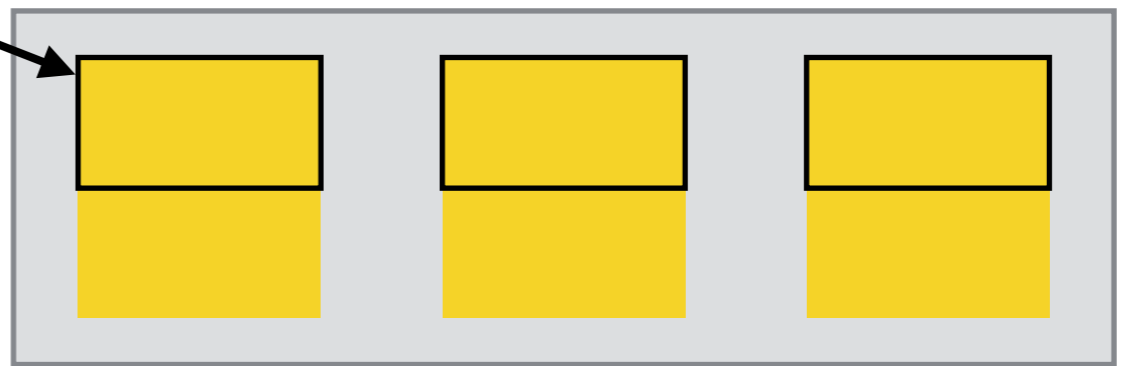


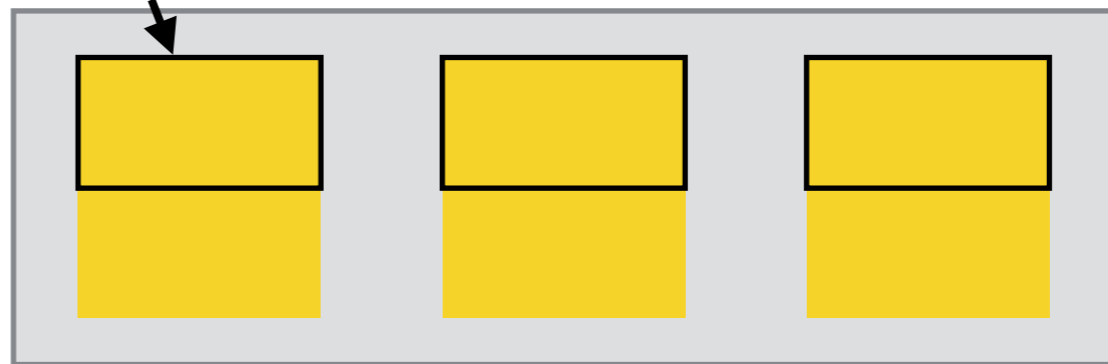


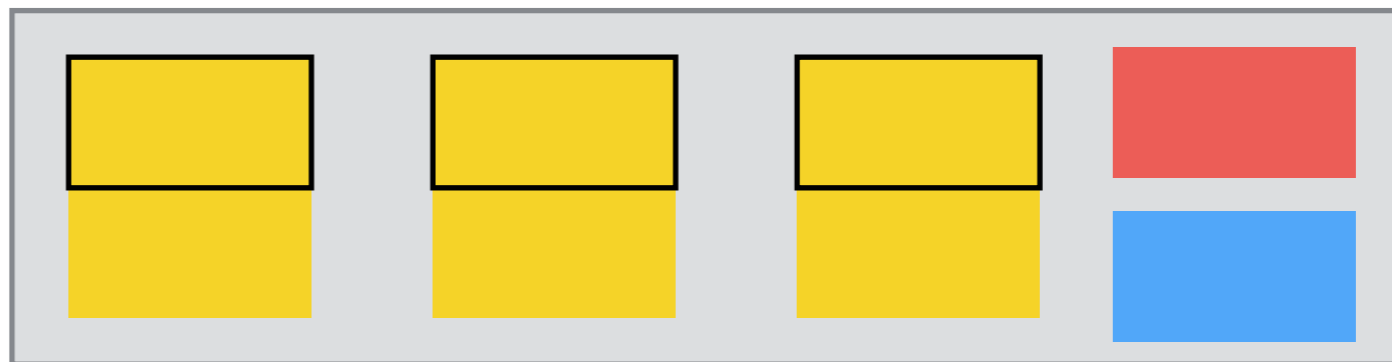
code



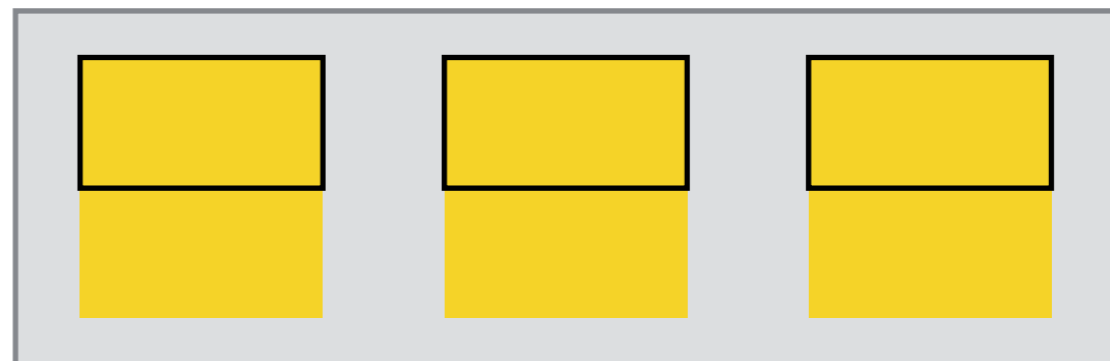
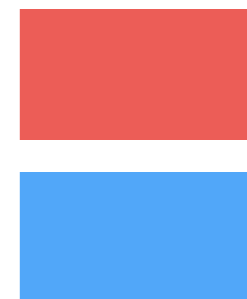
memory





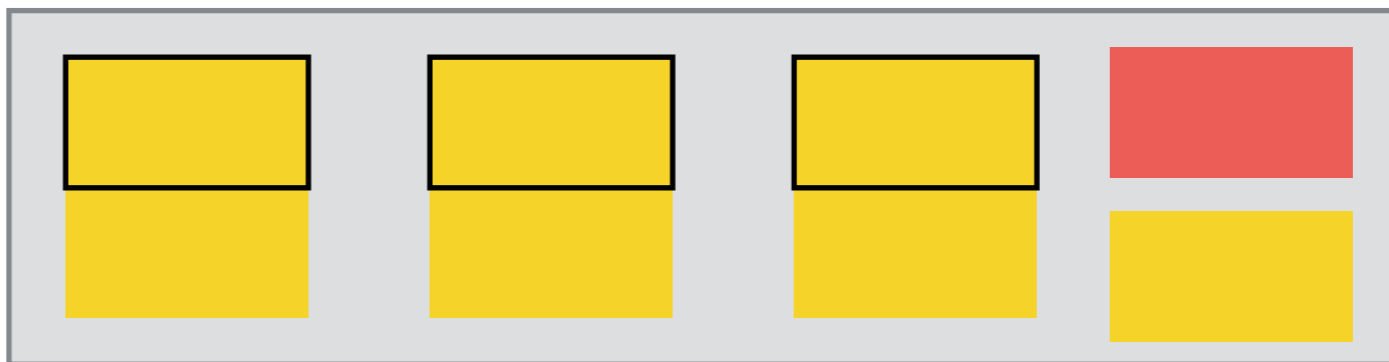


implementation

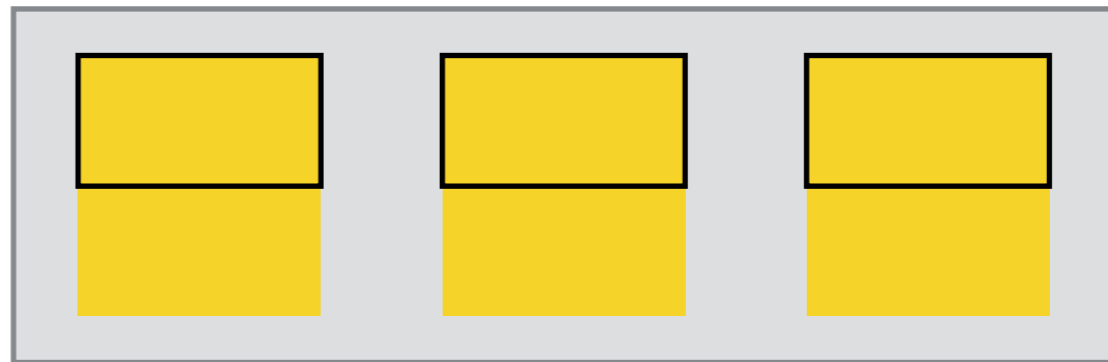
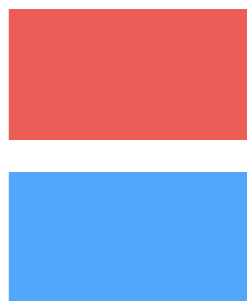




specification

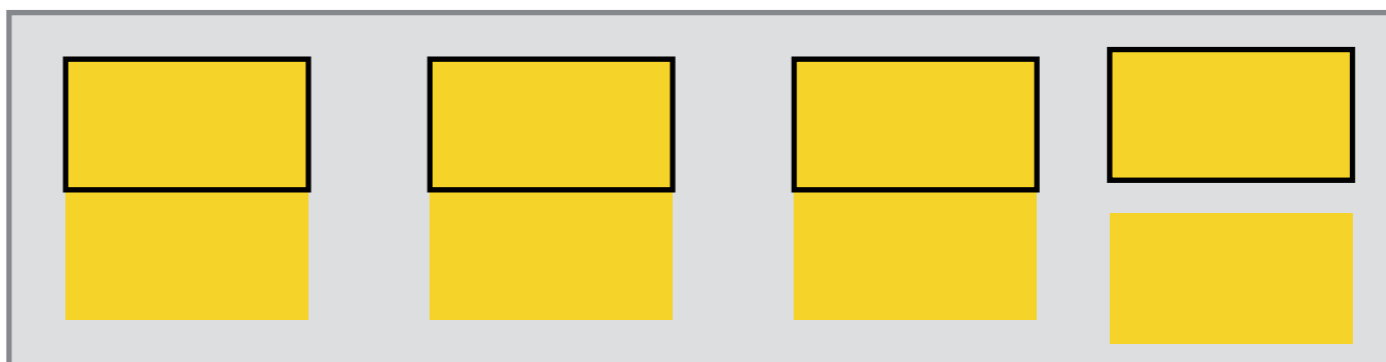


implementation

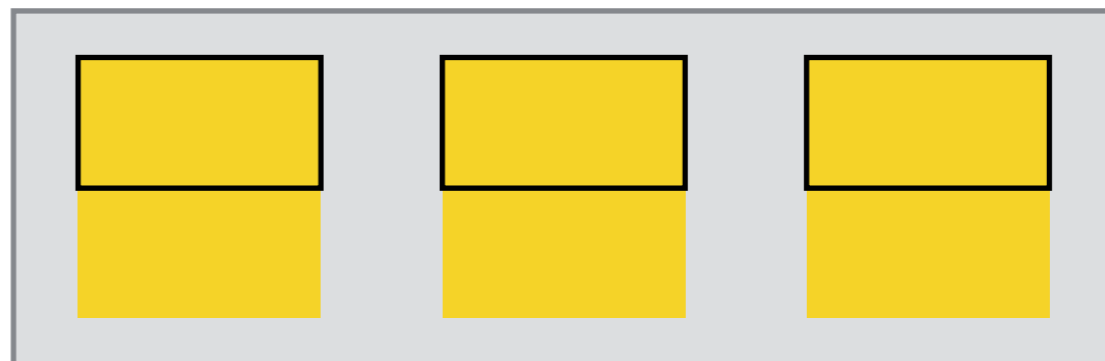
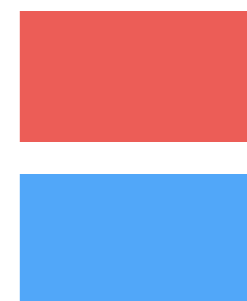




specification



implementation

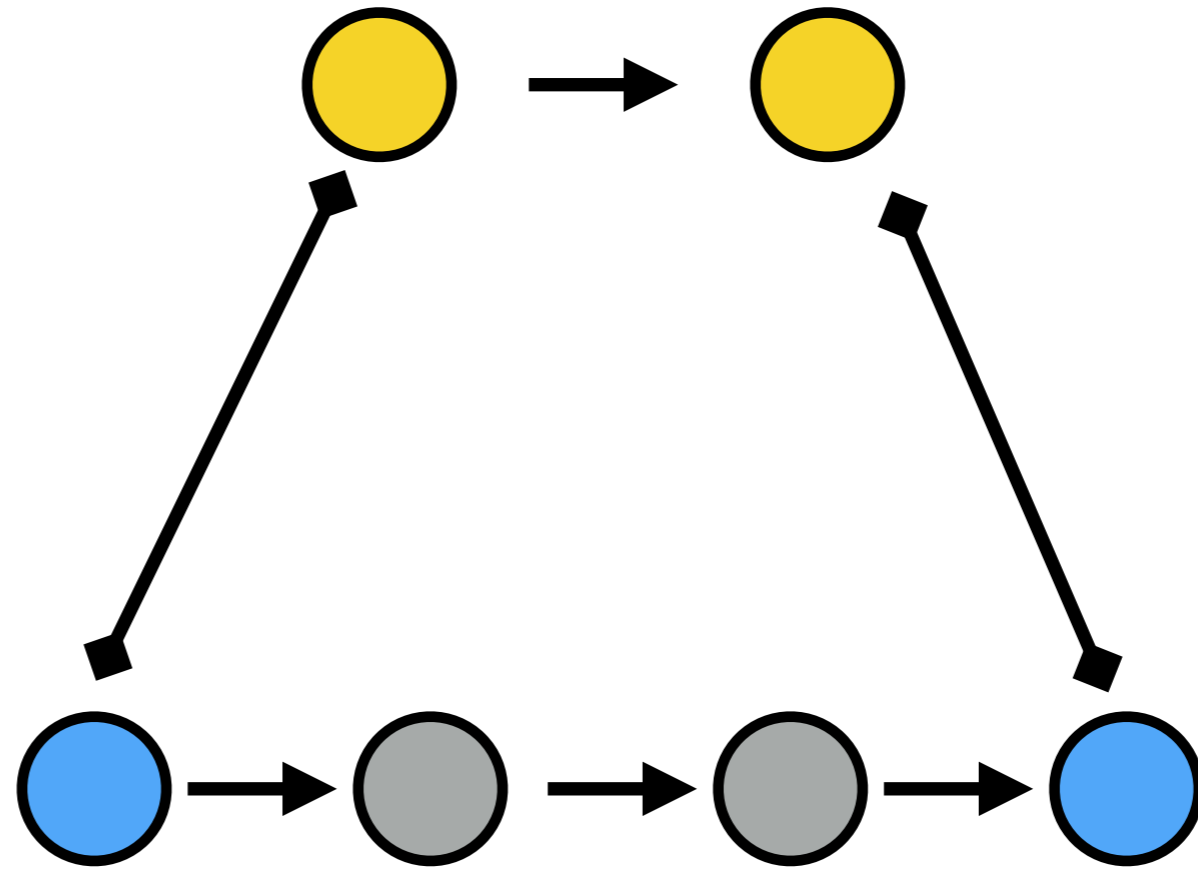


simulation proof

specification



implementation



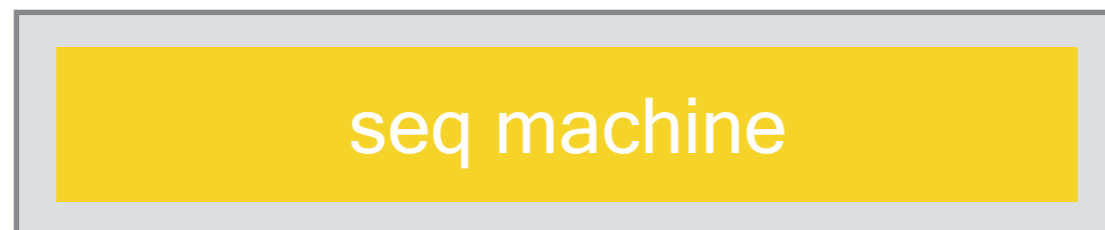


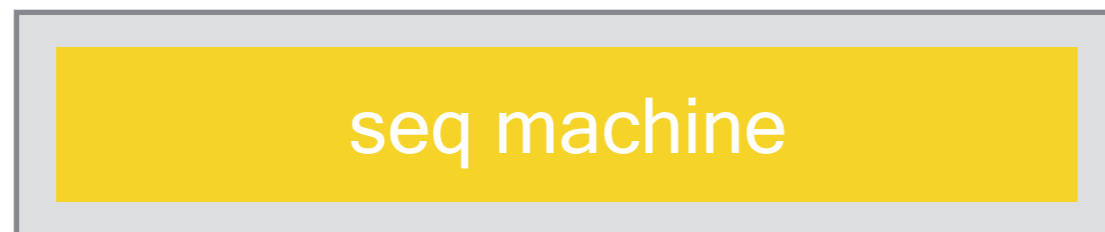
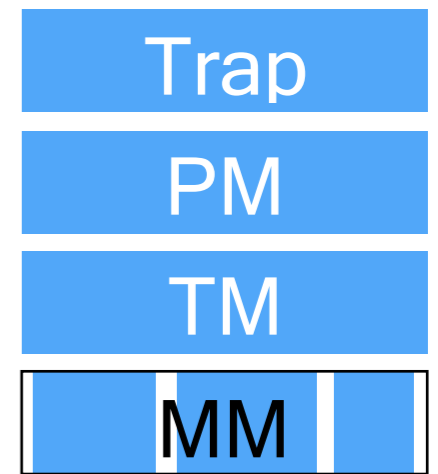
verify a **sequential** kernel

[POPL'15]

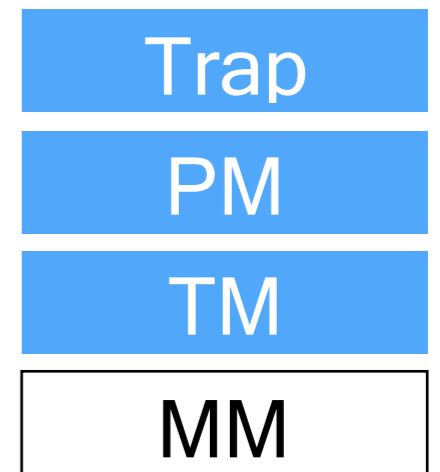
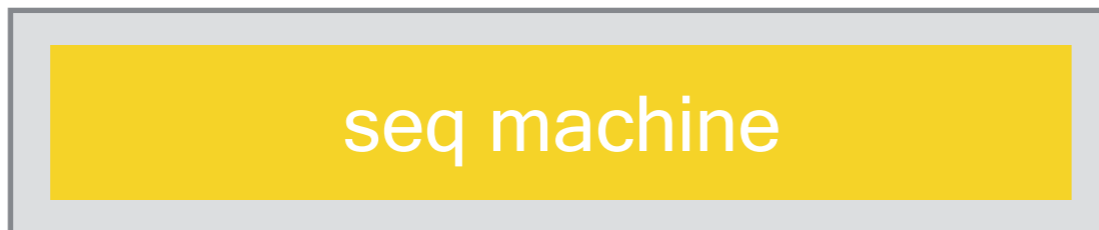
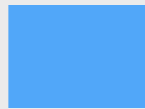
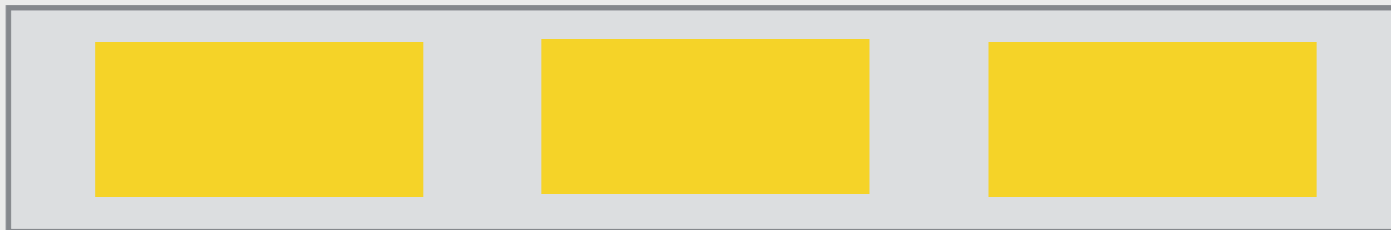


kernel





memory management





trap

Trap

proc

PM

thread

TM

mem

seq machine

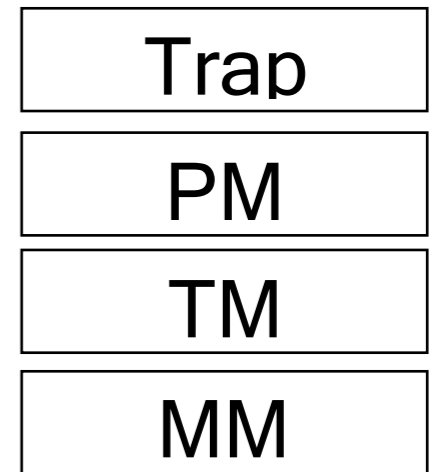
Trap

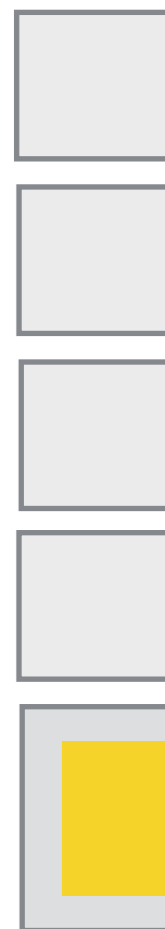
PM

TM

MM

verified sequential kernel





TSysCallLayer

(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/rcv/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

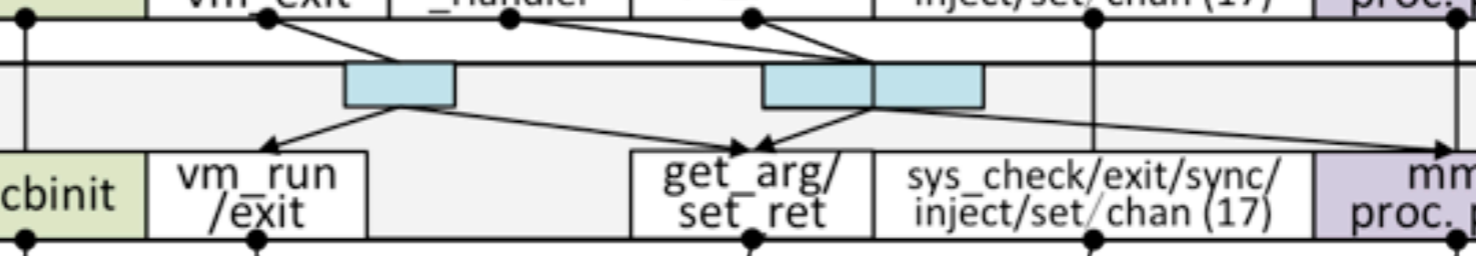
(mm/proc/virt.abs)

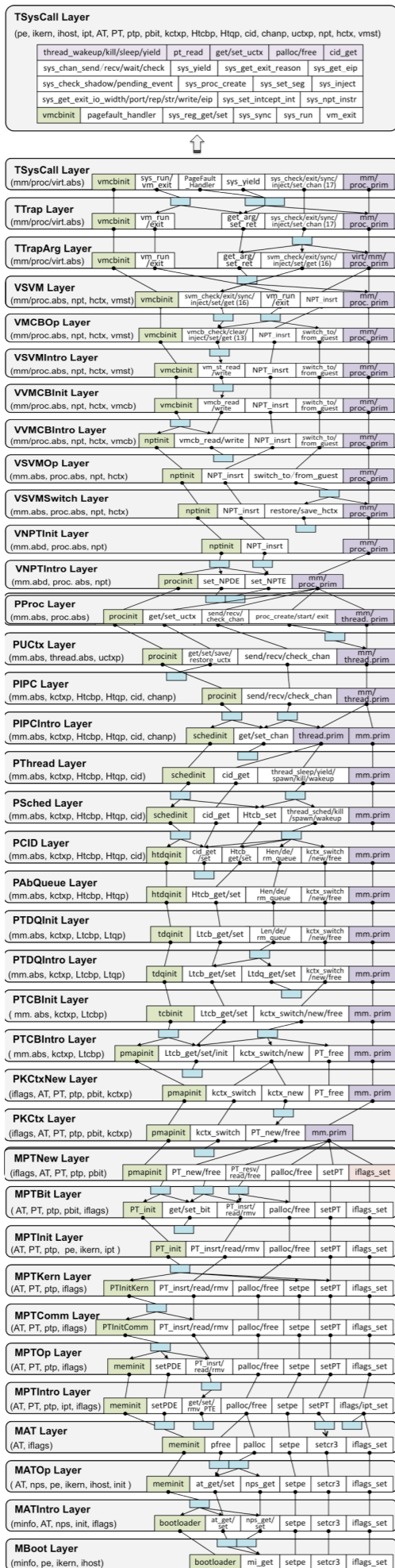
vmcbinit	sys_run/vm_exit	PageFault_Handler	sys_yield	sys_check/exit/sync/inject/set/chan(17)	mm/proc_prim
----------	-----------------	-------------------	-----------	---	--------------

TTrap Layer

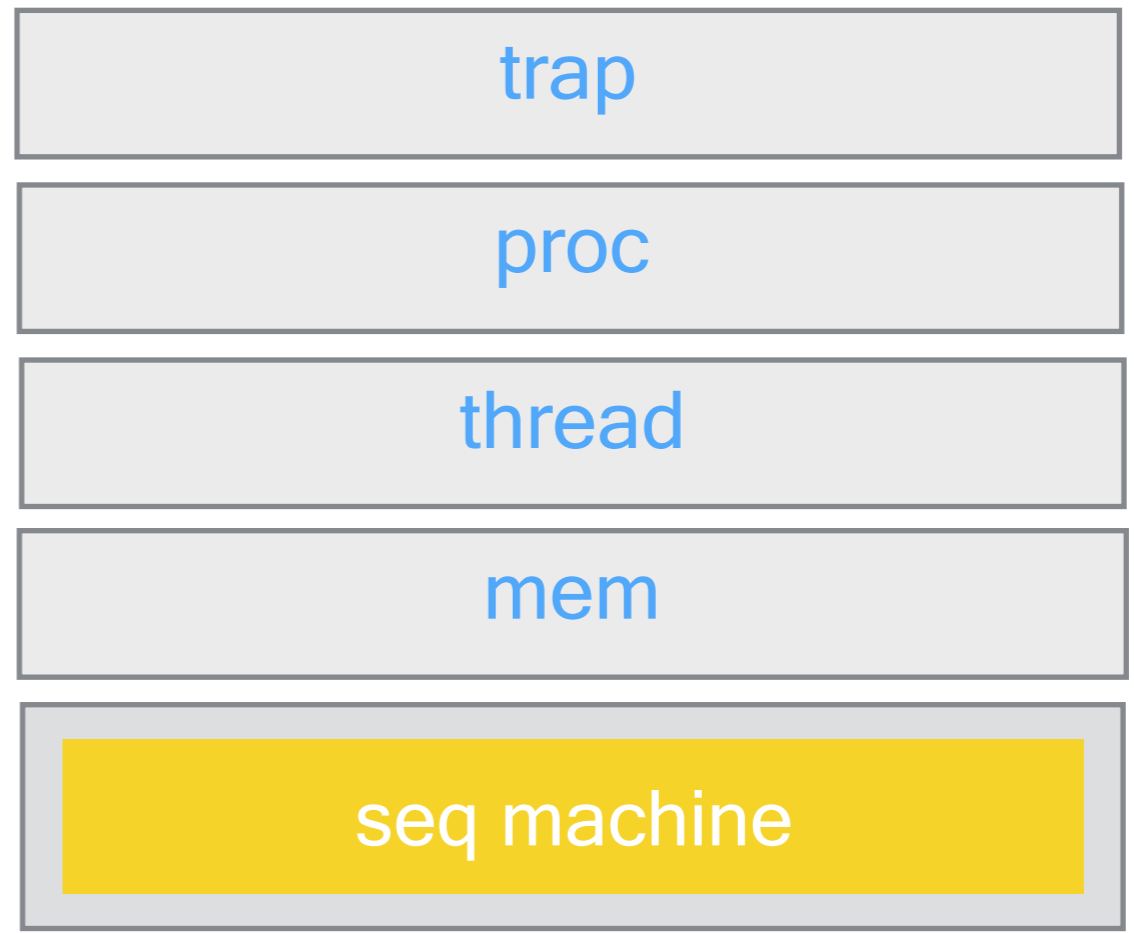
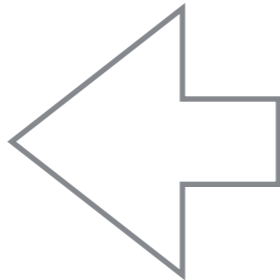
(mm/proc/virt.abs)

vmcbinit	vm_run/exit	get_arg/set_ret	sys_check/exit/sync/inject/set/chan(17)	mm/proc_prim
----------	-------------	-----------------	---	--------------



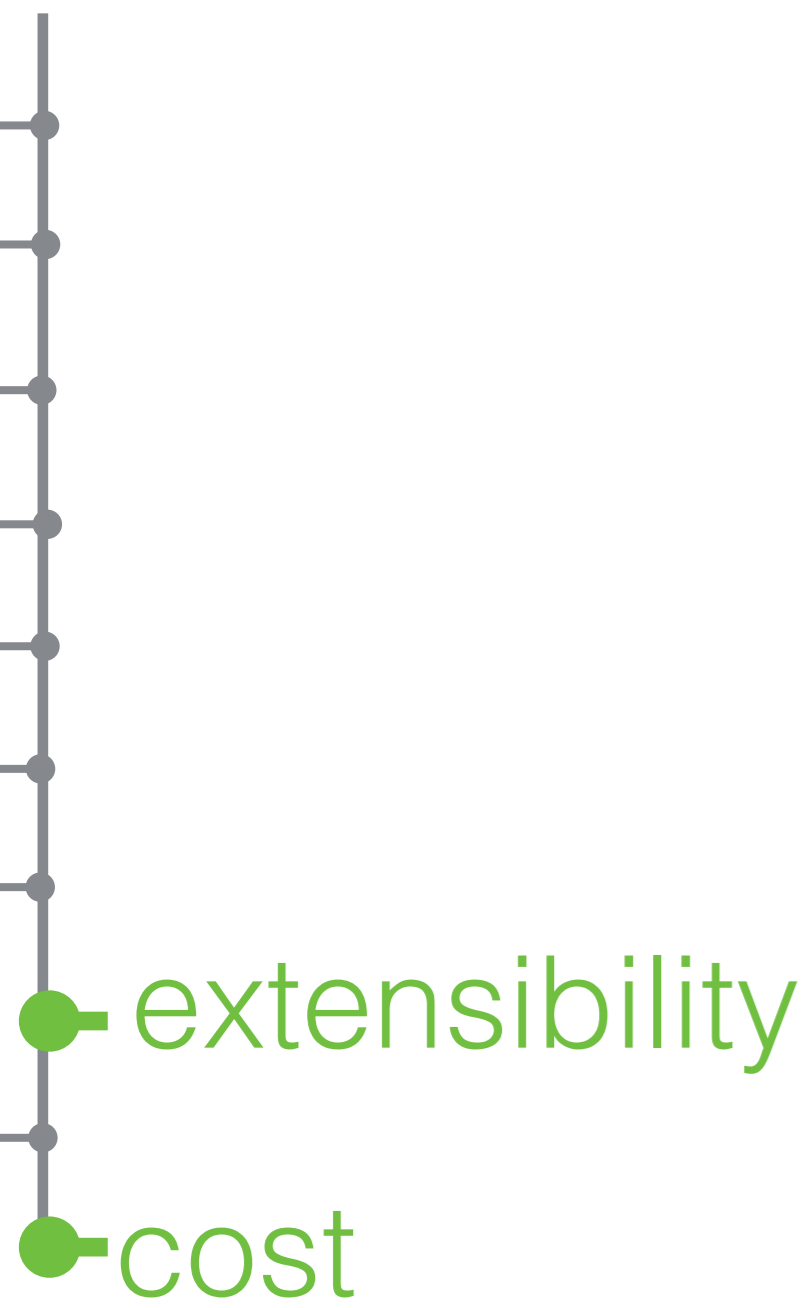


1 person year
(tools construct excluded)

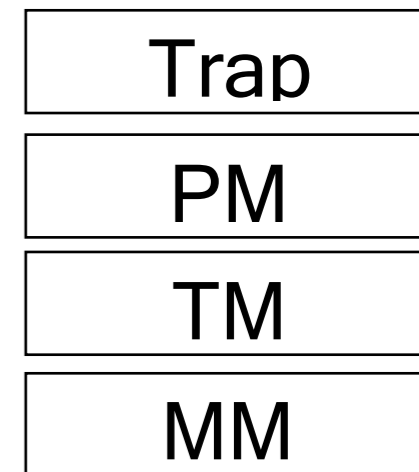


cost

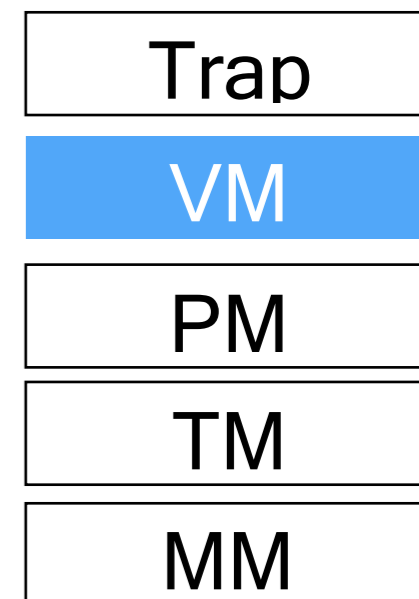
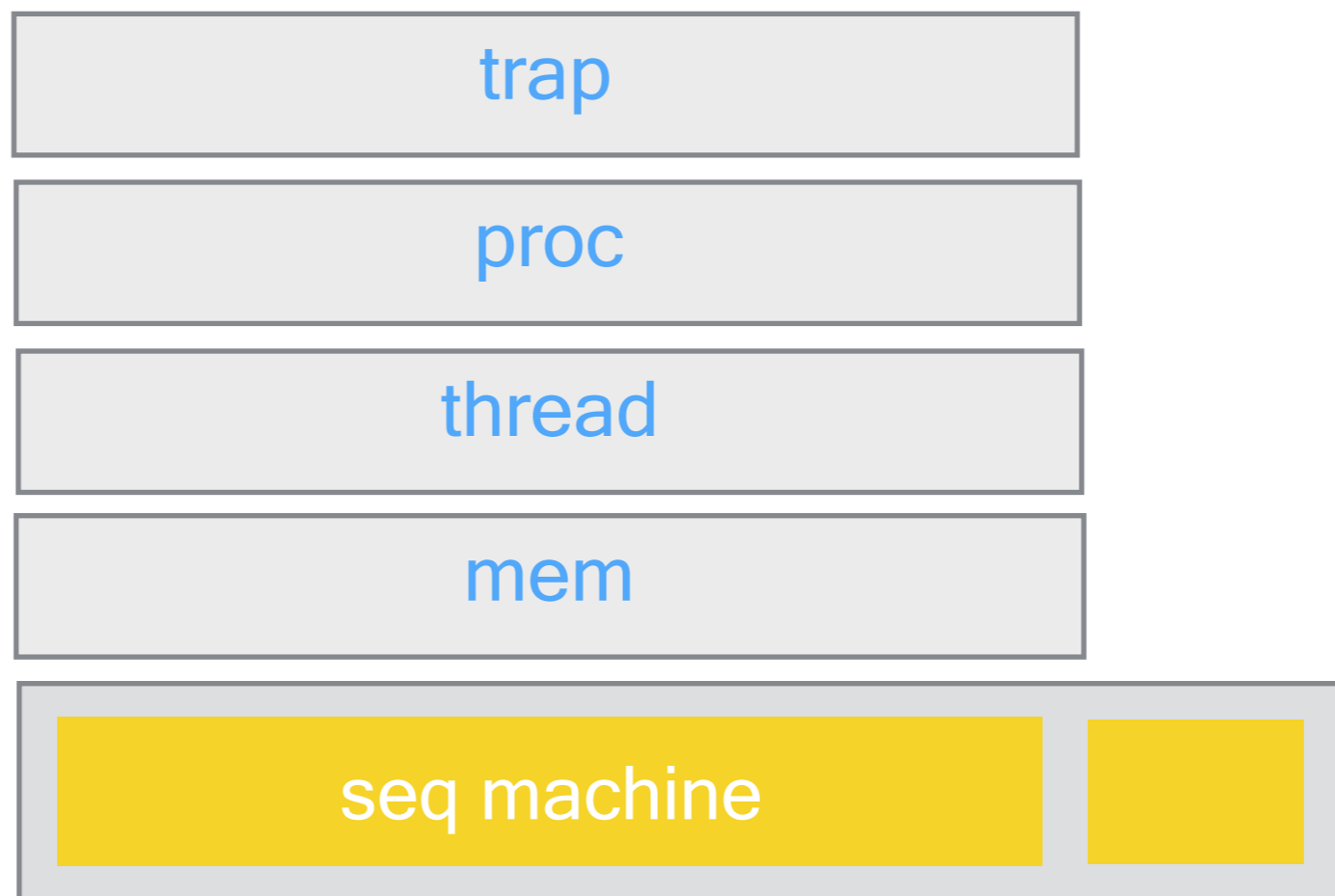
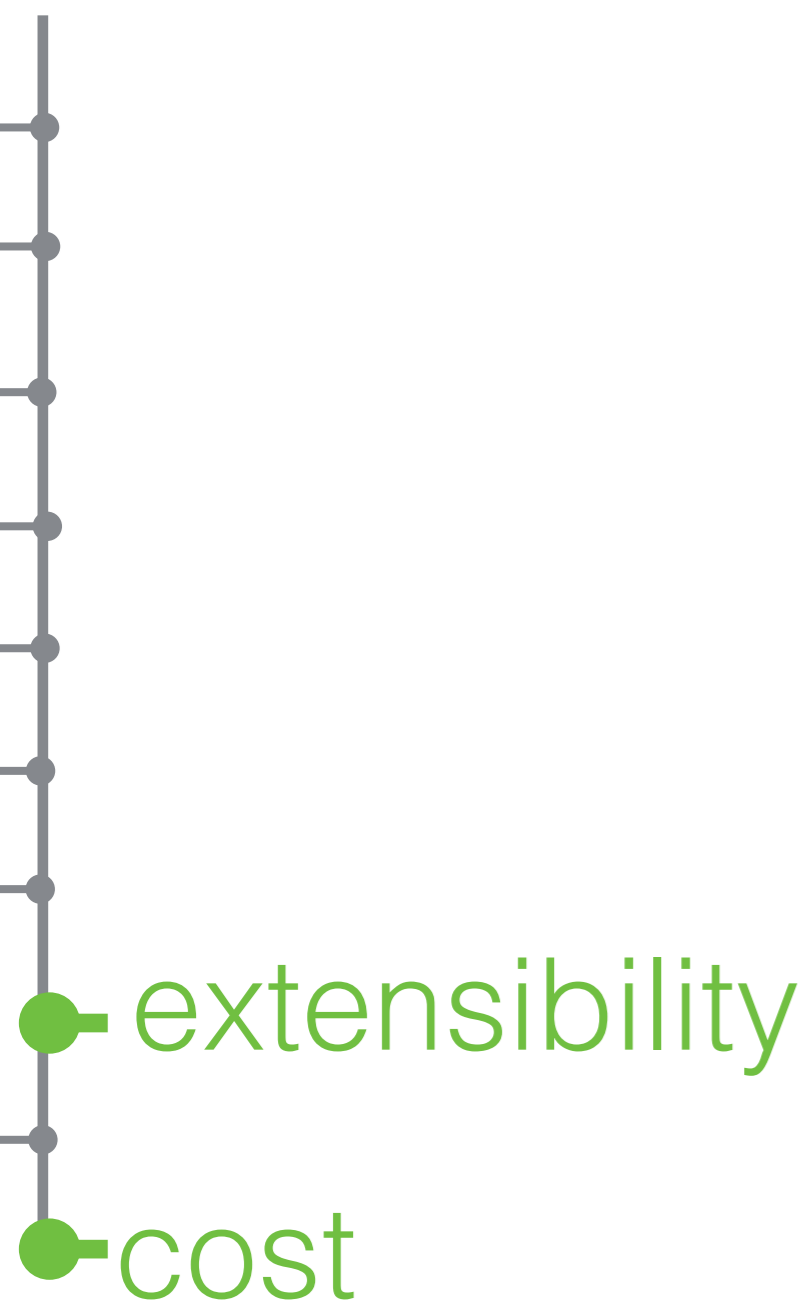
contributions



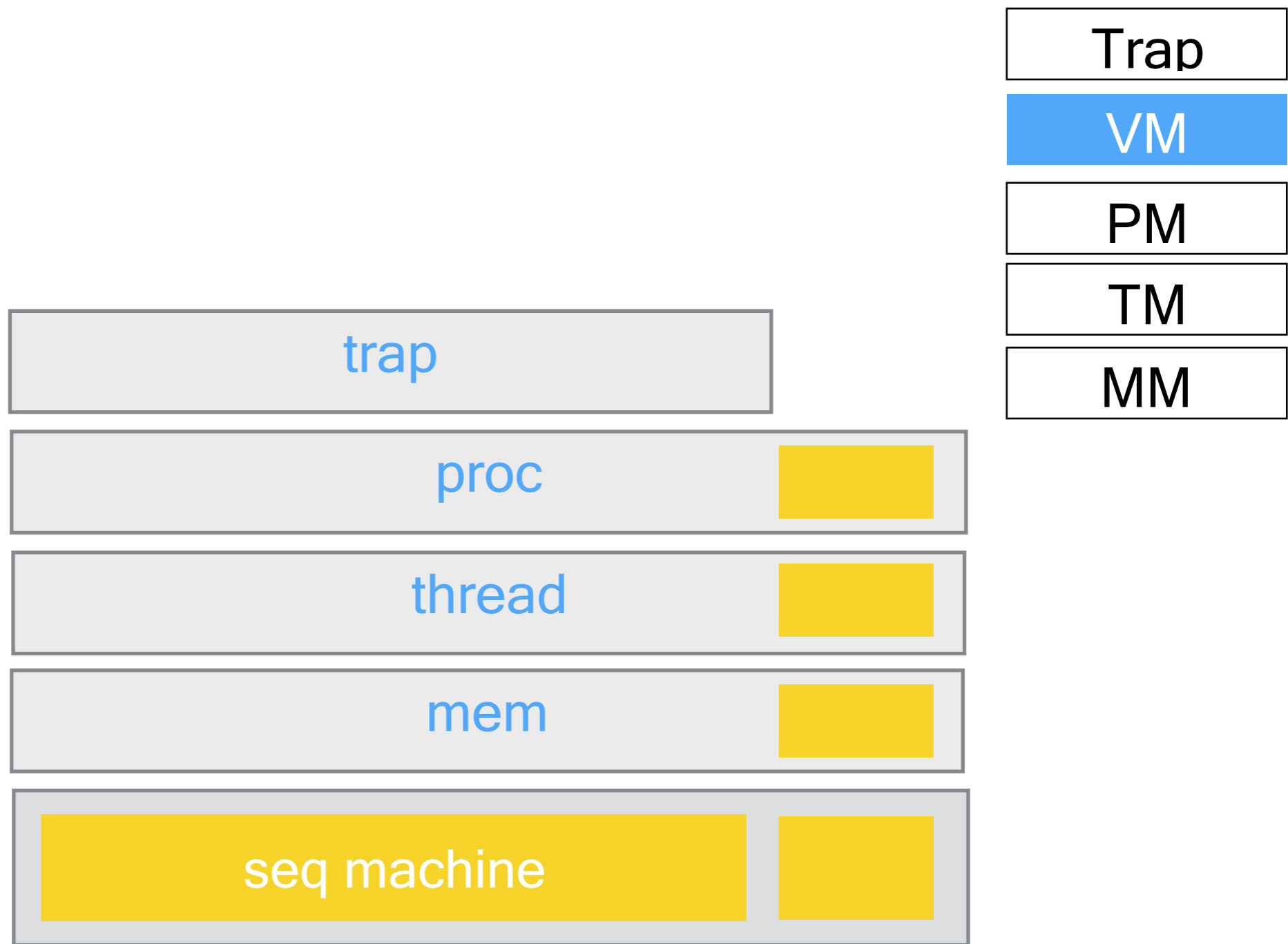
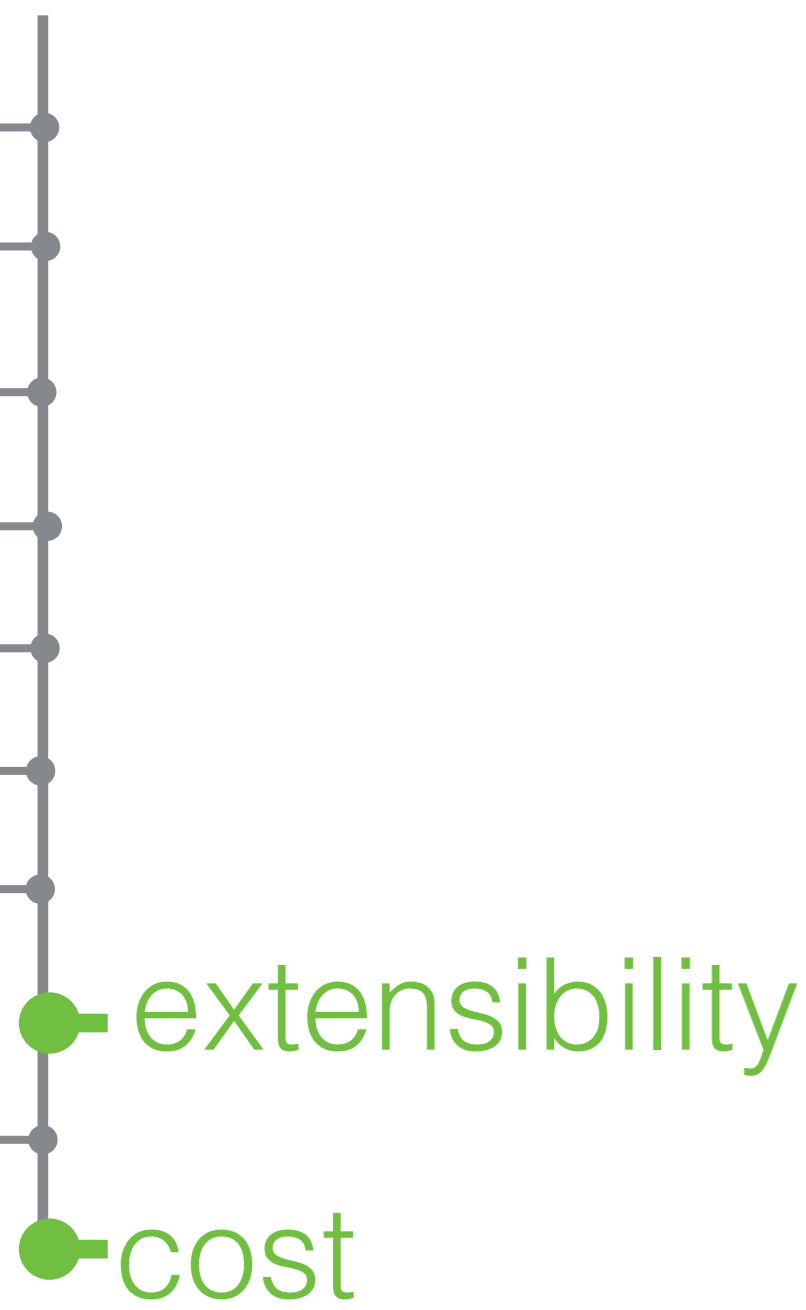
VM



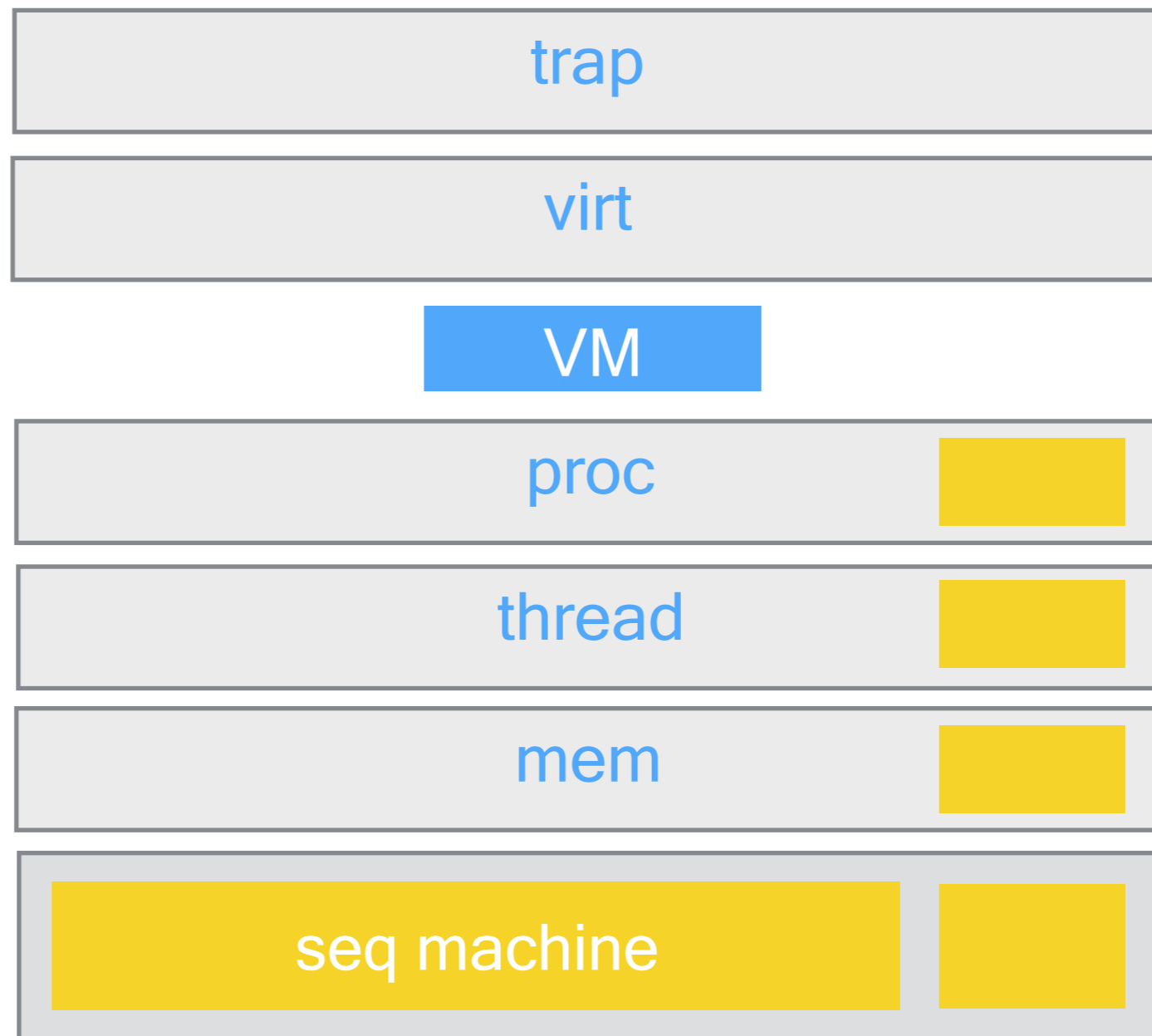
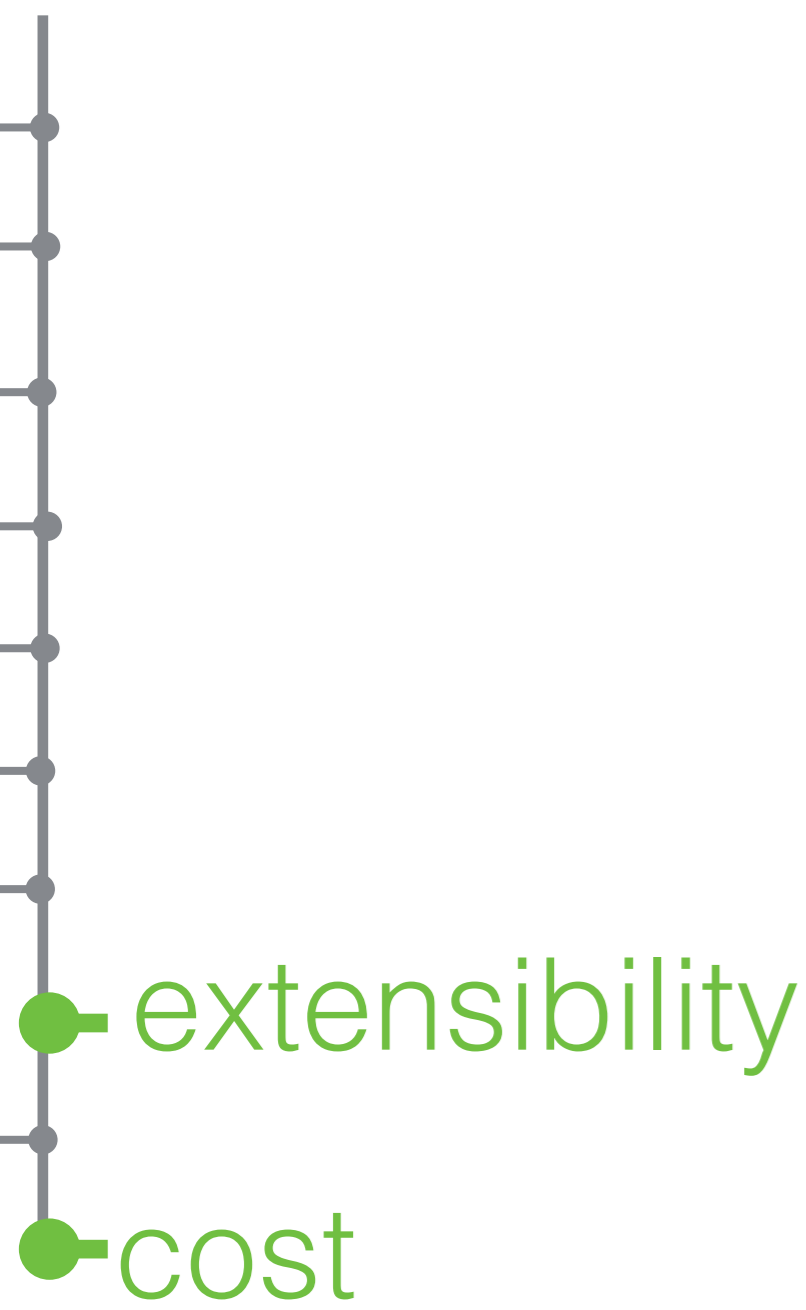
contributions



contributions

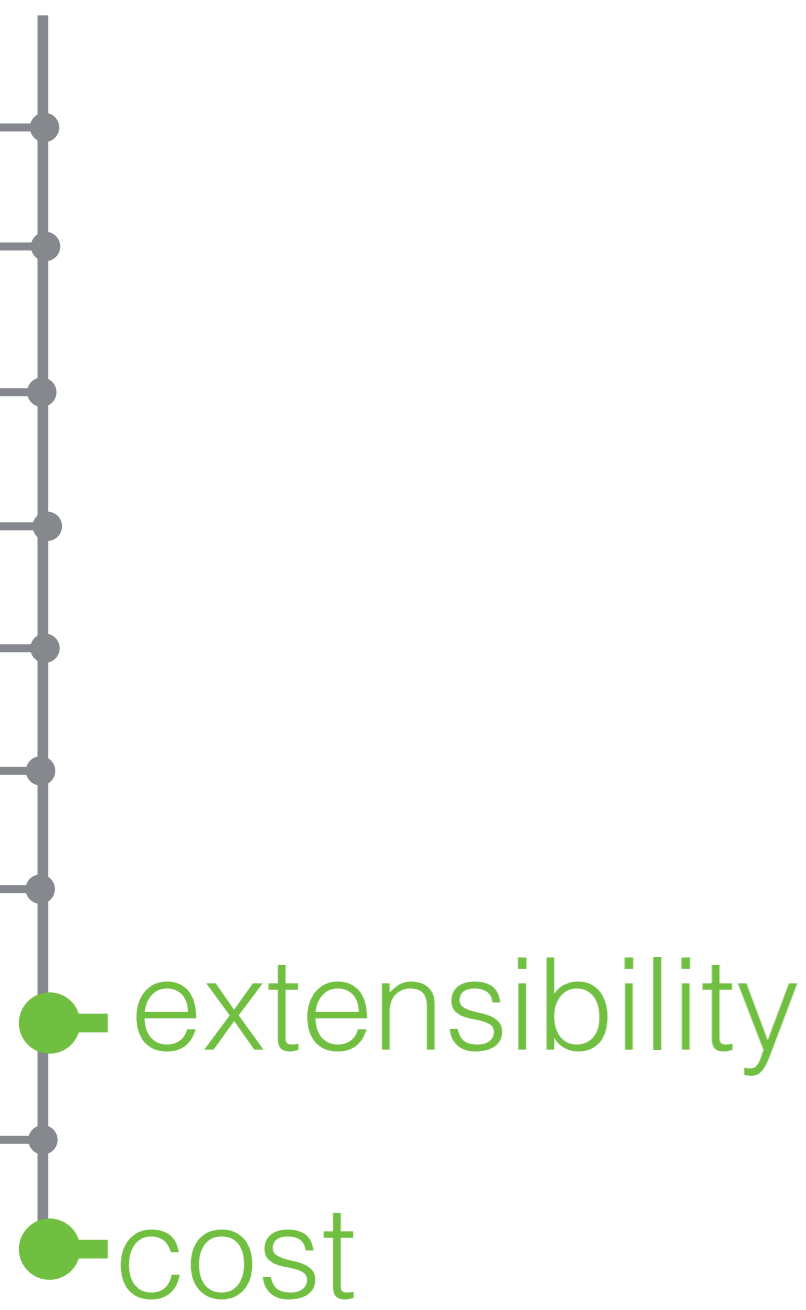


contributions

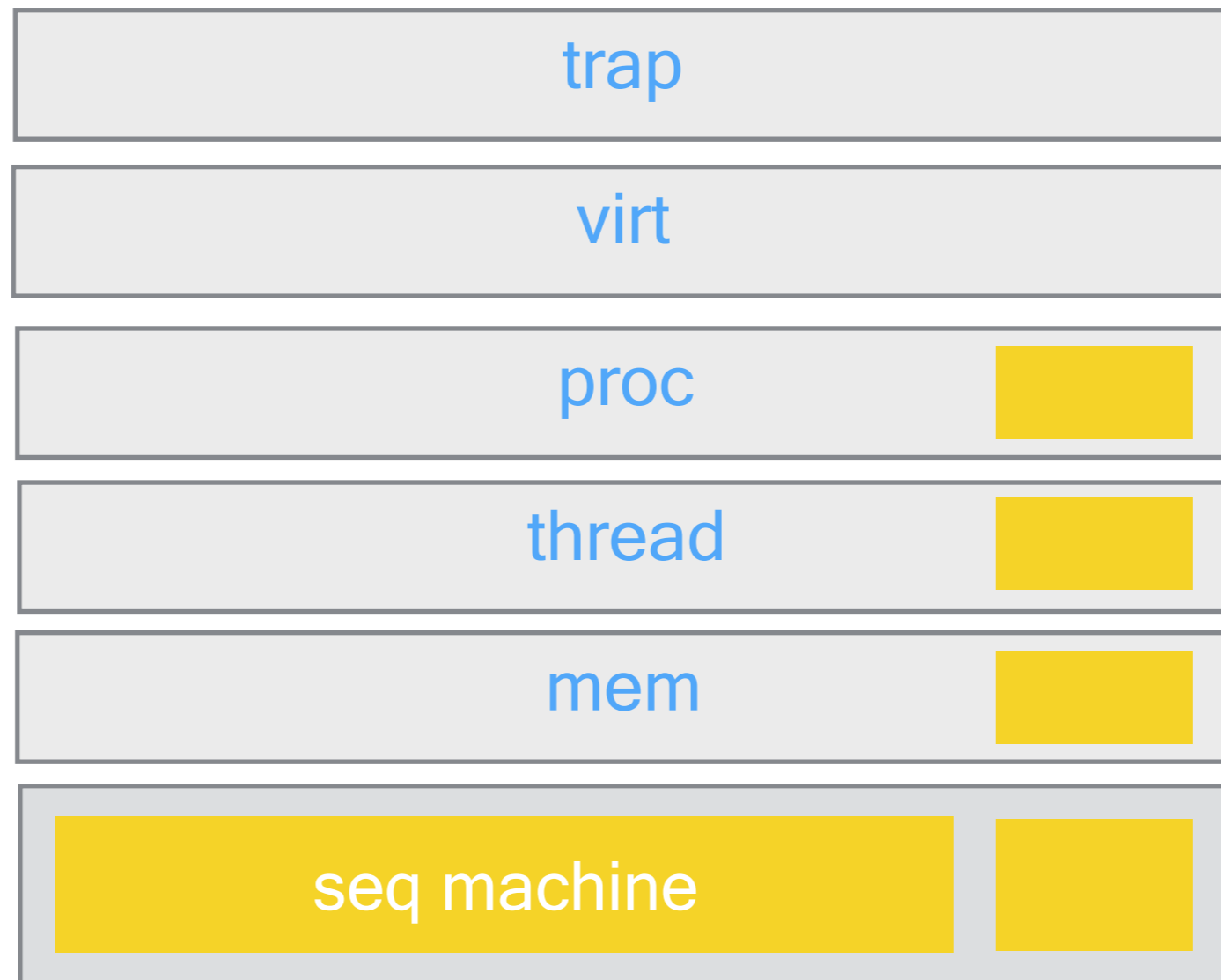


- Trap
- VM
- PM
- TM
- MM

contributions

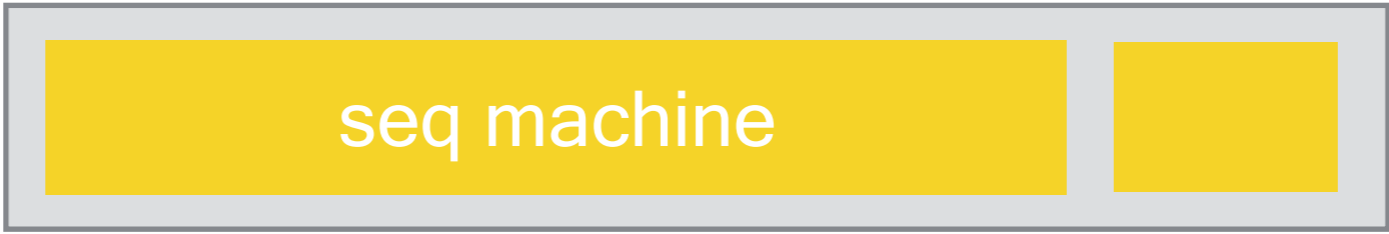


verified hypervisor



- Trap
- VM
- PM
- TM
- MM

contributions



extensibility is the key to support
concurrency

support concurrency

contributions

trap

virt

proc

thread

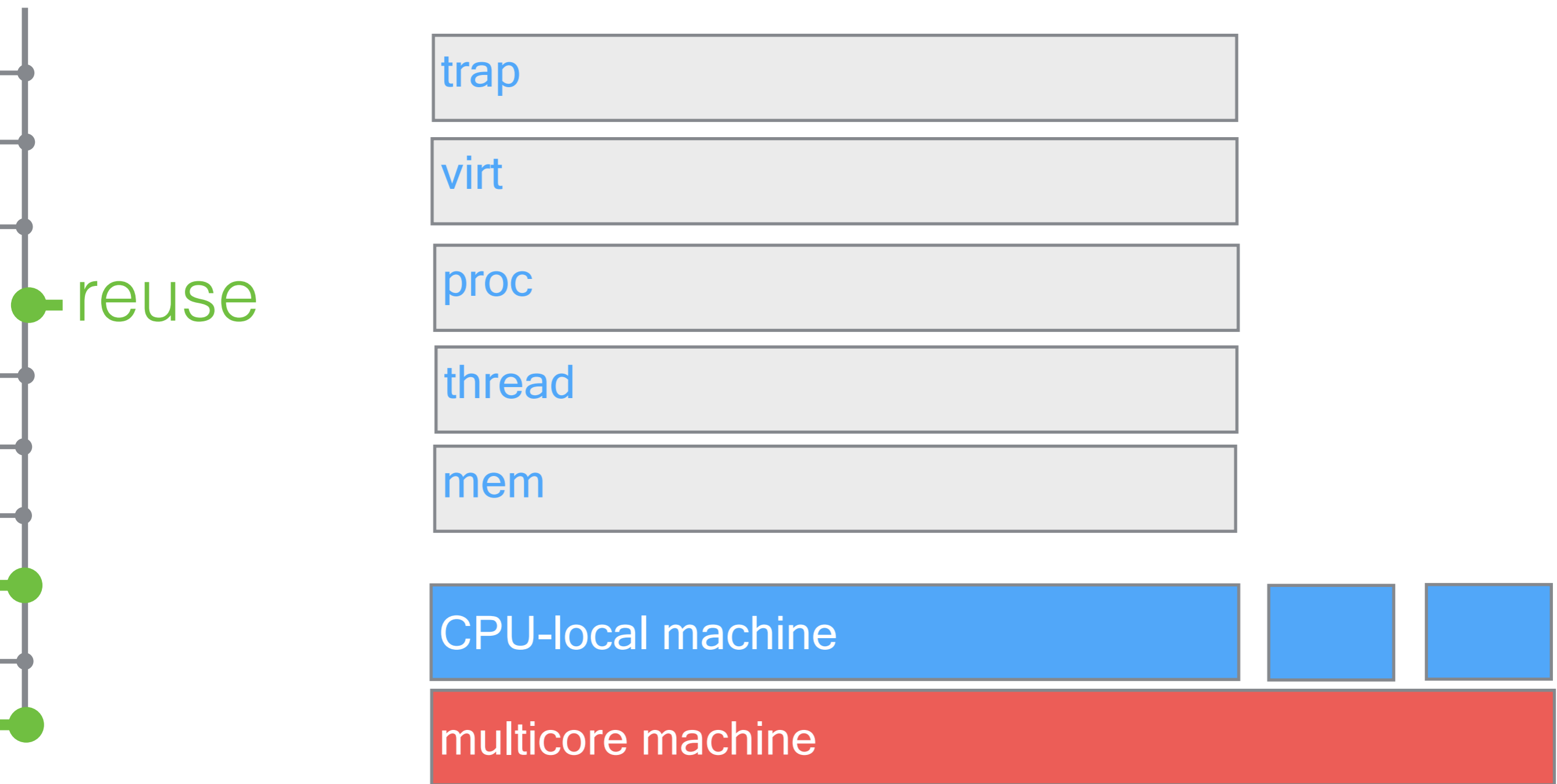
mem

seq machine

multicore machine



contributions



contributions



trap

virt

proc

thread

mem

spin-lock

CPU-local machine

multicore machine

contributions

reuse
mix of 3

trap

virt

proc

thread-local machine

thread

mem

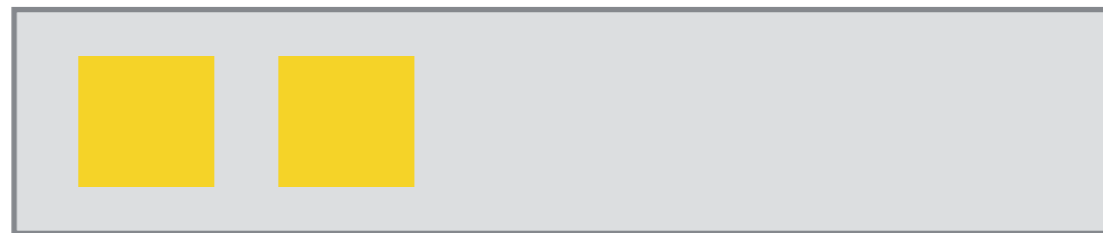
spin-lock

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



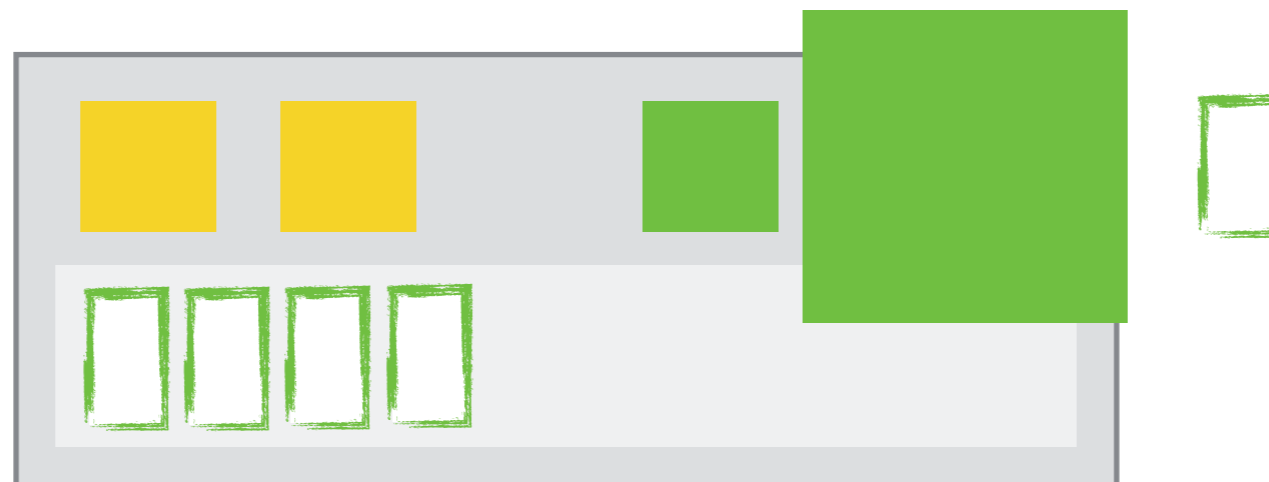
atomic objects



logical log

a sequence of events

certified concurrent layers



trap

virt

proc

thread

thread

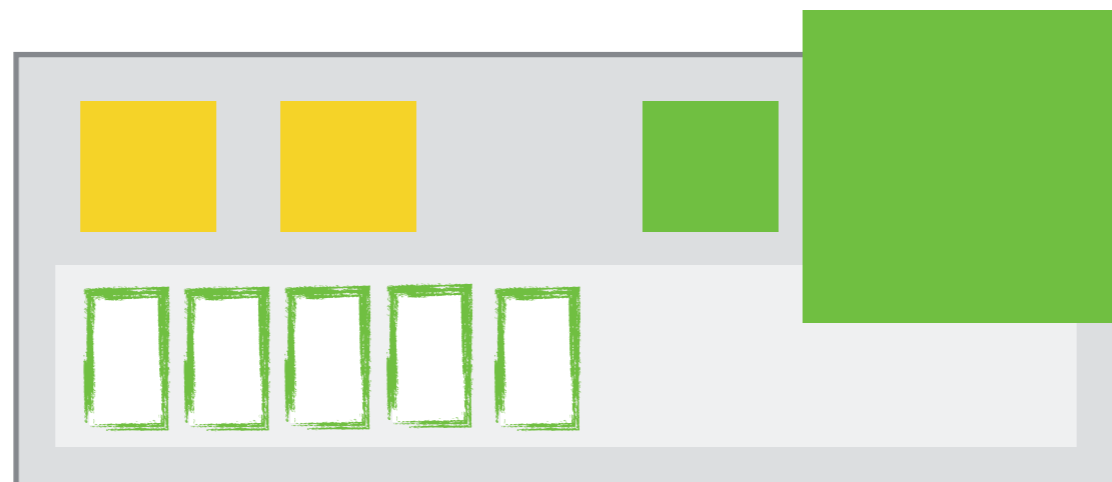
mem

spin-lo

CPU-lo

multico

certified concurrent layers



trap

virt

proc

thread

thread

mem

spin-lo

CPU-Id

multico

certified concurrent layers



trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico



share 



trap

virt

proc

thread

thread

mem

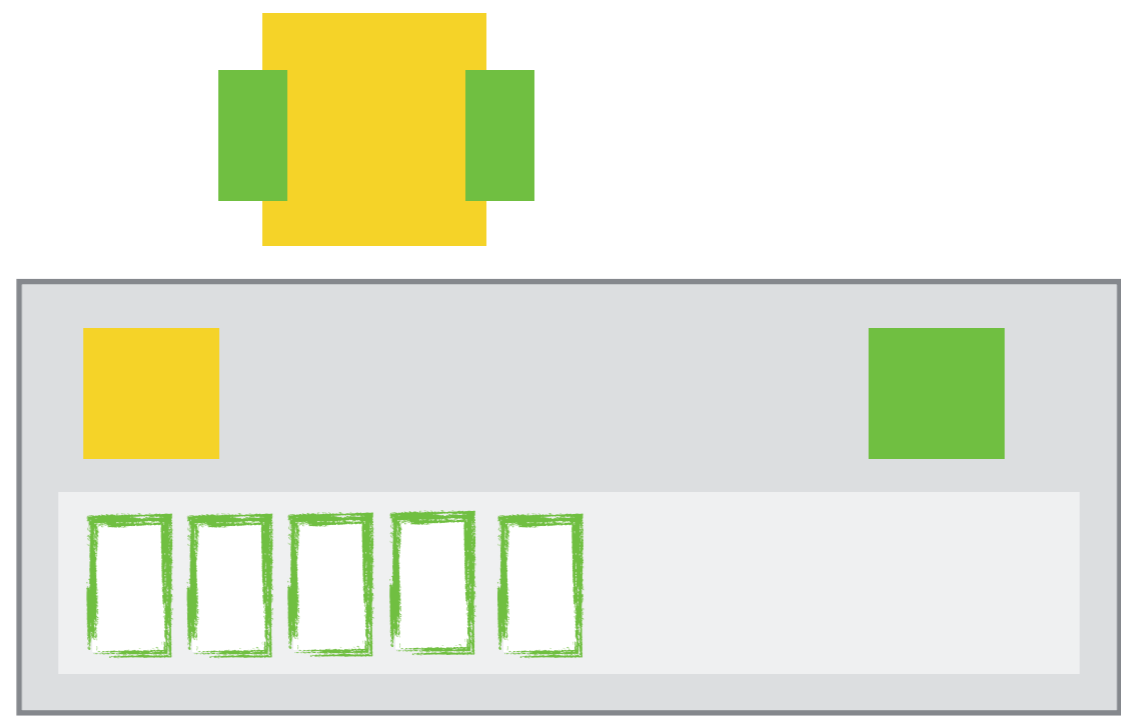
spin-lo

CPU-lo

multico



fine-grained lock



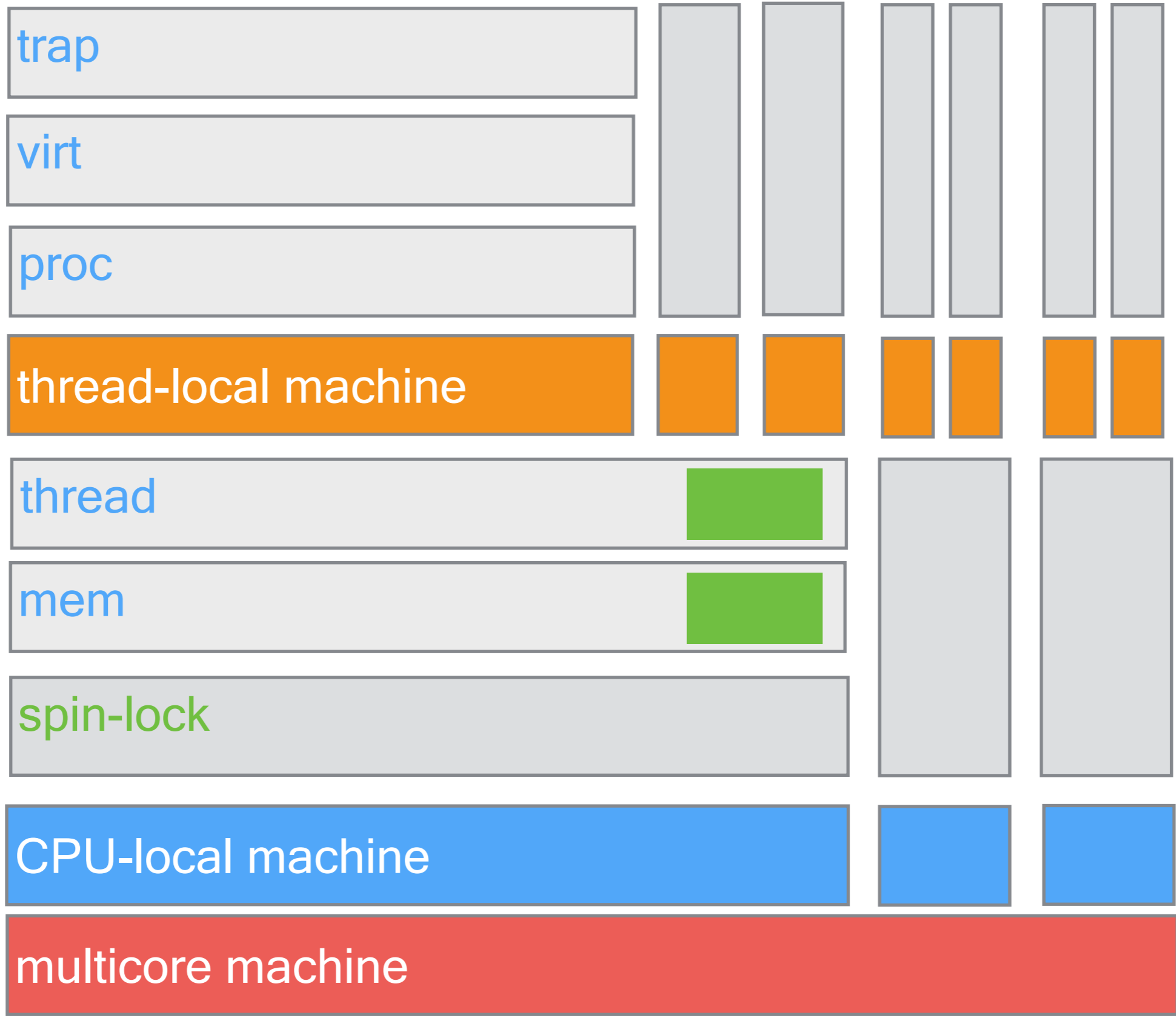
- trap
- virt
- proc
- thread
- thread
- mem
- spin-lo
- CPU-lo
- multico



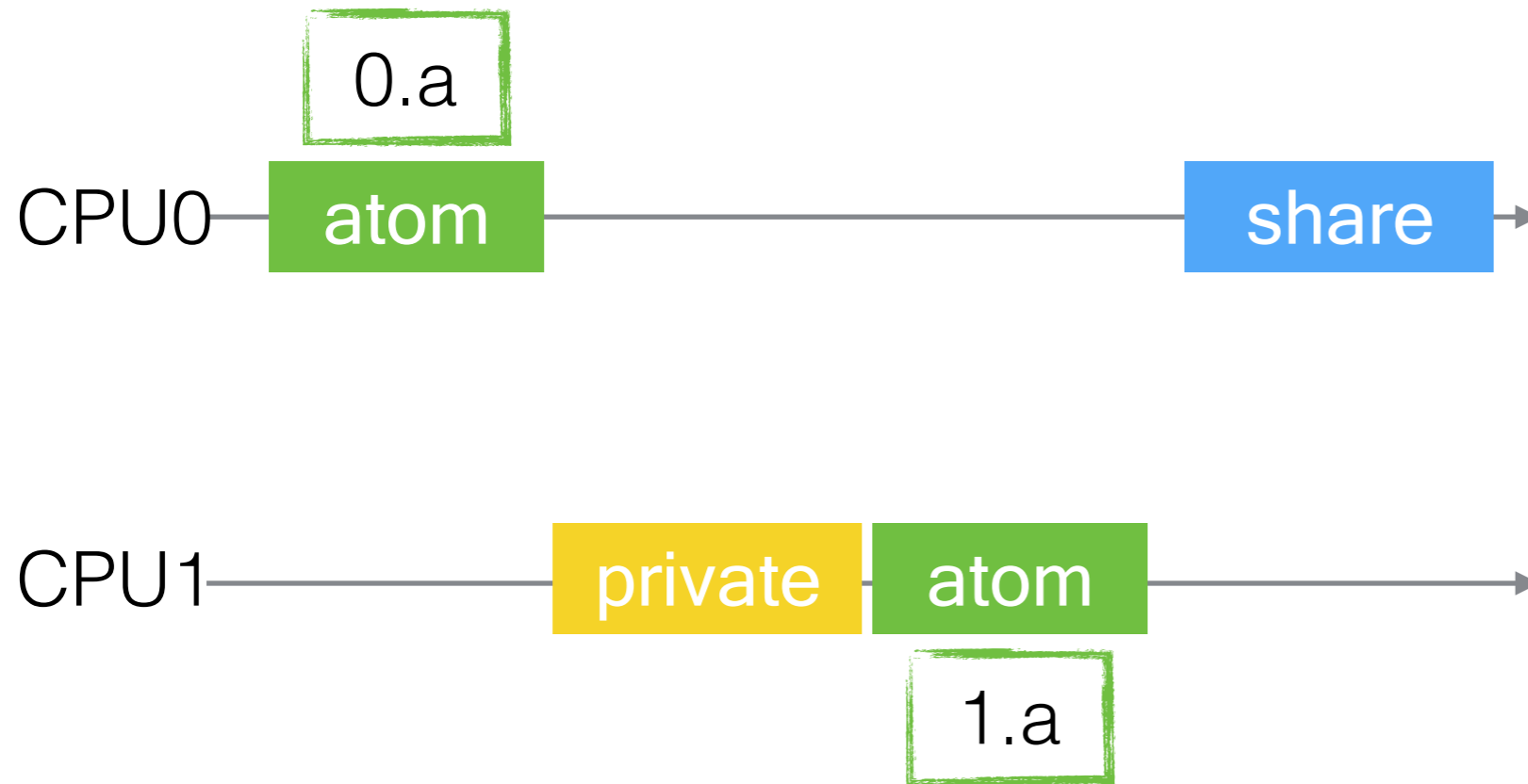
fine-grained lock



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



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



trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multicore machine

step 0: raw x86 multicore model



logical log



multicore machine

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

step 0: raw x86 multicore model

non-determinism

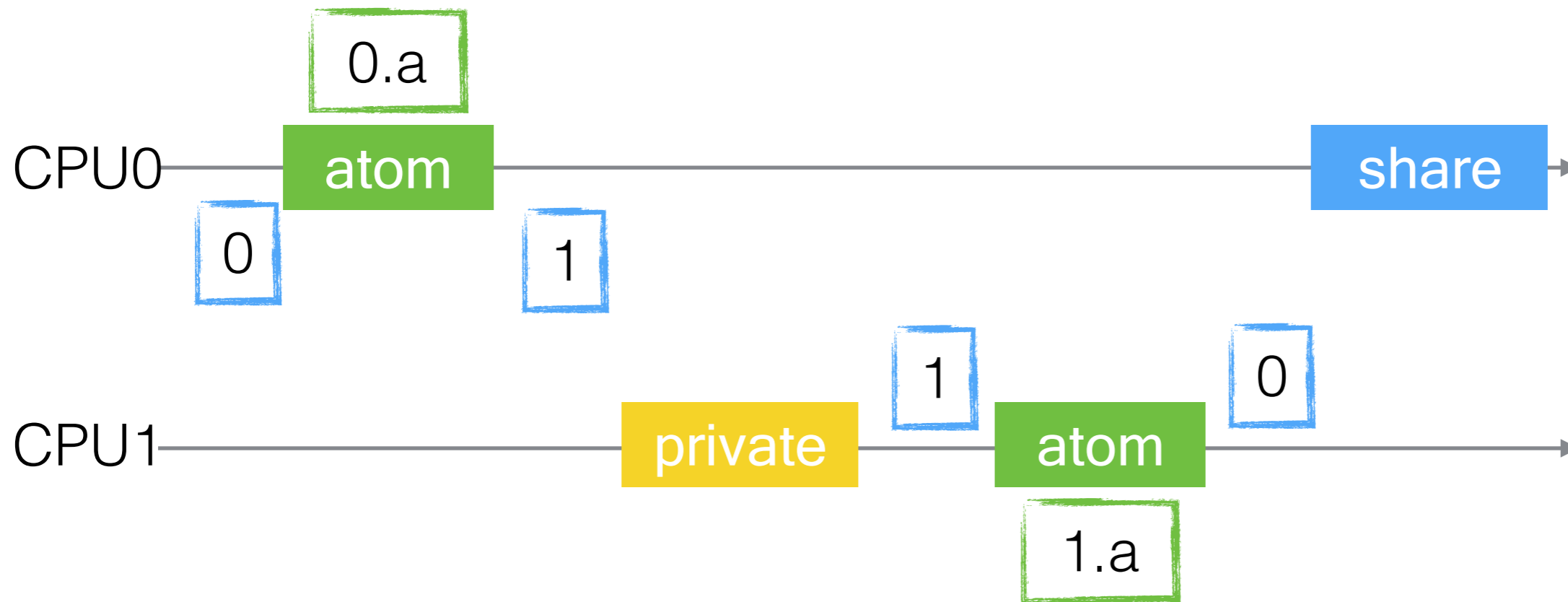


multicore machine

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

step 0: raw x86 multicore model

non-determinism

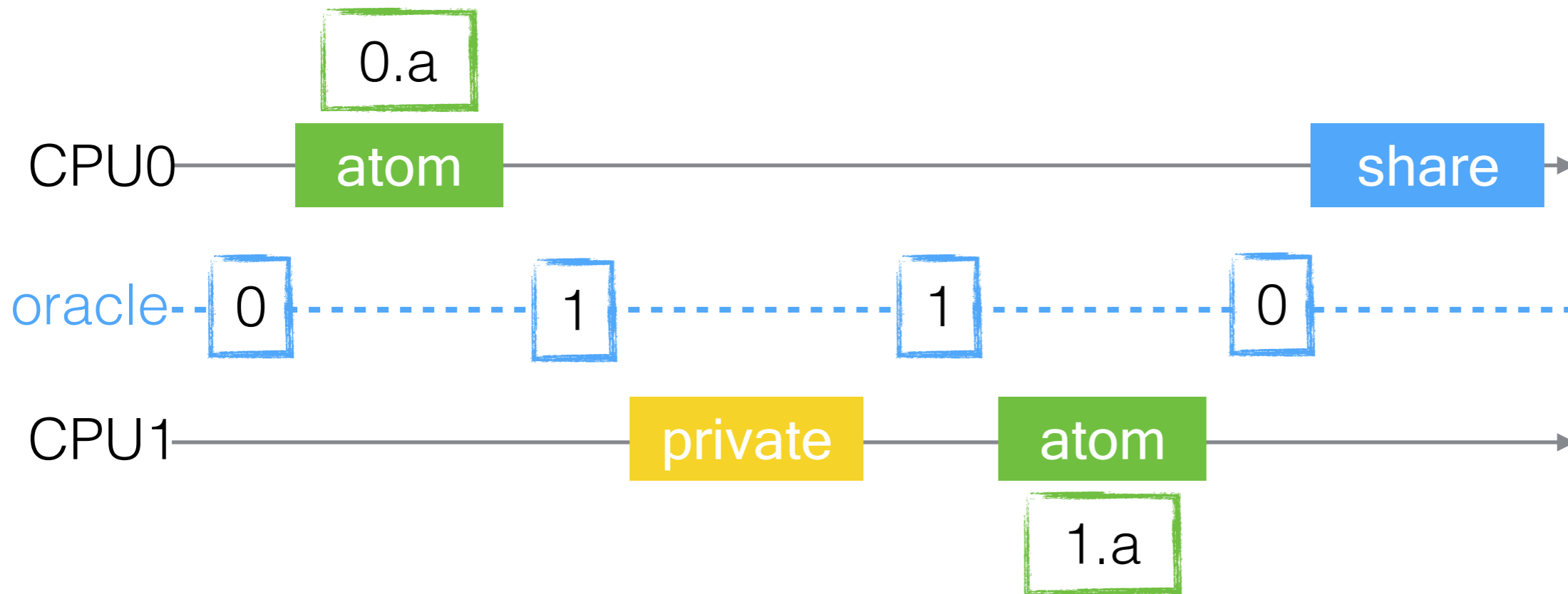


multicore machine

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

step 0: raw x86 multicore model

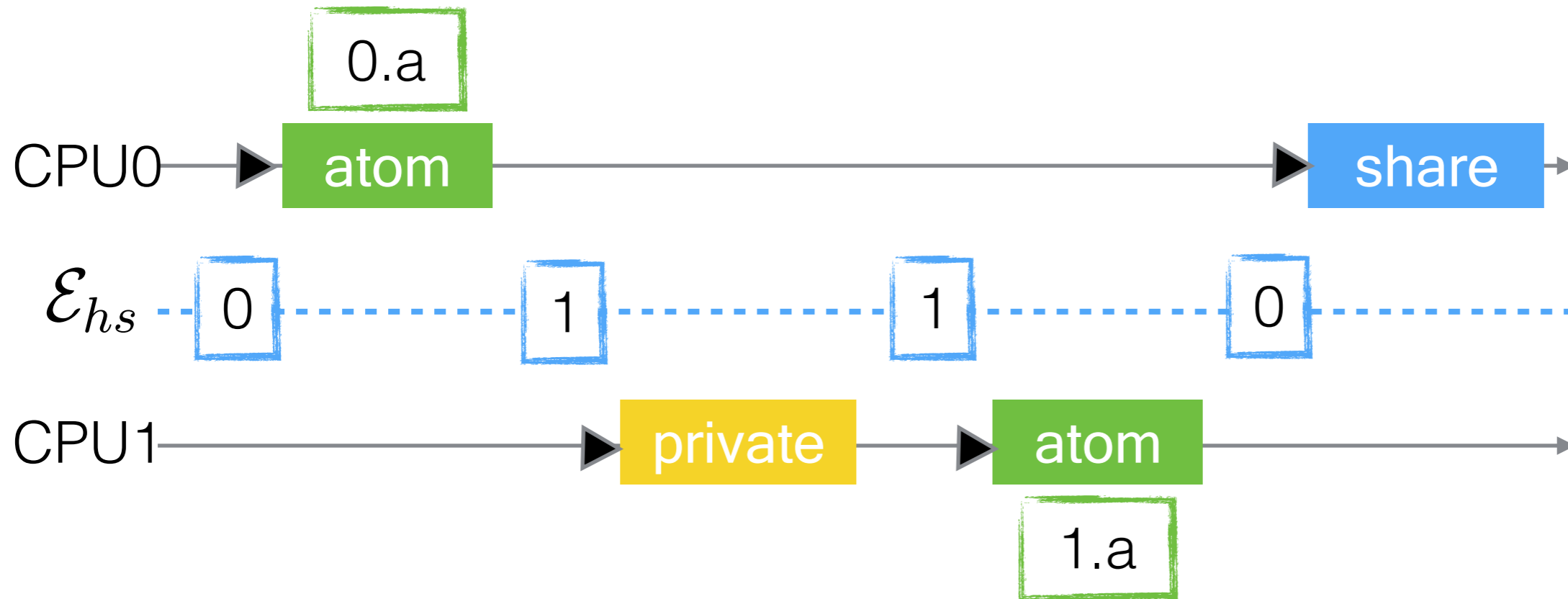
non-determinism



multicore machine

step 1: hardware scheduler \mathcal{E}_{hs}

purely logical

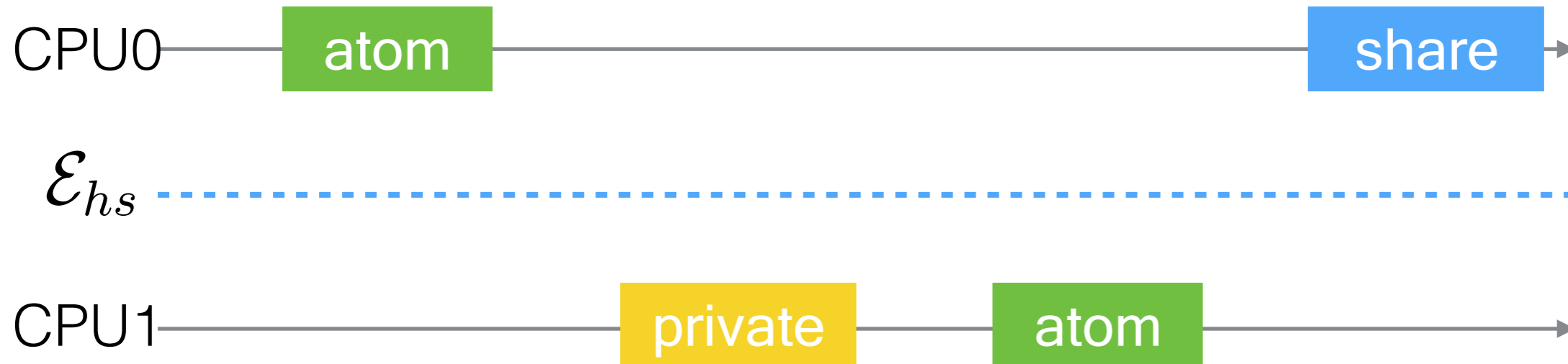


multicore machine

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

step 1: hardware scheduler \mathcal{E}_{hs}

purely logical



multicore machine

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

step 1: hardware scheduler
purely logical

$$\forall \mathcal{E}_{hs}$$

trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multicore machine

step 1: hardware scheduler



trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

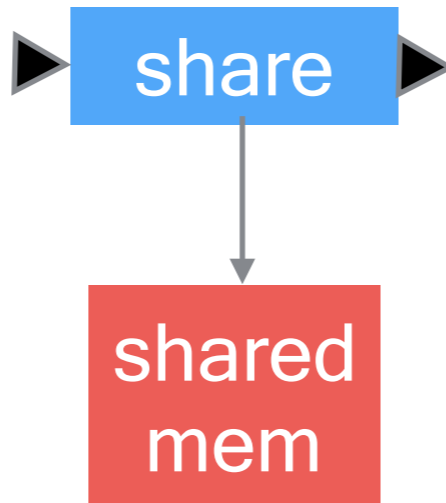
$\forall \mathcal{E}_{hs}$

machine with hardware scheduler

multicore machine

step 2: push/pull model

CPU0



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

$\forall \mathcal{E}_{hs}$

machine with hardware scheduler

multicore machine

step 2: push/pull model



logical copy

shared mem

$\forall \mathcal{E}_{hs}$

machine with hardware scheduler

multicore machine

trap

virt

proc

thread

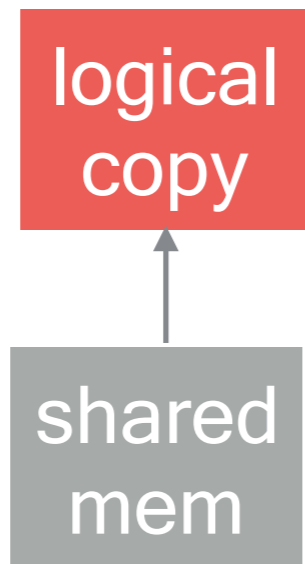
thread

mem

spin-lo

CPU-lo

step 2: push/pull model



$\forall \mathcal{E}_{hs}$

machine with hardware scheduler

multicore machine

trap

virt

proc

thread

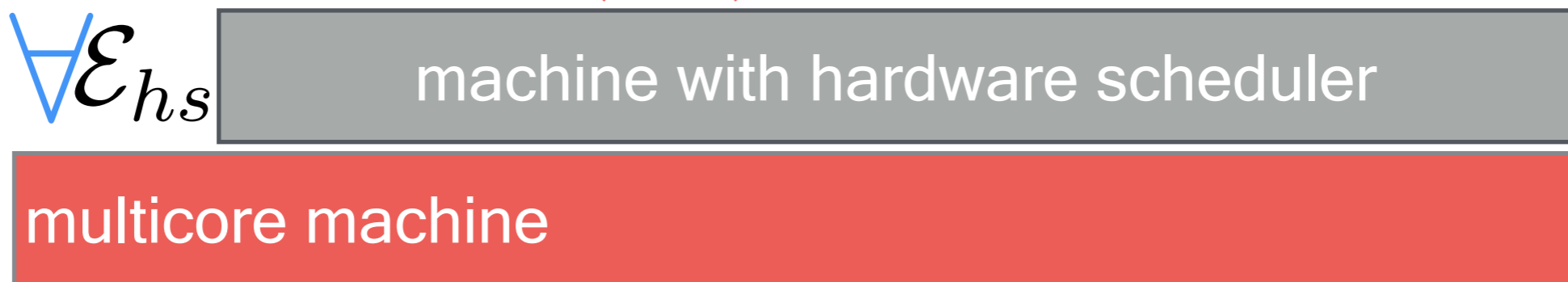
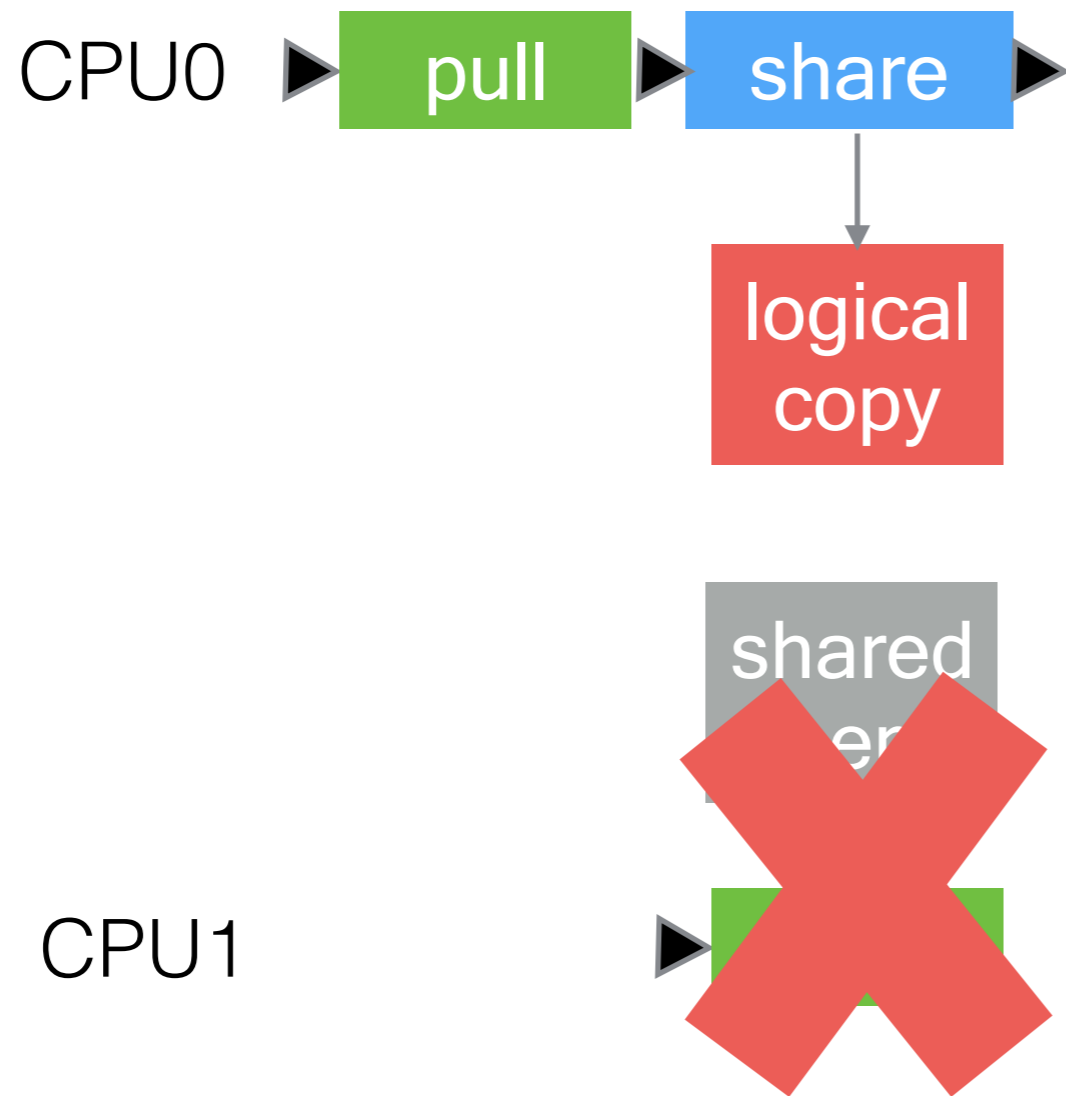
thread

mem

spin-lo

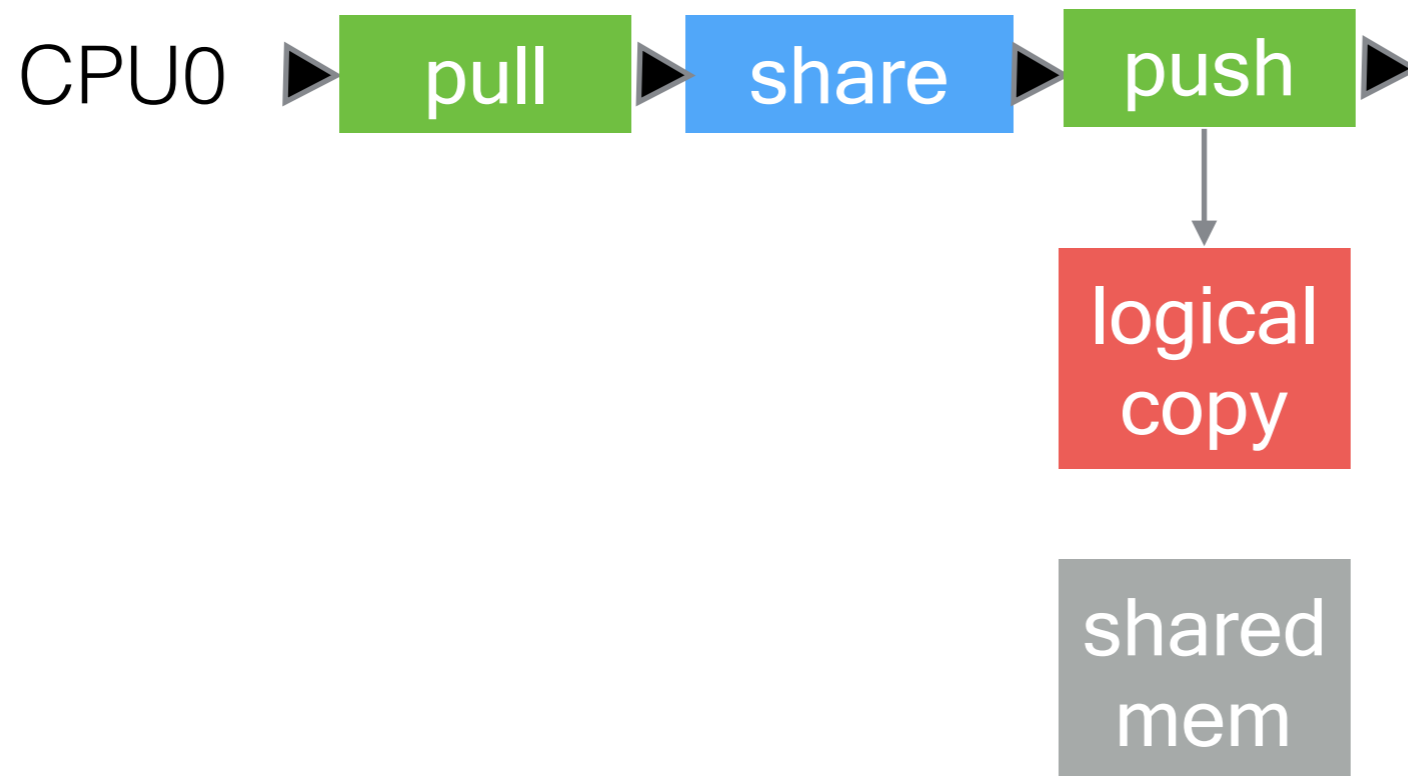
CPU-lo

step 2: push/pull model



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

step 2: push/pull model

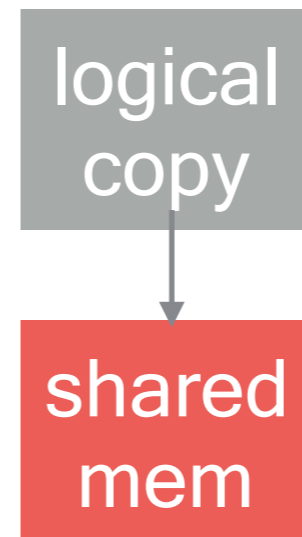


- trap
- virt
- proc
- thread
- thread
- mem
- spin-lo
- CPU-lo

$\forall \mathcal{E}_{hs}$ machine with hardware scheduler

multicore machine

step 2: push/pull model



$\forall \mathcal{E}_{hs}$

machine with hardware scheduler

multicore machine

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



$\forall \mathcal{E}_{hs}$

multicore machine

machine with local copy

machine with hardware scheduler

trap

virt

proc

thread

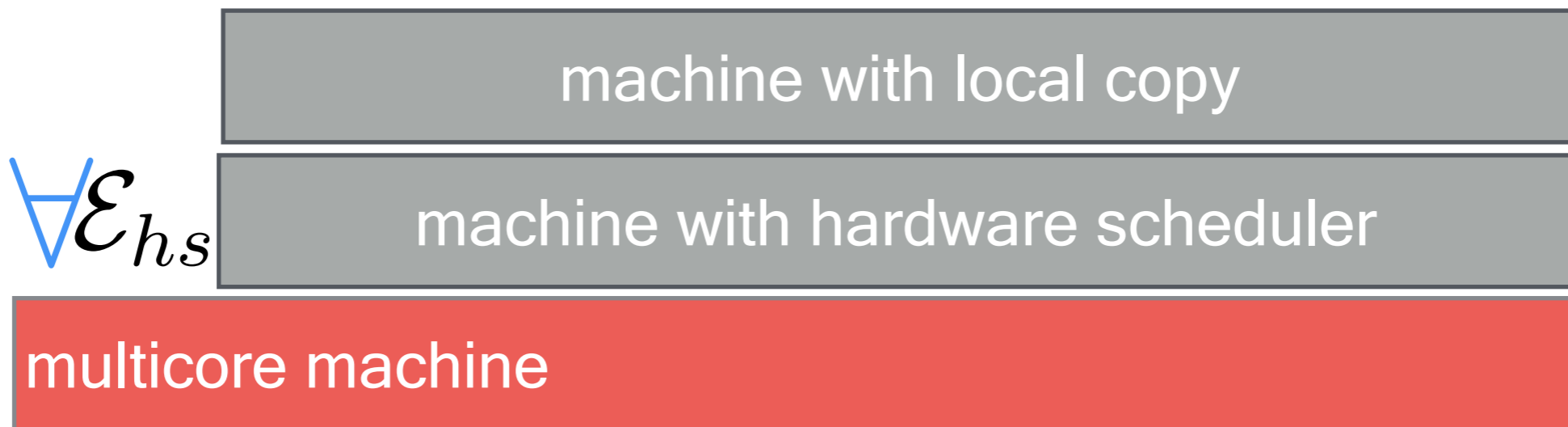
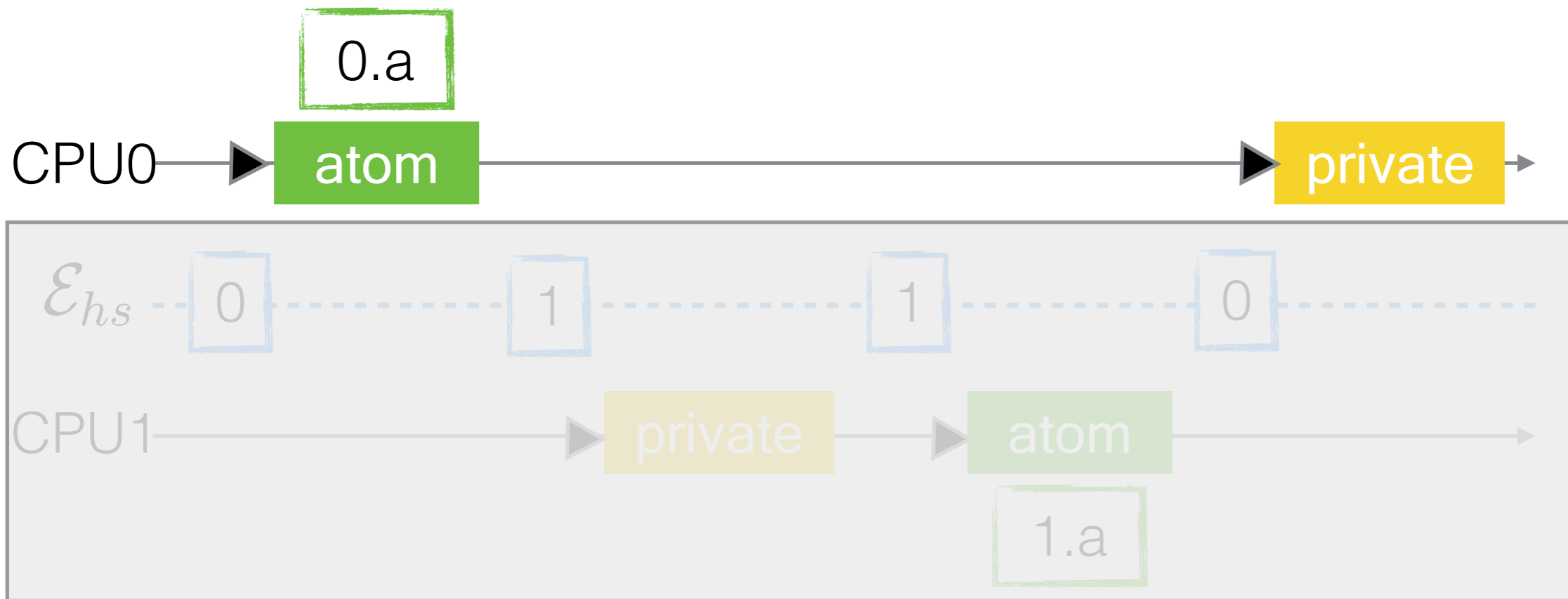
thread

mem

spin-lo

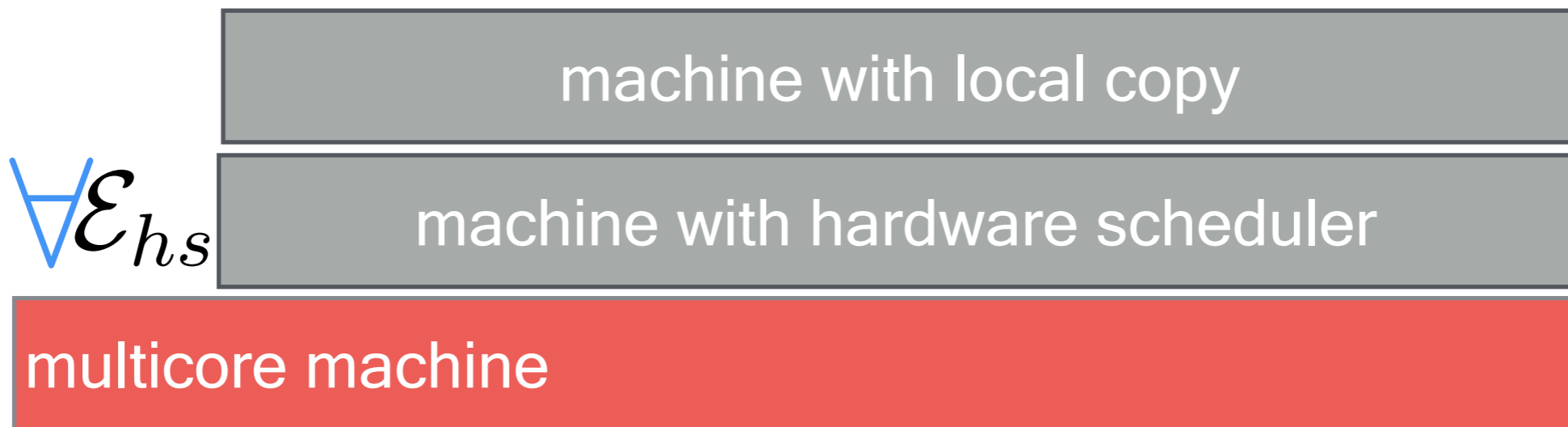
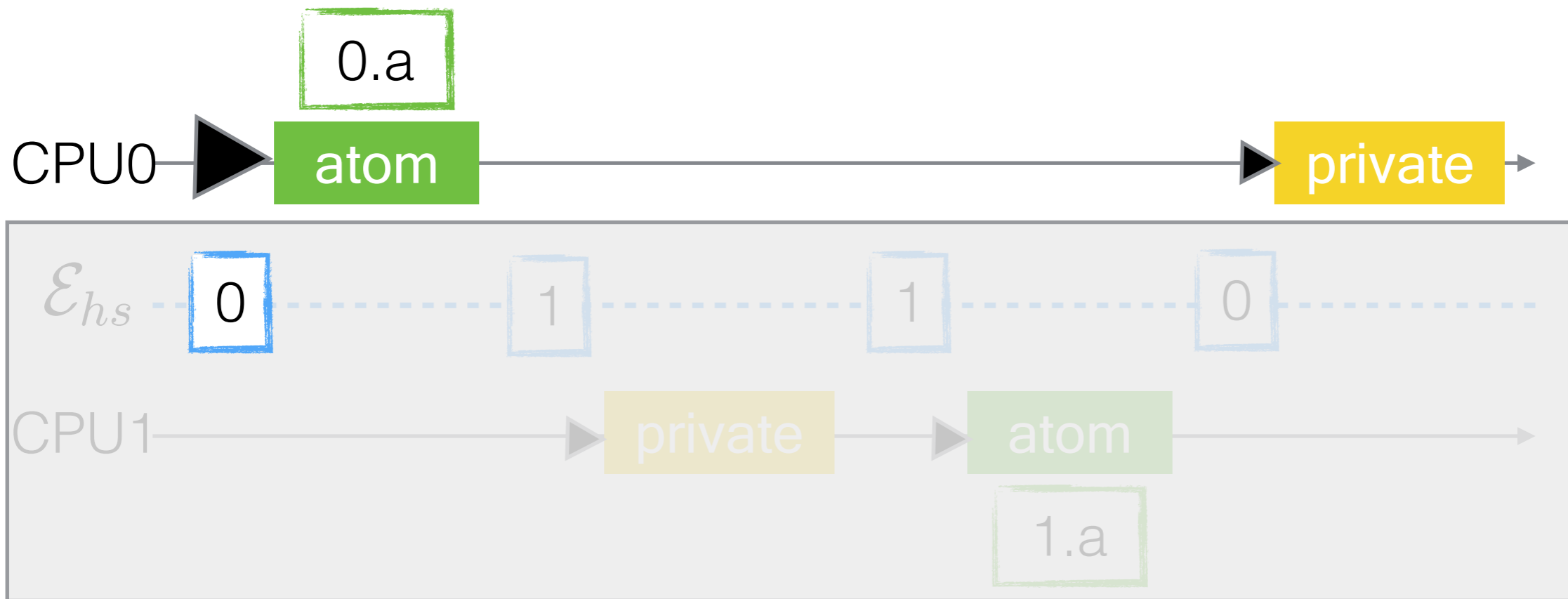
CPU-lo

step 3: per-CPU machine



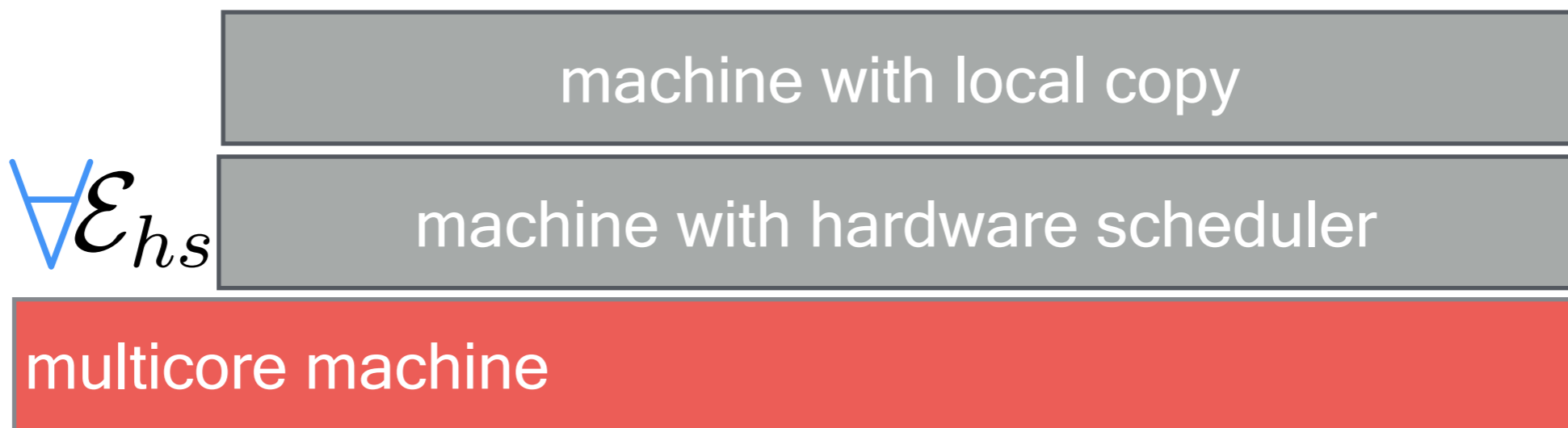
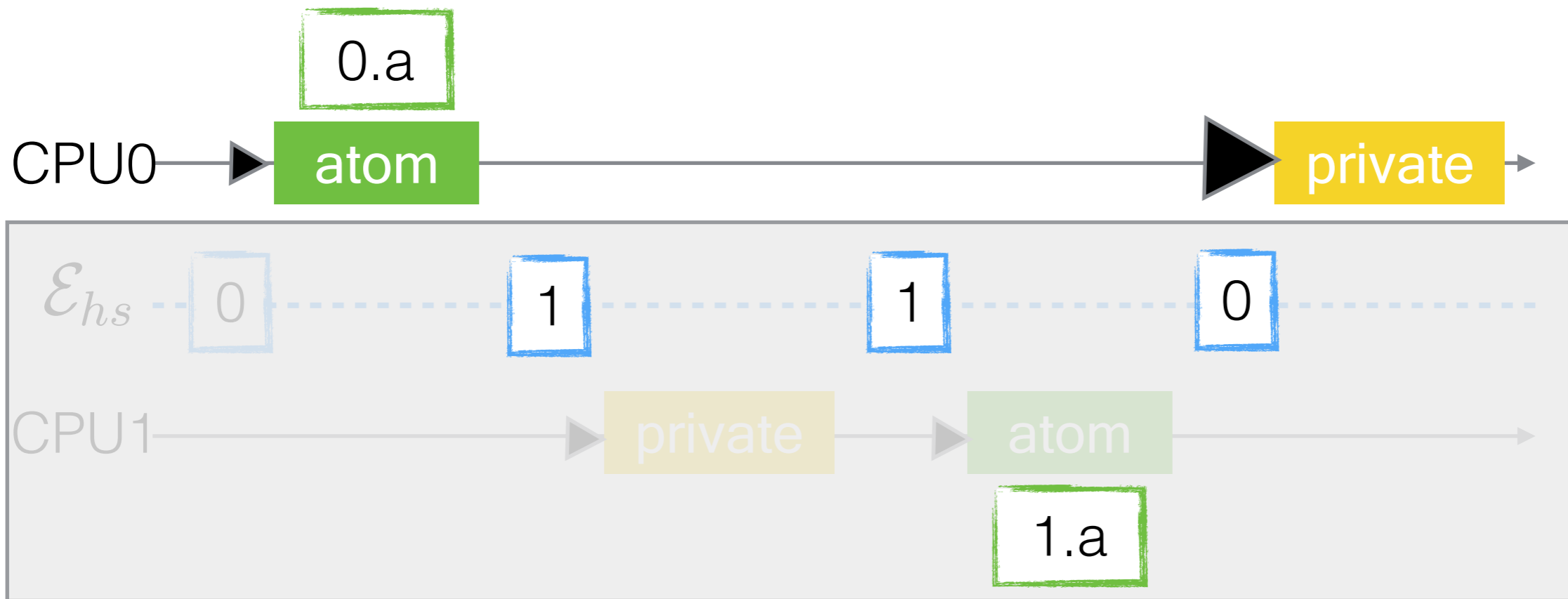
- trap
- virt
- proc
- thread
- thread
- mem
- spin-lo
- CPU-lo

step 3: per-CPU machine



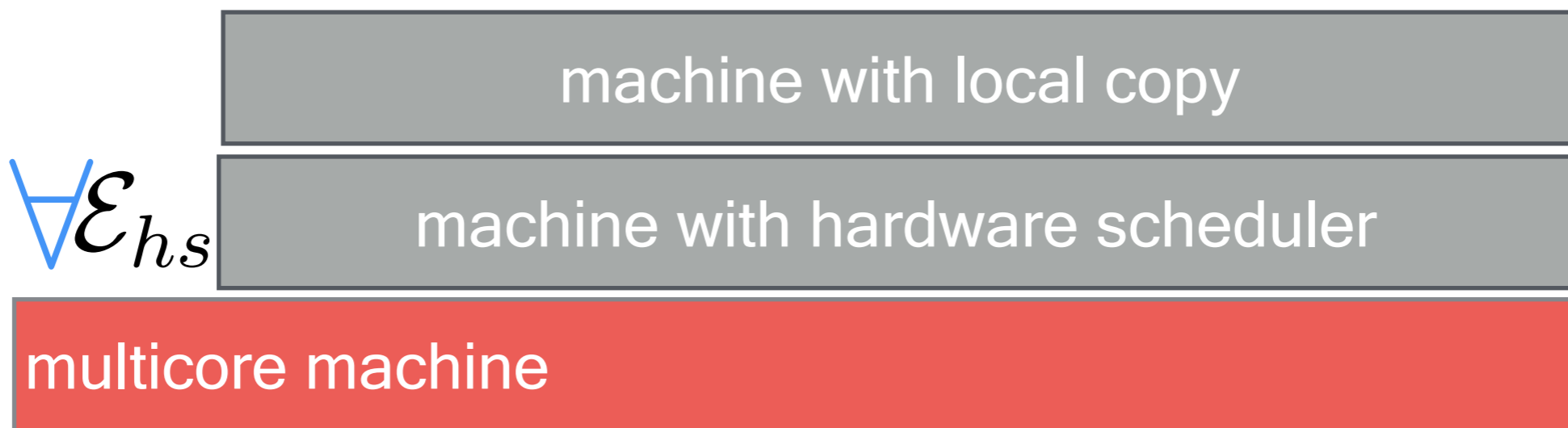
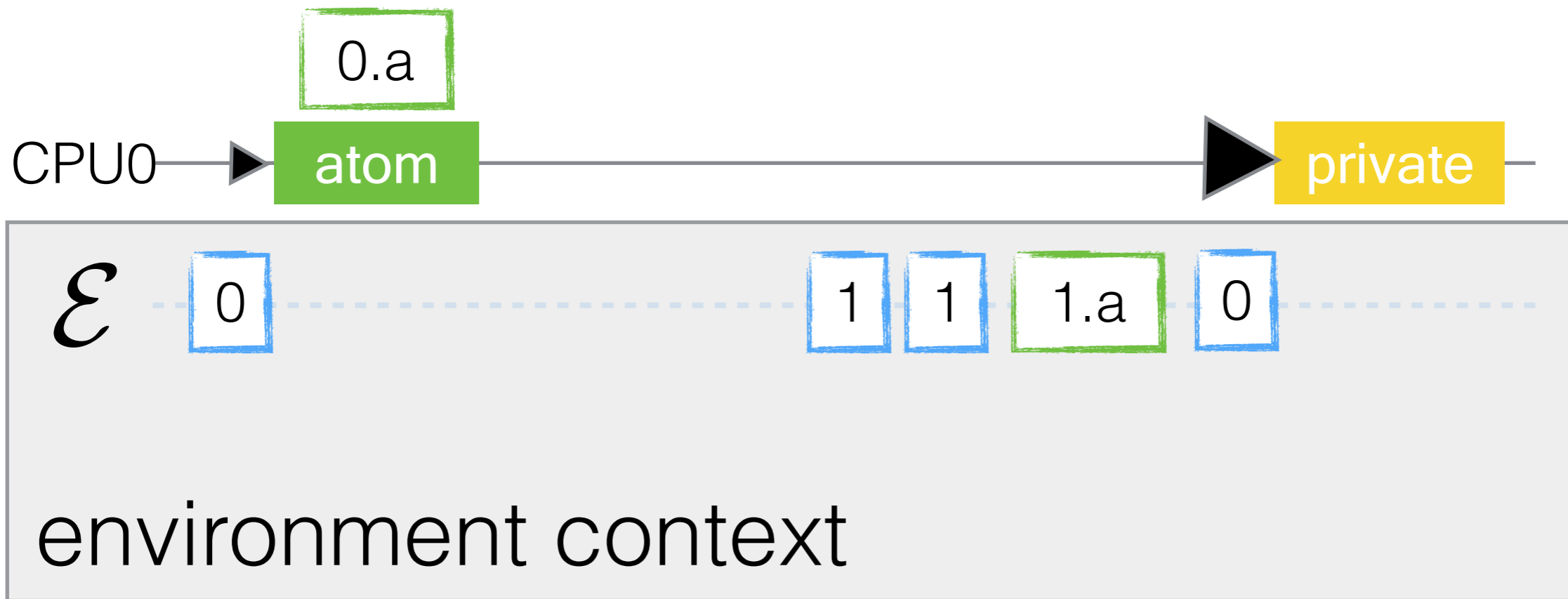
- trap
- virt
- proc
- thread
- thread
- mem
- spin-lo
- CPU-lo

step 3: per-CPU machine



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

step 3: per-CPU machine



- trap
- virt
- proc
- thread
- thread
- mem
- spin-loc
- CPU-loc



$\forall \mathcal{E}_{hs}$

multicore machine

CPU i machine

CPU j machine

machine with local copy

machine with hardware scheduler

trap

virt

proc

thread

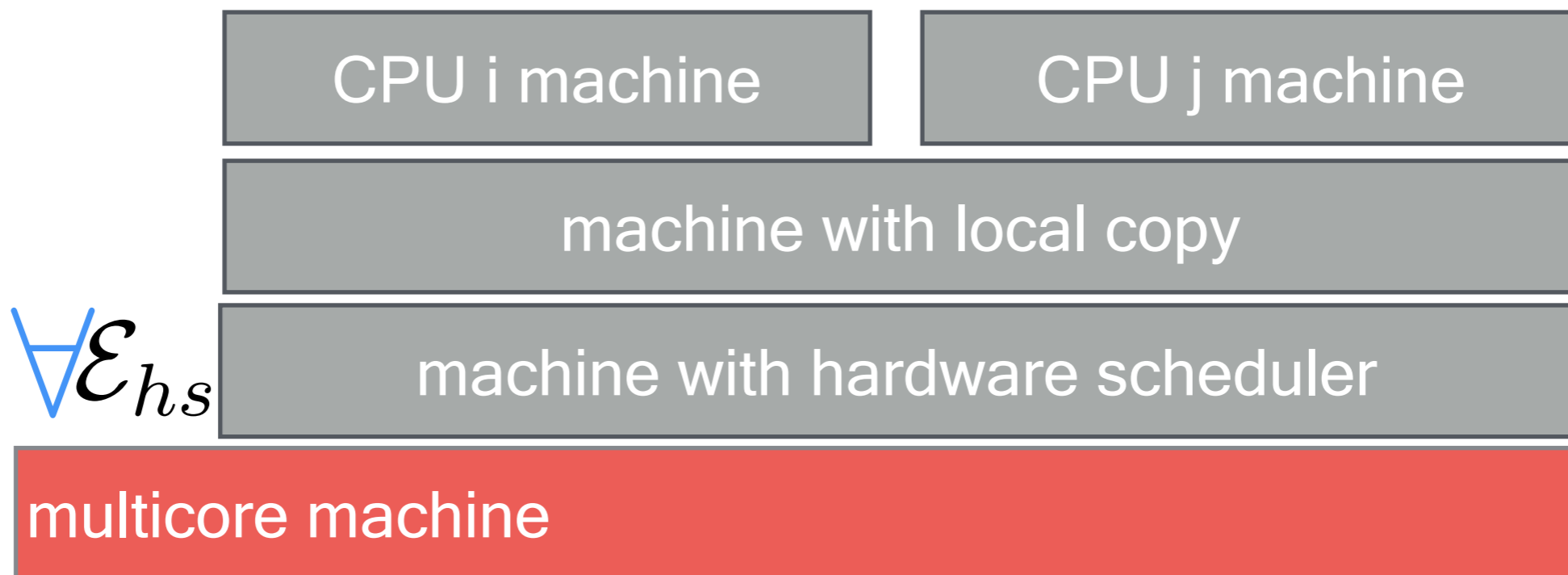
thread

mem

spin-lo

CPU-lo

step 4: remove unnecessary interleaving



trap

virt

proc

thread

thread

mem

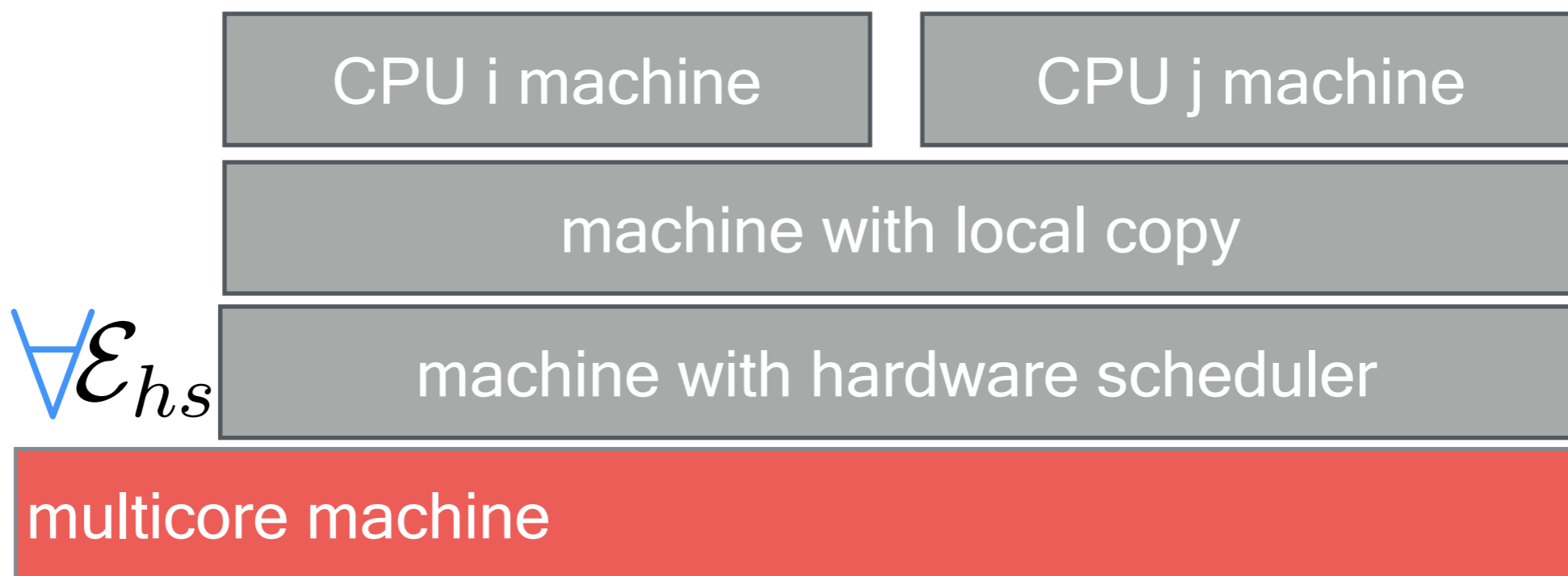
spin-lo

CPU-lo

step 4: remove unnecessary interleaving



shuffle



trap

virt

proc

thread

thread

mem

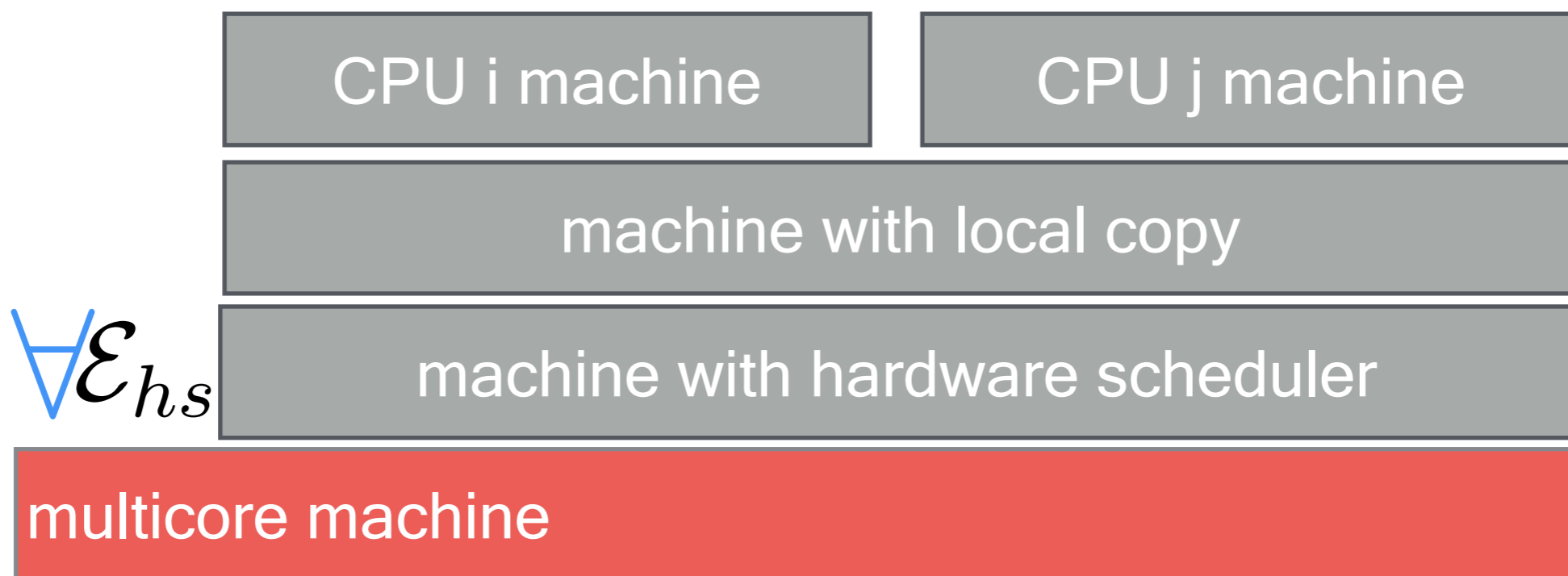
spin-lo

CPU-lo

step 4: remove unnecessary interleaving



merge



trap

virt

proc

thread

thread

mem

spin-lo

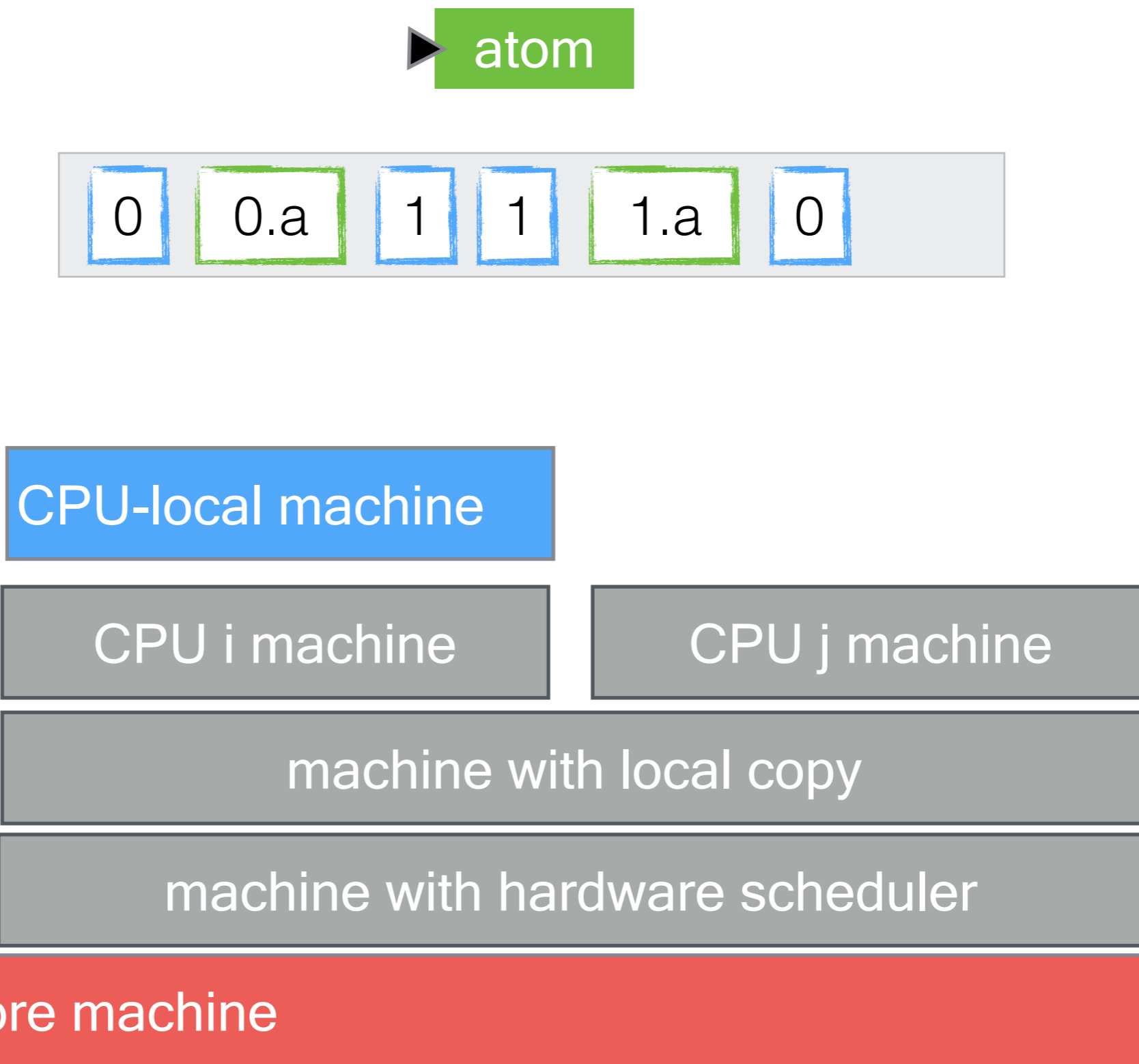
CPU-lo

contributions

reuse

$\forall \mathcal{E}_{hs}$

multicore machine



trap

virt

proc

thread

thread

mem

spin-lo



trap

virt

proc

thread

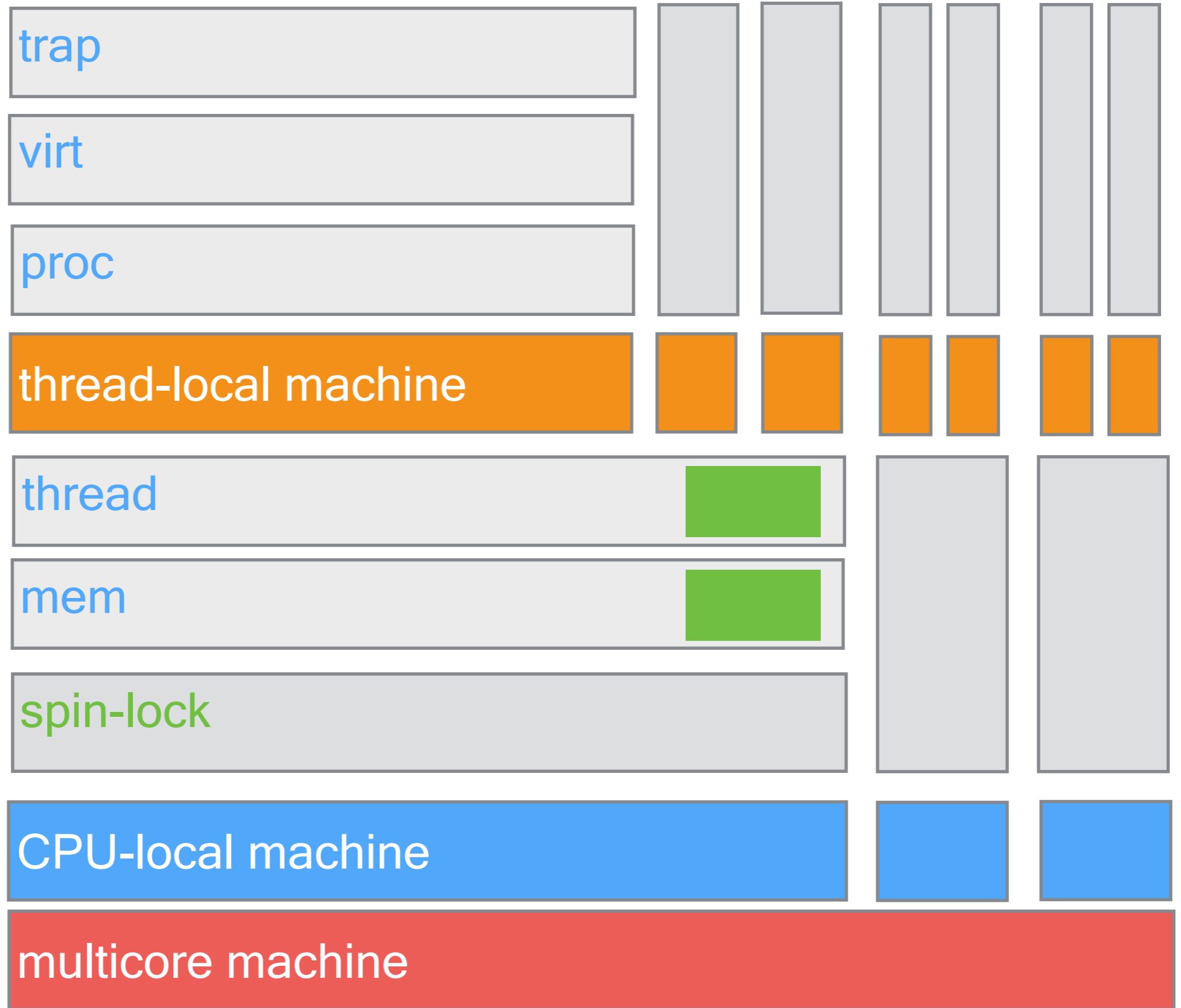
thread

mem

spin-lo

CPU-lo

multico



acq-lock specification



safely
pull

logical
copy



spin-lock

- trap
- virt
- proc
- thread
- thread
- mem
- CPU-Id
- multico

acq-lock specification



safely
pull

logical
copy

pull will
eventually return

trap

virt

proc

thread

thread

mem

CPU-Id

multico

spin-lock

acq-lock specification

mutual
exclusion

logical
copy

liveness

spin-lock

trap

virt

proc

thread

thread

mem

CPU-Id

multico



ticket lock

mutual exclusion + liveness

```
void acq_lock (uint i)
{
  uint t = FAI_ticket (i)
  while (get_now (i) != t)
  {}
  pull (i);
}
```

FAI
ticket

spin-lock

trap

virt

proc

thread

thread

mem

CPU-Id

multico

mutual exclusion + liveness

```
void acq_lock (uint i)
{
  uint t = FAI_ticket (i);
  while (get_now (i)
  {
  }
  pull (i);
}
```

FAI
ticket

get
now

spin-lock

trap

virt

proc

thread

thread

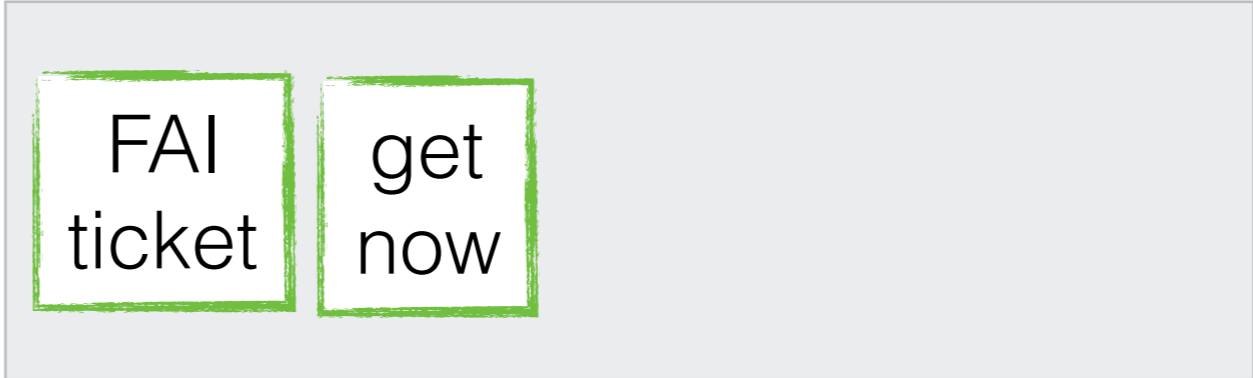
mem

CPU-Id

multico

mutual exclusion + liveness

```
void acq_lock (uint i)
{
  uint t = FAI_ticket (i);
  while (get_now (
  { }
  pull (i);
}
```

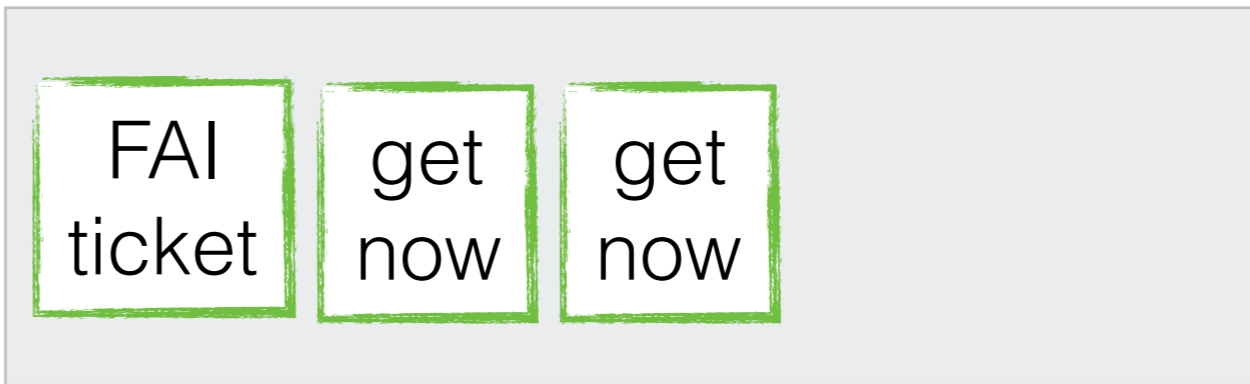


spin-lock

- trap
- virt
- proc
- thread
- thread
- mem
- CPU-Id
- multico

mutual exclusion + liveness

```
void acq_lock (uint i)
{
  uint t = FAI_ticket (i);
  while (get_now (i) != t)
  {}
  pull (i)
}
```

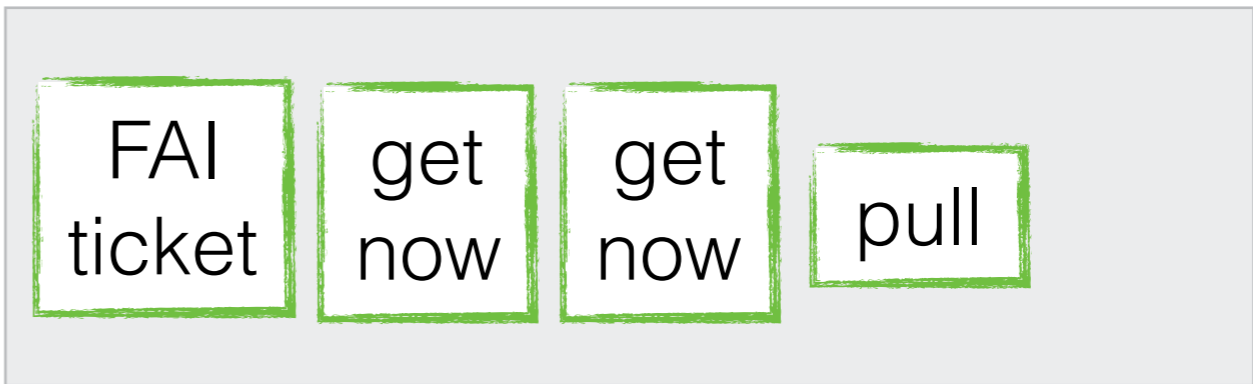


spin-lock

- trap
- virt
- proc
- thread
- thread
- mem
- CPU-Id
- multico

mutual exclusion + liveness

```
void acq_lock (uint i)
{
  uint t = FAI_ticket (i);
  while (get_now (i) != t)
  {}
  pull (i);
}
```



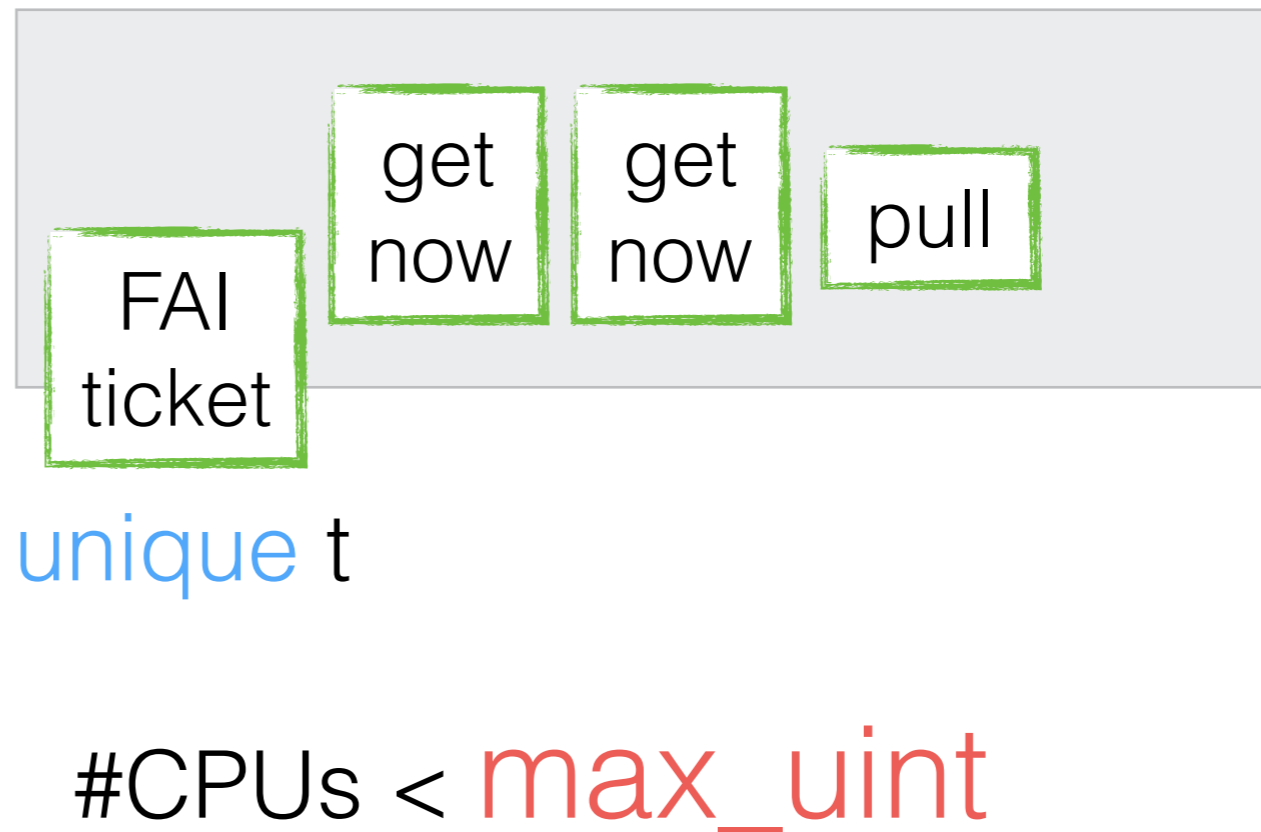
spin-lock

- trap
- virt
- proc
- thread
- thread
- mem
- CPU-Id
- multico

mutual exclusion

+ liveness

```
void acq_lock (uint i)
{
  uint t = FAI_ticket (i);
  while (get_now (i) != t)
  {}
  pull (i);
}
```



spin-lock

trap

virt

proc

thread

thread

mem

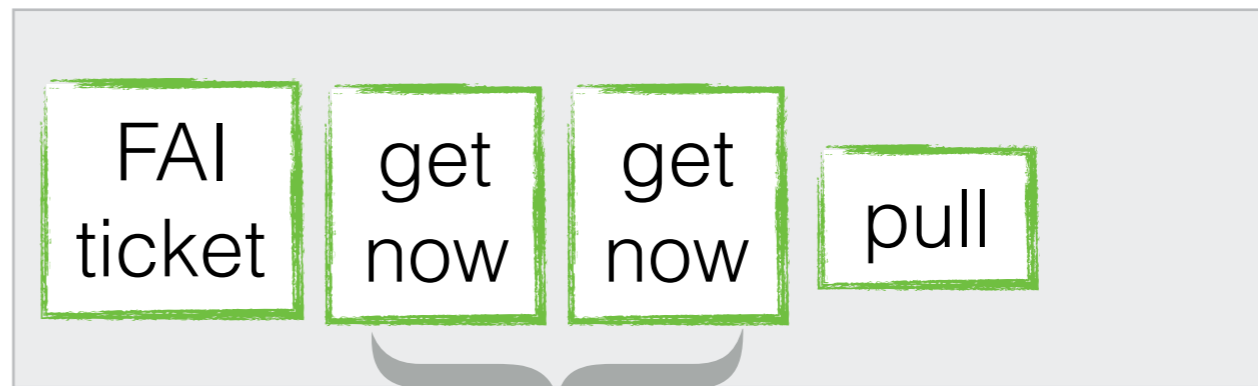
CPU-Id

multico

mutual exclusion + liveness

liveness

```
void acq_lock (uint i)
{
  uint t = FAI_ticket (i);
  while (get_now (i) != t)
  {}
  pull (i);
}
```



#CPUs is bounded

a fair scheduler

lock holders will release lock

spin-lock

trap

virt

proc

thread

thread

mem

CPU-Id

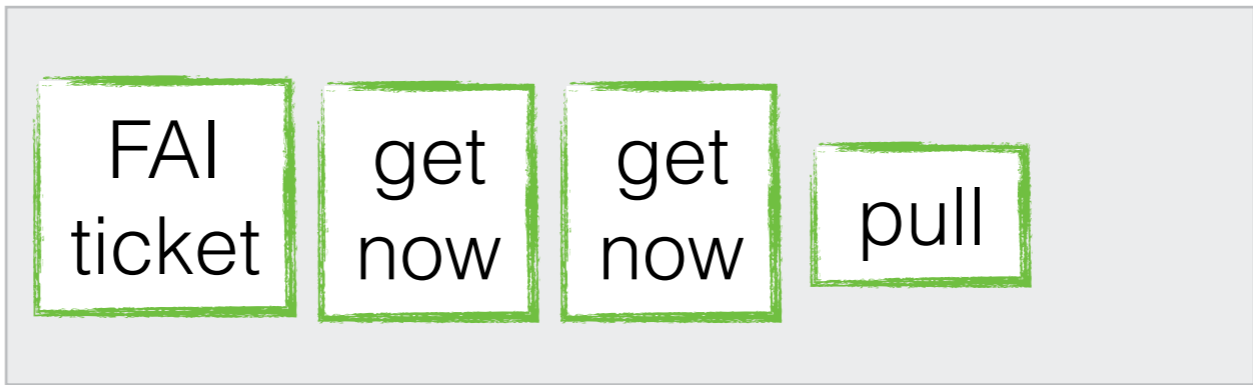
multico



acq_lock

acq_lock

acq
lock



spin-lock

trap

virt

proc

thread

thread

mem

CPU-Id

multico



acq_lock

acq_lock

spin-lock



acq lock

trap

virt

proc

thread

thread

mem

CPU-Id

multico



trap

virt

proc

thread

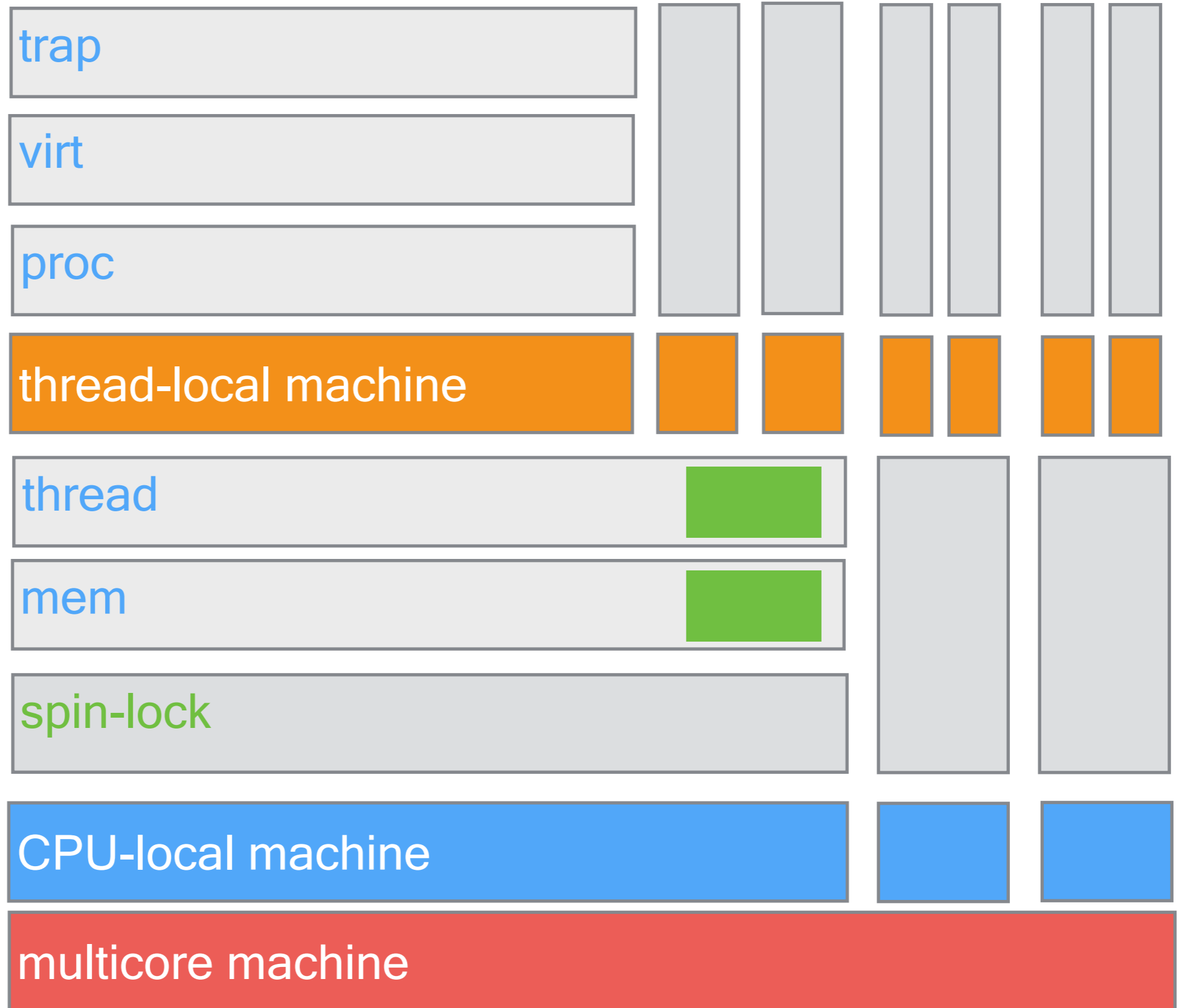
thread

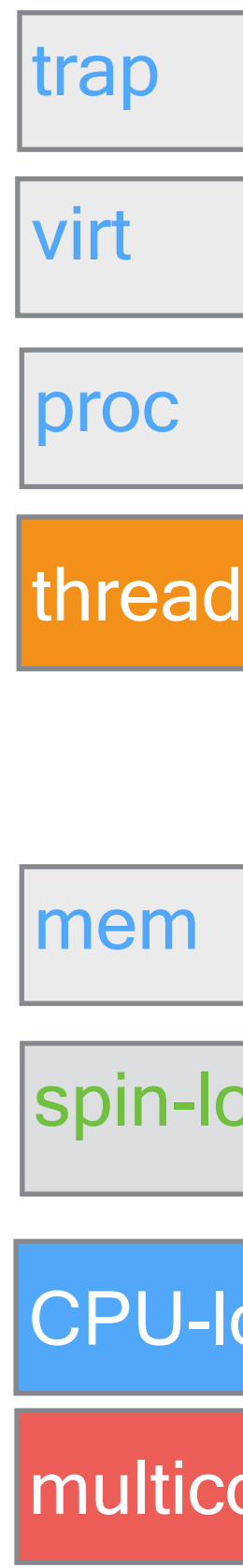
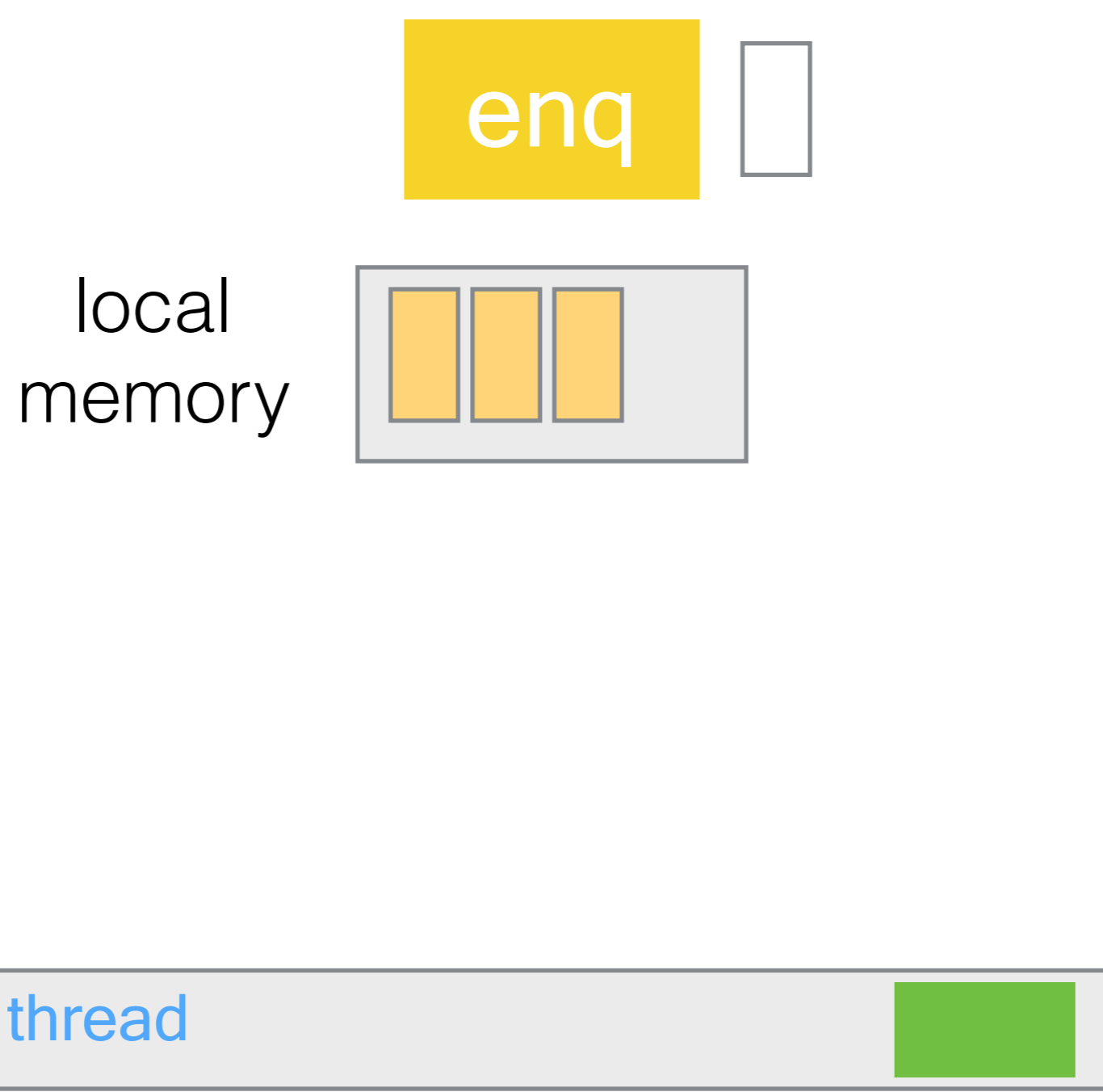
mem

spin-lo

CPU-lo

multico



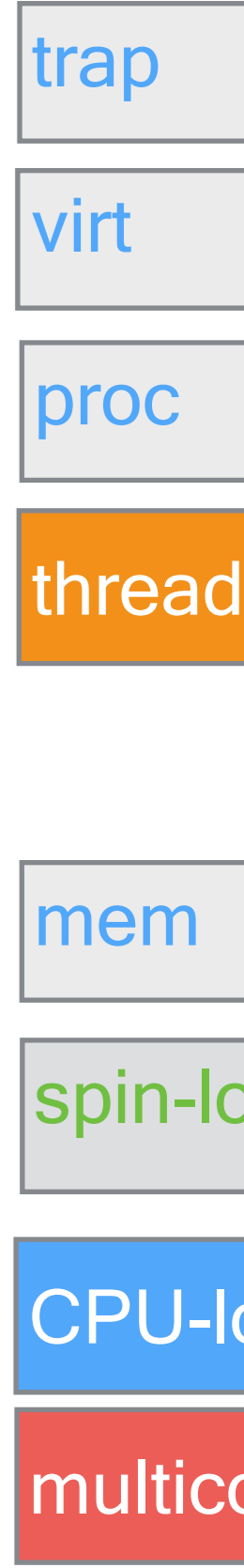




local
memory



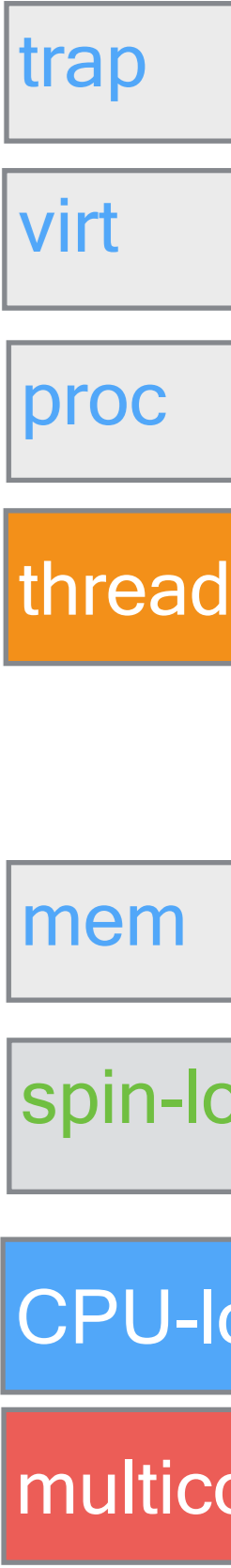
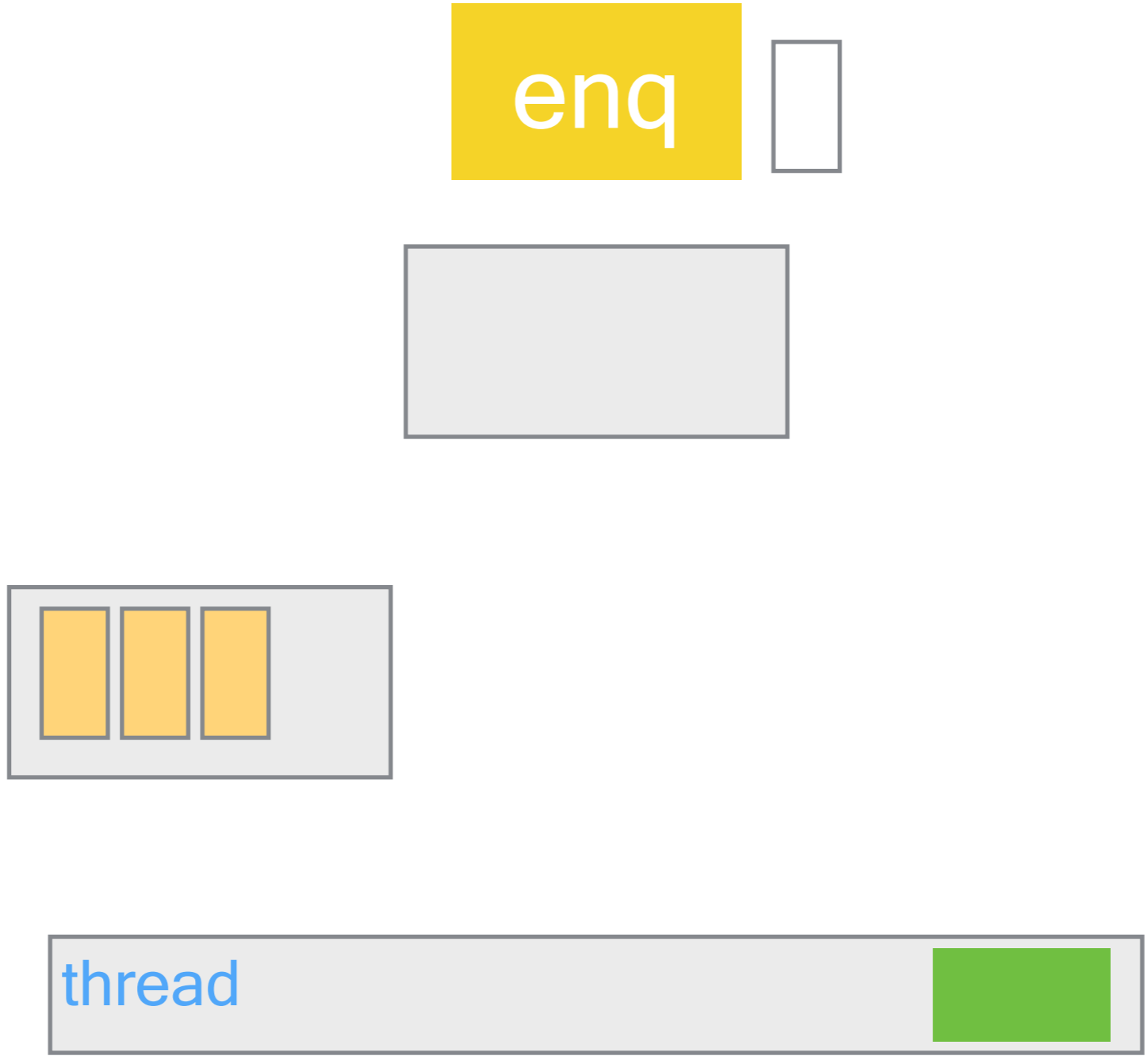
enq





logical
copy

shared
memory



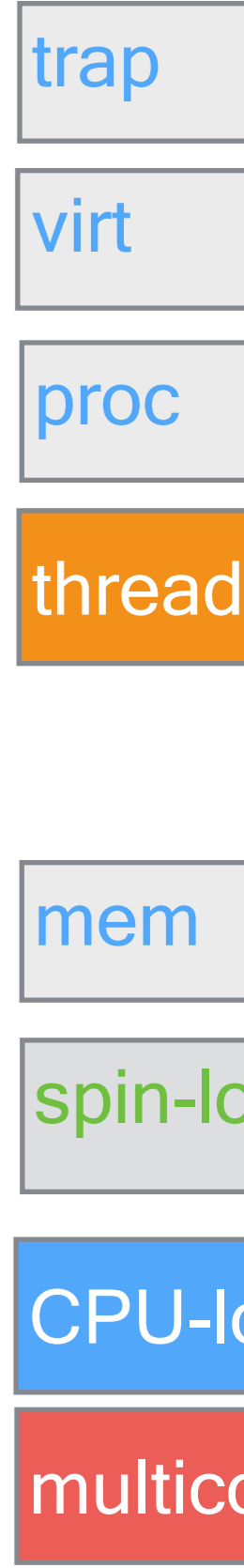
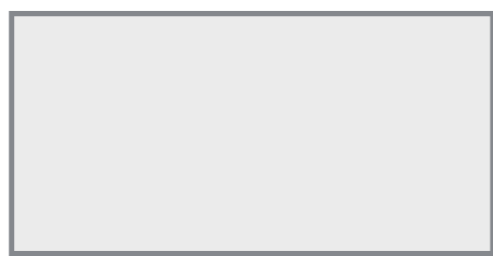


logical copy

shared memory

acq lock

enq



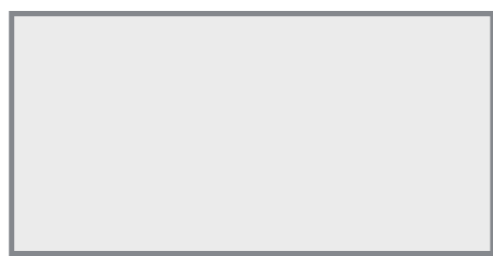


logical copy

shared memory

acq lock

enq



trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico



logical copy

shared memory

acq lock

enq

rel lock



trap

virt

proc

thread

mem

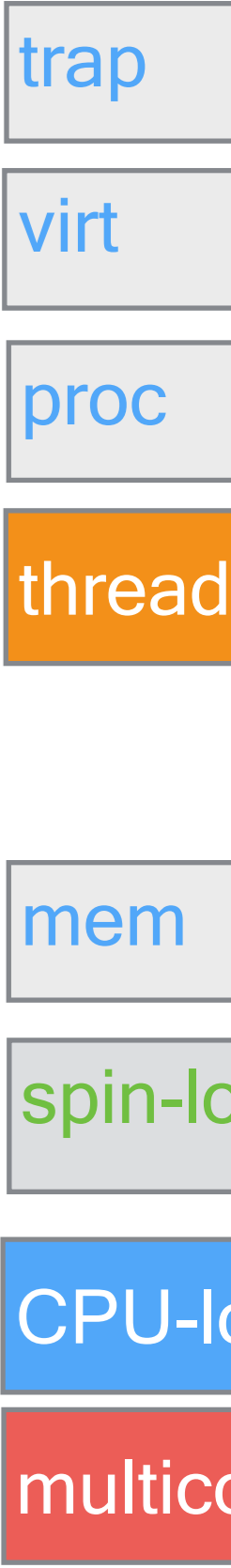
spin-lo

CPU-lo

multico



shared
memory

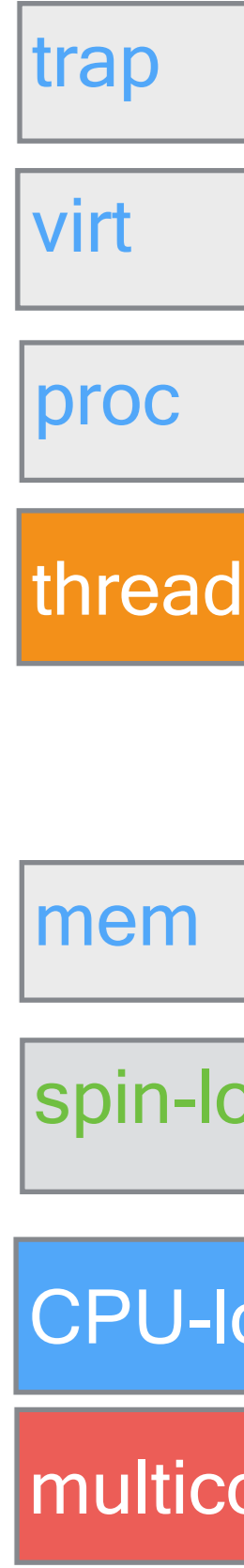




shared
memory



enq





trap

virt

proc

thread

thread

mem

spin-lo

CPU-Id

multico

contributions



```
void yield ()  
{  
  uint t = tid();  
  ...  
  ▶ enq (t, rdq());  
  uint s = ▶ deq (rdq());  
  ...  
  context_switch (t, s)  
}
```

thread-local machine

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

contributions



asm&C
CompcertX

```
void yield ()  
{  
  uint t = tid();  
  ...  
  ▶ enq (t, rdq());  
  uint s = ▶ deq (rdq());  
  context_switch s)  
}
```

thread-local machine

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico

contribution

software scheduler



mix of 3

yield

sleep

wakeup

thread-local machine

trap

virt

proc

thread

mem

spin-lo

CPU-lo

multico



trap

virt

proc

thread

thread

mem

spin-lo

CPU-lo

multico



proc

CV

IPC

multico

CPU-Id

spin-lo

mem

thread

thread

virt

trap

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

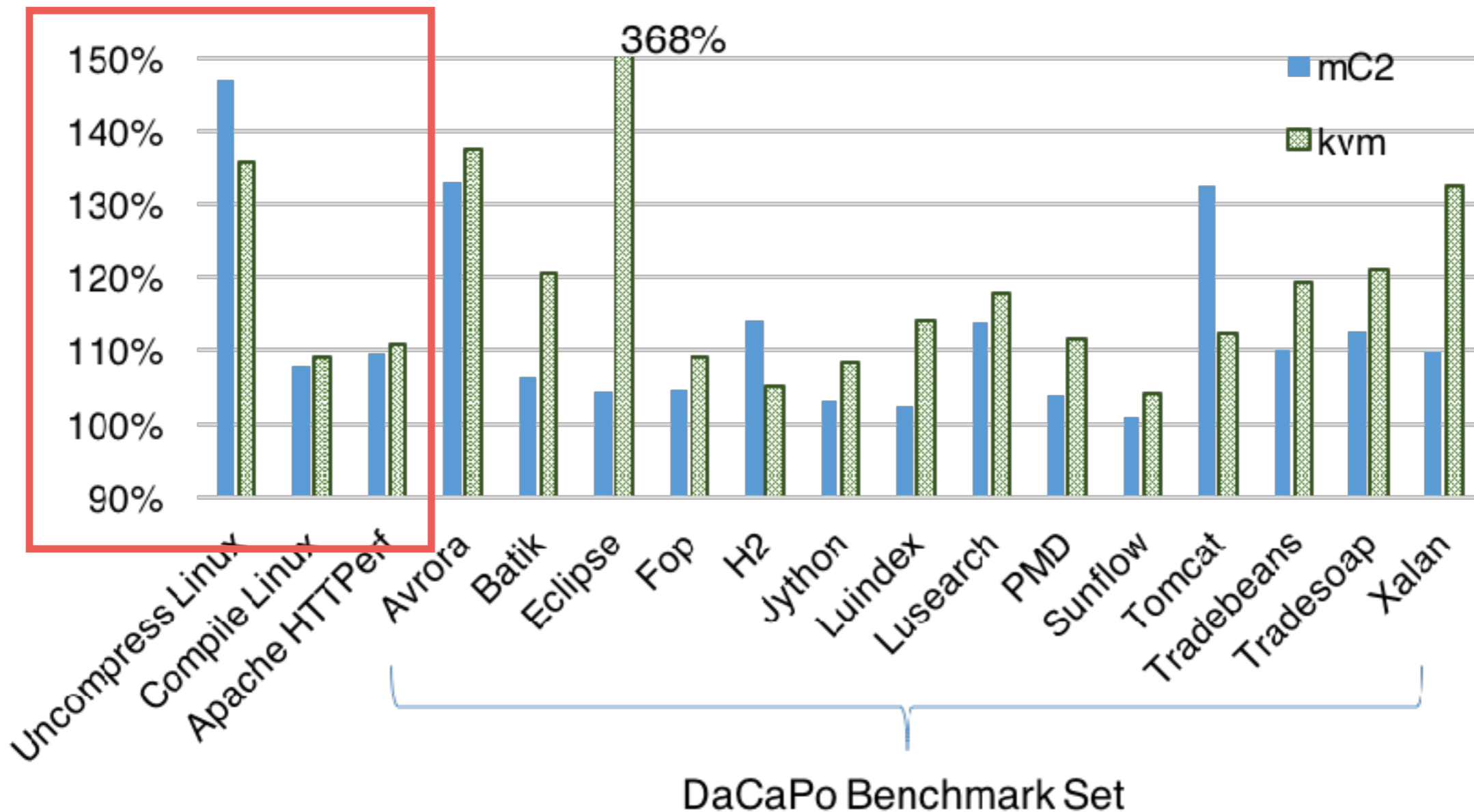
spin-lo

CPU-lo

multico

evaluation: performance

mC2 is comparable with **kvm**



trap

virt

proc

thread

thread

mem

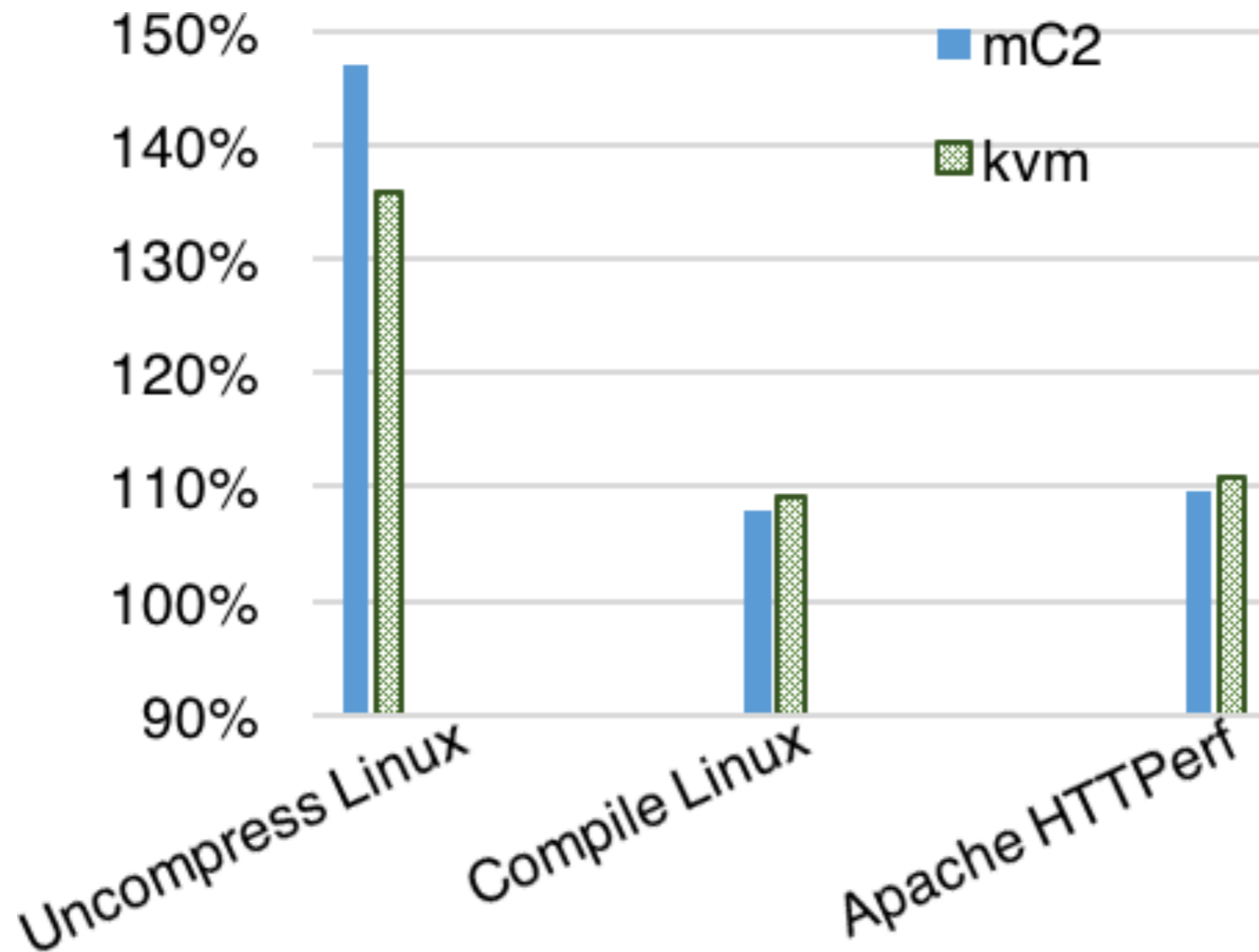
spin-lo

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

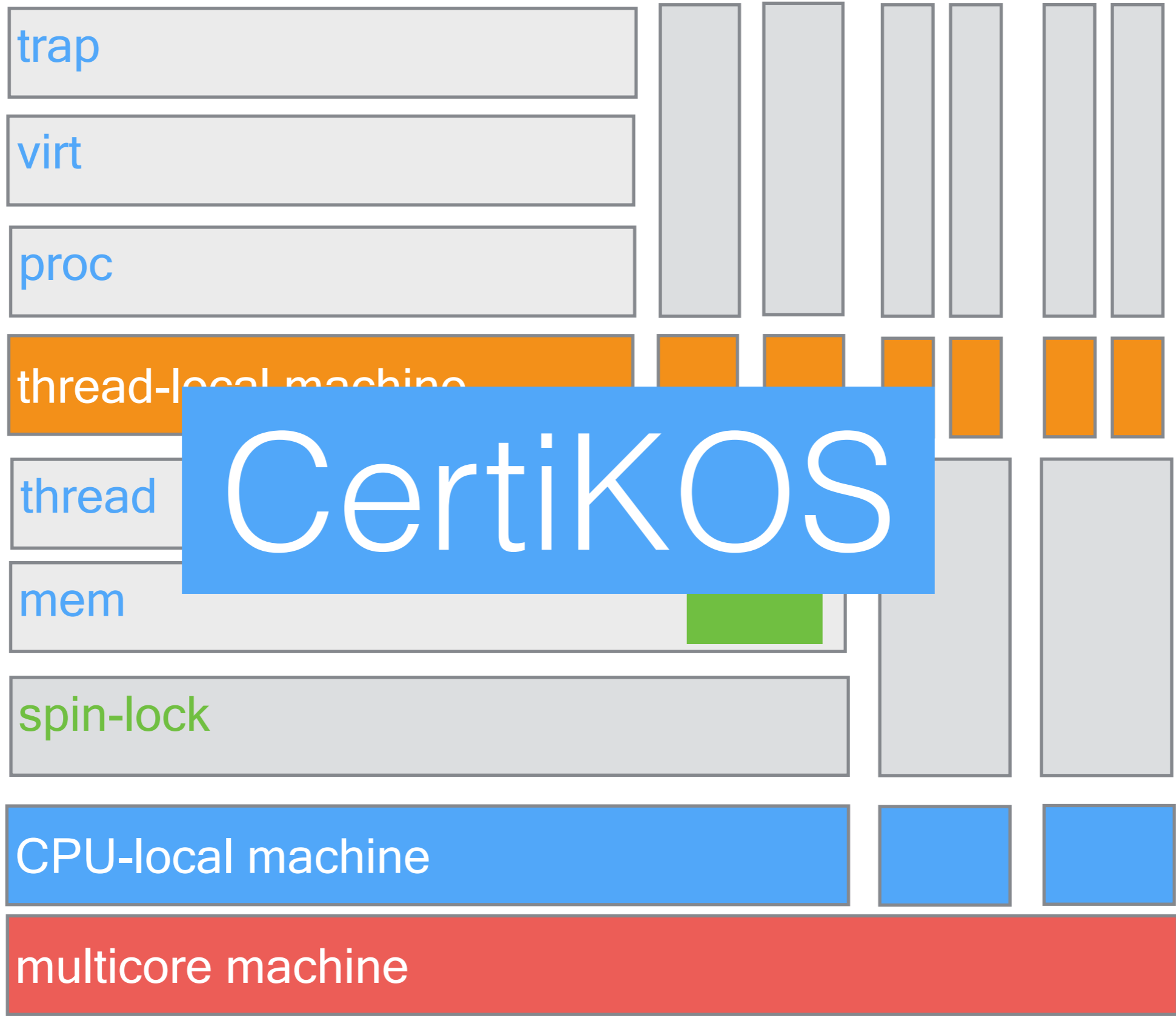
thread

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

mC2

the first formally verified
concurrent OS kernel.

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