



Enabling Performance on Modern Hardware with Practical Verification

Diogo Behrens



- ▶ Modern Hardware and Concurrent Programs
 - ▶ Enabling Performance with Practical Verification
 - ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
 - ▶ Wrap up and Outlook
- 

- ▶ Modern Hardware and Concurrent Programs
 - ▶ Enabling Performance with Practical Verification
 - ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
 - ▶ Wrap up and Outlook
- 

Concurrency is everywhere

Operating systems, databases and server applications resort to **multicore concurrency** to achieve high performance.



MariaDB



PostgreSQL

Concurrency is everywhere

Operating systems, databases and server applications resort to **multicore concurrency** to achieve high performance.

... but brings challenges!

- ▶ **Notoriously hard to get right**
 - ▶ Non-deterministic thread interleaving



MariaDB



Concurrency is everywhere

Operating systems, databases and server applications resort to **multicore concurrency** to achieve high performance.



MariaDB



... but brings challenges!

- ▶ **Notoriously hard to get right**
 - ▶ Non-deterministic thread interleaving
- ▶ **Traditional testing isn't sufficient**
 - ▶ Testing is **unlikely** to trigger subtle concurrency bugs
 - ▶ Production is **likely** to trigger them,
→ many instances and long executions
- ▶ **Reproducing concurrency bugs is hard**

Relaxed memory models (RMMs)

- ▶ Modern architectures becoming popular
eg, Arm, RISC-V

Relaxed memory models (RMMs)

- ▶ Modern architectures becoming popular
eg, Arm, RISC-V



Modern hardware makes matters worse

Relaxed memory models (RMMs)

- ▶ Modern architectures becoming popular
eg, Arm, RISC-V

Huawei Unveils Industry's Highest-Performance ARM-based CPU

Bringing

[Shenzhen, China, January 7, 2019] Today, Huawei unveiled its new ARM-based CPU, called Huawei Kunpeng 920, the new generation of ARM-based CPUs. It is designed for high performance, storage, and ARM-native application scenarios, and is a key component of the open, collaborative, and win-win ecosystem, to



AWS Graviton Processor

Enabling the best price performance in Amazon EC2

[Get Started with AWS Graviton-based EC2 Instances](#)

AWS Graviton processors are custom built by Amazon Web Services using 64-bit Arm Neoverse cores to deliver the best price performance for your cloud workloads running in Amazon EC2. Amazon EC2 provides the broadest and deepest portfolio of compute instances, including many that are powered by latest-generation Intel and AMD processors. AWS Graviton processors add even more choice to help customers optimize performance and cost for their workloads.

The first-generation AWS Graviton processors power Amazon EC2 A1 instances, the first ever Arm-based instances on AWS. These instances deliver significant cost savings over other general-purpose instances for scale-out applications such as web servers, containerized microservices, data/log processing, and other workloads that can run on smaller cores and fit within the available memory footprint.

Relaxed memory models (RMMs)

- Modern architectures becoming popular
eg, Arm, RISC-V

Huawei Unveils Industry's Highest-Performance ARM-based CPU

Bringing

[Shenzhen, China, January 7, 2019] Today, Huawei unveiled its new ARM-based CPU. Called Huawei Kunpeng 920, the new

AWS Graviton Processor

Enabling the best price performance in Amazon EC2

[Get Started with AWS Graviton-based EC2 Instances](#)

techcrunch.com

Microsoft launches the ARM-based Surface Pro X – TechCrunch

Zack Whittaker

3-3 minutes

At its annual Surface hardware event, [Microsoft](#) today announced the long-rumored ARM-based Surface, the first time Microsoft itself has launched a device with an ARM-based processor inside. The 13-inch device will use Microsoft's own custom SQ1 chip, based on [Qualcomm's](#) Snapdragon and an AI accelerator, making it the first Surface with an integrated AI engine. Microsoft and Qualcomm also worked on building custom-designed GPU cores for the Pro X, which will run Microsoft's version of Windows 10 for ARM.

The Pro X will be available on November 5, starting at \$999, and is now available for pre-order.

Web Services using 64-bit Arm Neoverse cores to power AI and ML workloads running in Amazon EC2. Amazon EC2 provides a range of instance types, including many that are powered by latest-generation processors add even more choice to help customers

Amazon EC2 A1 instances, the first ever Arm-based EC2 instances, that offer significant savings over other general-purpose instances for a wide range of microservices, data/log processing, and other workloads with a small memory footprint.

Modern hardware makes matters worse

Relaxed memory models (RMMs)

- ▶ Modern architectures becoming popular
eg, Arm, RISC-V

Huawei Unveils Industry's Highest-Performance ARM-based CPU

Bringing

[Shenzhen, China, January 7, 2019] Today, Huawei unveiled its latest ARM-based CPU. Called Huawei Kirin980, the new

AWS Graviton Processor

Enabling the best price performance in Amazon EC2

[Get Started with AWS Graviton-based EC2 Instances](#)

techcrunch.com

Microsoft launches the ARM-based Surface

Zack Whittaker

3-3 minutes

At its annual Surface hardware event, [Microsoft](#) today announced for the first time Microsoft itself has launched a device with an ARM-based processor. Microsoft's own custom SQ1 chip, based on [Qualcomm's](#) Snapdragon 8cx, will power the first Surface with an integrated AI engine. Microsoft and Qualcomm GPU cores for the Pro X, which will run Microsoft's version of Windows.

The Pro X will be available on November 5, starting at \$999, and is

Apple Silicon: M1 MacBook Pro, MacBook Air and Mac mini Now Available

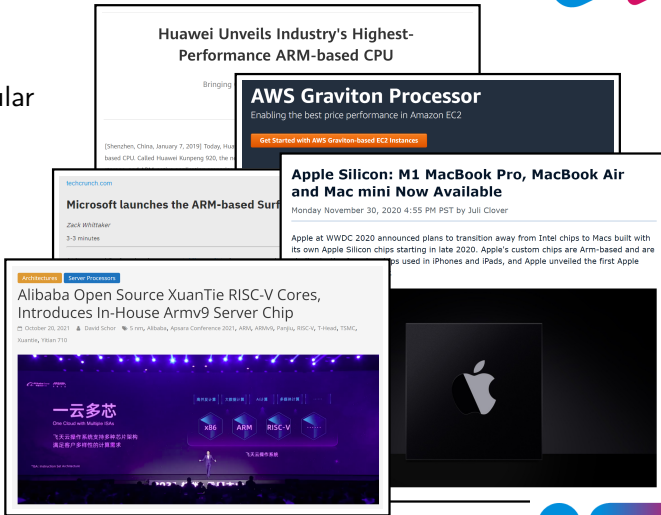
Monday November 30, 2020 4:55 PM PST by Juli Clover

Apple at WWDC 2020 announced plans to transition away from Intel chips to Macs built with its own Apple Silicon chips starting in late 2020. Apple's custom chips are Arm-based and are similar to the A-series chips used in iPhones and iPads, and Apple unveiled the first Apple Silicon Macs in November.



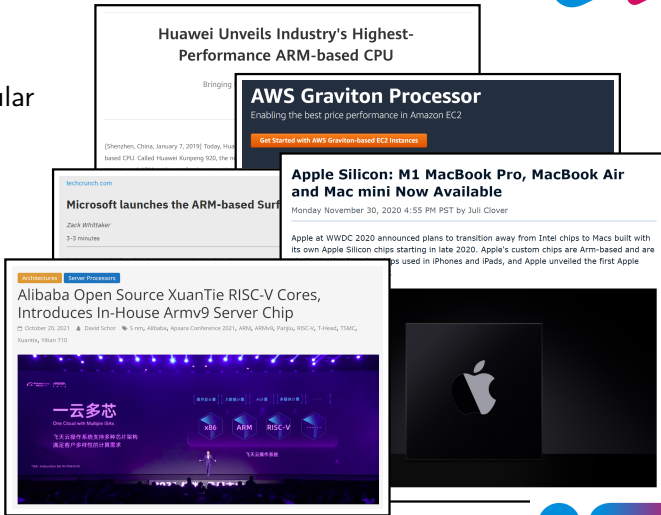
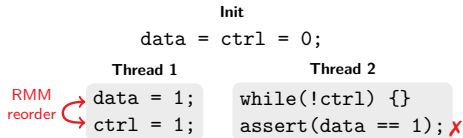
Relaxed memory models (RMMs)

- Modern architectures becoming popular
eg, Arm, RISC-V



Relaxed memory models (RMMs)

- ▶ Modern architectures becoming popular
eg, Arm, RISC-V
- ▶ **Aggressive reorderings** to improve performance
- ▶ **Much higher non-determinism**
- ▶ Careful use of **memory barriers**
(neither too many, nor too few)



Do RMM bugs really happen?

Bugs found with our tools



Potential hang: DPDK MCS lock

The lock had one missing release barrier. An Arm engineer replied to our patch: “Unfortunately, memory ordering questions are hard topics. I have been discussing this internally [...], hope to conclude soon.” More than 3 months to accept the 1-line patch.



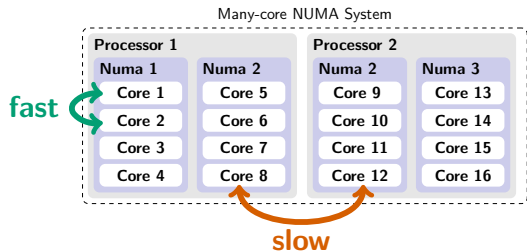
Crash: MariaDB lockfree hashtable

Due to missing barriers for Arm, data of a node can be accessed after the node has been deleted. That causes SEGFAULT crashes on high load workloads.

Modern hardware makes matters worse



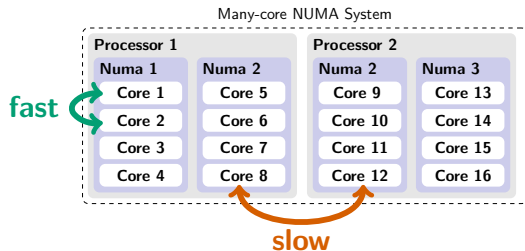
Modern hardware makes matters worse



Deep NUMA hierarchies

- Core **distance affects** shared-memory communication **performance**

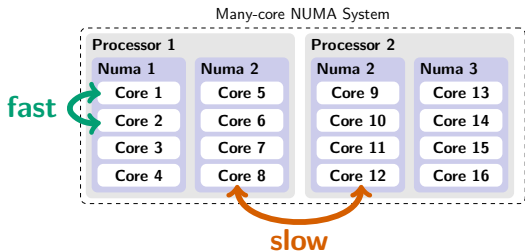
Modern hardware makes matters worse



Deep NUMA hierarchies

- ▶ Core **distance affects** shared-memory communication **performance**
- ▶ Concurrent algorithms must exploit that!

Modern hardware makes matters worse

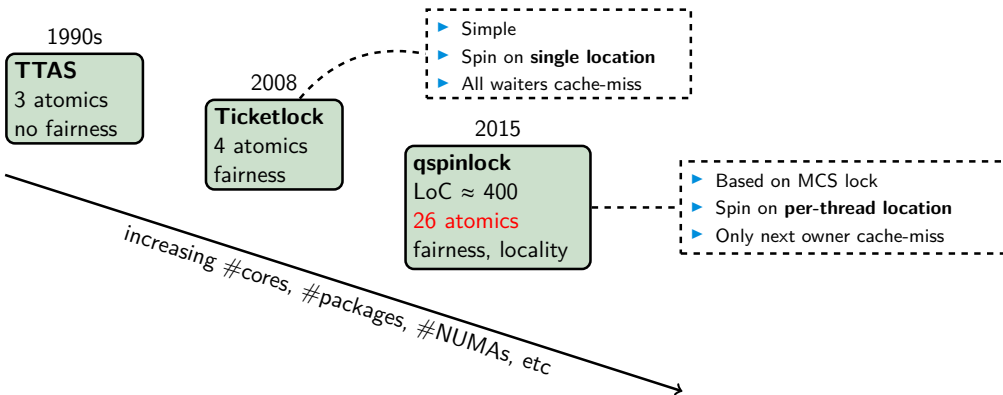


Deep NUMA hierarchies

- ▶ Core **distance affects** shared-memory communication **performance**
- ▶ Concurrent algorithms must exploit that!
- ▶ **Makes code even more complex**

As the hardware evolves, so does the concurrency control

Example of Linux spinlock over the years



"atomics" refers to racy accesses, ie, variables concurrently accessed by multiple cores

As the hardware evolves, so does the concurrency control

Example of Linux spinlock over the years

1990s
TTAS
3 atomics
no fairness

2008
Ticketlock
4 atomics
fairness

2015

qspinlock
LoC \approx 400
26 atomics
fairness, locality

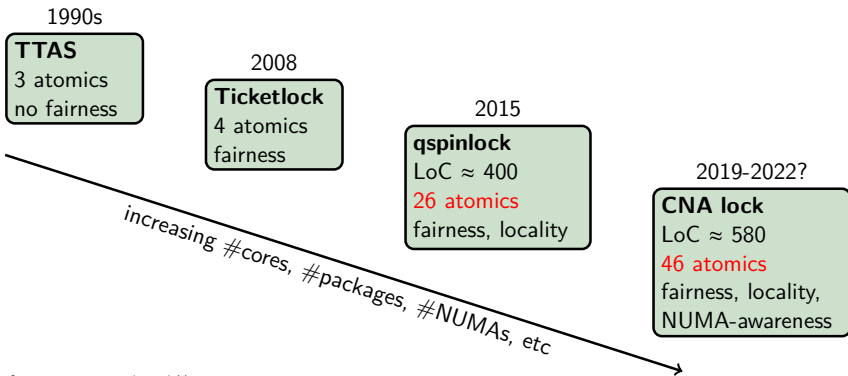
Version	Date	Correctness
Linux 4.4	2015/09/11	Not verified
Linux 4.5	2015/11/09	RMM bug (fixed in 4.16)
Linux 4.8	2016/06/03	RMM bug (fixed in 4.16)
Linux 4.16	2018/02/13	Not verified
Linux 5.6	2020/01/07	Not verified

increasing #cores, #packages, #NUMAs, etc

"atomics" refers to racy accesses, ie, variables concurrently accessed by multiple cores

As the hardware evolves, so does the concurrency control

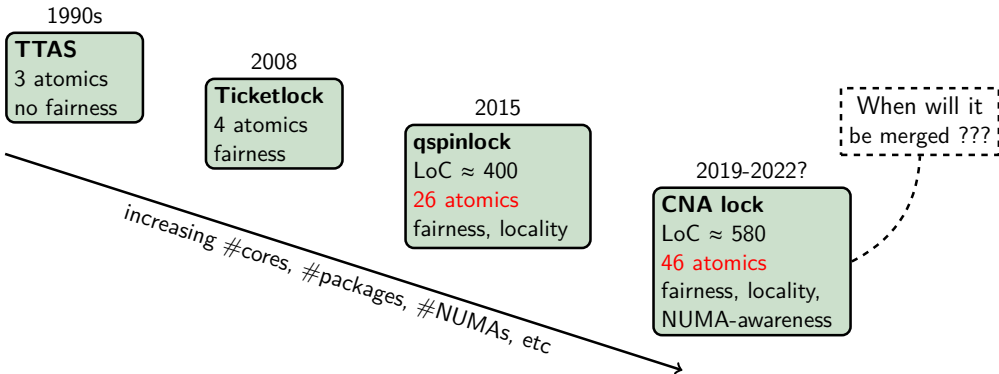
Example of Linux spinlock over the years



"atomics" refers to racy accesses, ie, variables concurrently accessed by multiple cores

As the hardware evolves, so does the concurrency control

Example of Linux spinlock over the years



"atomics" refers to racy accesses, ie, variables concurrently accessed by multiple cores

What can the industry do?



Solution 1: overprotection

- ▶ Spray the code with memory barriers
- ▶ Simplify design as much as possible
- ▶ Cost: **performance impact**



Solution 2: expert design

- ▶ Hire highly-skilled engineers
- ▶ Carefully design and implement
- ▶ Cost: **error-prone, low maintainability**

Solution 3:
Enabling Performance with **Practical** Verification

What is **Practical** Verification?



In short: we don't know!

What is **Practical** Verification?

We know what formal verification is!

- ▶ Goal: **show correctness**
- ▶ Examples:
 - ▶ interactive theorem proving
 - ▶ model checking

We know what formal verification is!

- ▶ Goal: **show correctness**
- ▶ Examples:
 - ▶ interactive theorem proving
 - ▶ model checking
- ▶ In contrast to testing, **exhaustive**:
 - ▶ All possible interleavings
 - ▶ All possible reorderings
 - ▶ All possible contexts, ...

We know what formal verification is!

- ▶ Goal: **show correctness**
- ▶ Examples:
 - ▶ interactive theorem proving
 - ▶ model checking
- ▶ In contrast to testing, **exhaustive**:
 - ▶ All possible interleavings
 - ▶ All possible reorderings
 - ▶ All possible contexts, ...
- ▶ Problems:
 - ▶ **High level of expertise** required
 - ▶ Often **unrealistic assumptions** (eg, hardware simplifications)

What is Practical Verification?

We know what formal verification is!

- ▶ Goal: **show correctness**
- ▶ Examples:
 - ▶ interactive theorem proving
 - ▶ model checking
- ▶ In contrast to testing, **exhaustive**:
 - ▶ All possible interleavings
 - ▶ All possible reorderings
 - ▶ All possible contexts, ...
- ▶ Problems:
 - ▶ **High level of expertise** required
 - ▶ Often **unrealistic assumptions** (eg, hardware simplifications)

And practical verification?

- ▶ Goal: **increase confidence**
- ▶ Exhaustive as long as **it pays off**
- ▶ A concurrency counterpart of testing
- ▶ **Ideally any developer** can use it:
 - ▶ Push-button, fully-automated
 - ▶ No expertise required

What is Practical Verification?

We know what formal verification is!

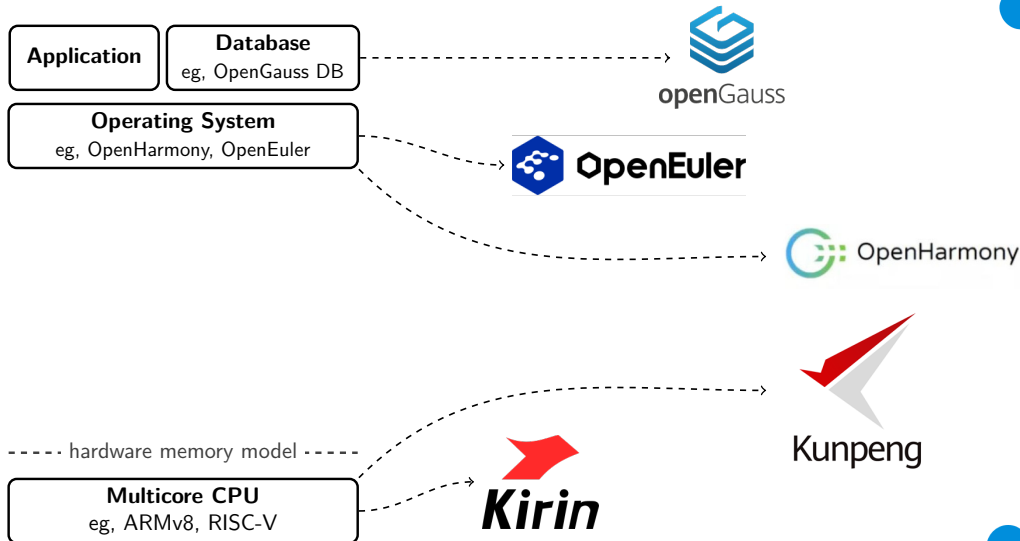
- ▶ Goal: **show correctness**
- ▶ Examples:
 - ▶ interactive theorem proving
 - ▶ model checking
- ▶ In contrast to testing, **exhaustive**:
 - ▶ All possible interleavings
 - ▶ All possible reorderings
 - ▶ All possible contexts, ...
- ▶ Problems:
 - ▶ **High level of expertise** required
 - ▶ Often **unrealistic assumptions** (eg, hardware simplifications)

And practical verification?

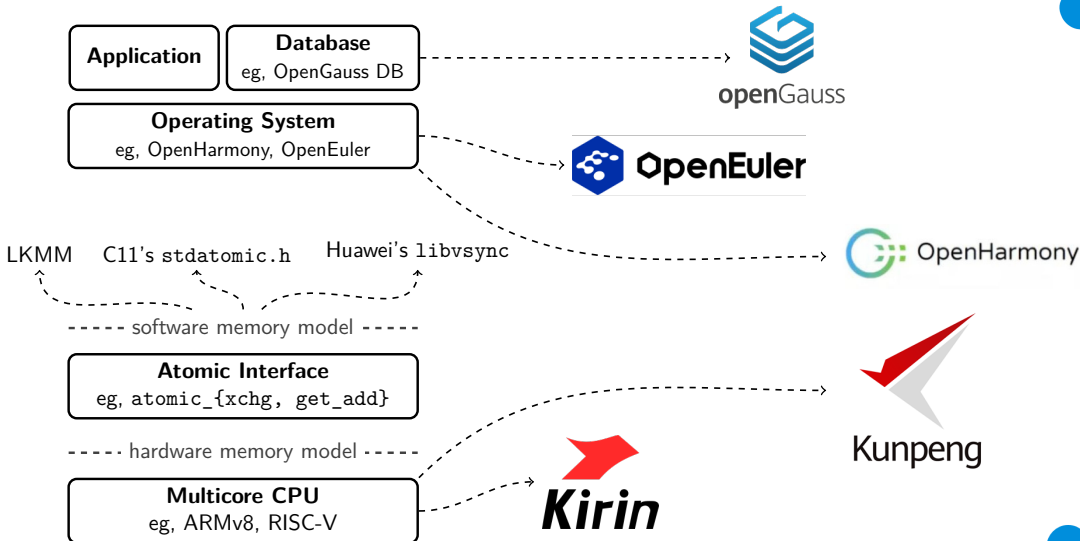
- ▶ Goal: **increase confidence**
- ▶ Exhaustive as long as **it pays off**
- ▶ A concurrency counterpart of testing
- ▶ **Ideally any developer** can use it:
 - ▶ Push-button, fully-automated
 - ▶ No expertise required
- ▶ **In practice, a few experts** must help
 - ▶ Defining what to check (properties)
 - ▶ Improving tool **scalability**
 - ▶ Making methods **more realistic**

Where do we apply practical verification?

Practical Verification Scope in Huawei



Practical Verification Scope in Huawei



Application**Database**

eg, OpenGauss DB

Operating System

eg, OpenHarmony, OpenEuler

**Synchronization
Primitives**eg, spinlock,
mutex, RCU

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Synchronous with locks

```
run_thread() {  
    // some parallel work  
    initialize_data();  
    do_something_locally();  
  
    // critical section  
    lock_acquire(&lock);  
    sequential_code();  
    lock_release(&lock);  
  
    // some more parallel work  
    do_something_locally();  
}
```

Practical Verification Scope in Huawei

Application

Database

eg, OpenGauss DB

Operating System

eg, OpenHarmony, OpenEuler

Synchronization Primitives

eg, spinlock, mutex, RCU

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Synchronous with locks

```
run_thread() {  
    // some parallel work  
    initialize_data();  
    do_something_locally();  
  
    // critical section  
    lock_acquire(&lock);  
    sequential_code();  
    lock_release(&lock);  
  
    // some more parallel work  
    do_something_locally();  
}
```



DPDK
DATA PLANE DEVELOPMENT KIT



Security. Performance. Proof.

Practical Verification Scope in Huawei

Application

Database
eg, OpenGauss DB

Operating System

eg, OpenHarmony, OpenEuler

Synchronization Primitives

eg, spinlock, mutex, RCU

Lockless Data Structures

eg, list, queue, hashtable, etc

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Synchronous with locks

```
run_thread() {  
    // some parallel work  
    initialize_data();  
    do_something_locally();  
  
    // critical section  
    lock_acquire(&lock);  
    sequential_code();  
    lock_release(&lock);  
  
    // some more parallel work  
    do_something_locally();  
}
```

Asynchronous with data structures

```
thread1() {  
    do_something_locally();  
    queue_enq(q, data);  
    do_something_else_locally();  
}  
thread2() {  
    data_t *data;  
    data = queue_deq(q);  
    if (data != NULL)  
        do_something_locally(data);  
}
```



DPDK
DATA PLANE DEVELOPMENT KIT



Security. Performance. Proof.

Practical Verification Scope in Huawei

Application

Database
eg, OpenGauss DB

Operating System

eg, OpenHarmony, OpenEuler

Synchronization Primitives

eg, spinlock, mutex, RCU

Lockless Data Structures

eg, list, queue, hashtable, etc

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Synchronous with locks

```
run_thread() {
    // some parallel work
    initialize_data();
    do_something_locally();

    // critical section
    lock_acquire(&lock);
    sequential_code();
    lock_release(&lock);

    // some more parallel work
    do_something_locally();
}
```

Asynchronous with data structures

```
thread1() {
    do_something_locally();
    queue_enq(q, data);
    do_something_else_locally();
}

thread2() {
    data_t *data;
    data = queue_deq(q);
    if (data != NULL)
        do_something_locally(data);
}
```



MariaDB



DPDK
DATA PLANE DEVELOPMENT KIT



Security. Performance. Proof.

Practical Verification Scope in Huawei

Operating System

eg, OpenHarmony, OpenEuler

Synchronization Primitives

eg, spinlock, mutex, RCU

Lockless Data Structures

eg, list, queue, hashtable, etc

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Synchronous with locks

```
run_thread() {  
    // some parallel work  
    initialize_data();  
    do_something_locally();  
  
    // critical section  
    lock_acquire(&lock);  
    sequential_code();  
    lock_release(&lock);  
  
    // some more parallel work  
    do_something_locally();  
}
```

Asynchronous with data structures

```
thread1() {  
    do_something_locally();  
    queue_enq(q, data);  
    do_something_else_locally();  
}  
thread2() {  
    data_t *data;  
    data = queue_deq(q);  
    if (data != NULL)  
        do_something_locally(data);  
}
```

Adhoc synchronization

```
if (atomic_add(counter) > MAX)  
    do_something(data);
```

Practical Verification Scope in Huawei

Operating System

eg, OpenHarmony, OpenEuler

Synchronization Primitives

eg, spinlock, mutex, RCU

Lockless Data Structures

eg, list, queue, hashtable, etc

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Synchronous with locks

```
run_thread() {  
    // some parallel work  
    initialize_data();  
    do_something_locally();  
  
    // critical section  
    lock_acquire(&lock);  
    sequential_code();  
    lock_release(&lock);  
  
    // some more parallel work  
    do_something_locally();  
}
```

Asynchronous with data structures

```
thread1() {  
    do_something_locally();  
    queue_enq(q, data);  
    do_something_else_locally();  
}  
thread2() {  
    data_t *data;  
    data = queue_deq(q);  
    if (data != NULL)  
        do_something_locally(data);  
}
```

Adhoc synchronization

```
if (atomic_add(counter) > MAX)  
    do_something(data);
```

► Sequential code

- It can be tested! Good news!
- Even the critical section code

► Concurrent code

- Cannot be easily tested
- Barrier placement is hard
- Hardware-awareness is hard
- Scope of practical verification!

Synchronous with locks


```
run_thread() {  
    // some parallel work  
    initialize_data();  
    do_something_locally();  
  
    // critical section  
    lock_acquire(&lock);  
    sequential_code();  
    lock_release(&lock);  
  
    // some more parallel work  
    do_something_locally();  
}
```

Asynchronous with data structures

```
thread1() {  
    do_something_locally();  
    queue_enq(q, data);  
    do_something_else_locally();  
}  
thread2() {  
    data_t *data;  
    data = queue_deq(q);  
    if (data != NULL)  
        do_something_locally(data);  
}
```

Adhoc synchronization

```
if (atomic_add(counter) > MAX)  
    do_something(data);
```


- ▶ Modern Hardware and Concurrent Programs
 - ▶ Enabling Performance with Practical Verification
 - ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
 - ▶ Wrap up and Outlook
- 

- ▶ Modern Hardware and Concurrent Programs
- ▶ Enabling Performance with Practical Verification
- ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
- ▶ Wrap up and Outlook

Application**Database**
eg, OpenGauss DB**Operating System**
eg, OpenHarmony, OpenEuler**Synchronization
Primitives**
eg, spinlock,
mutex, RCU**Lockless Data
Structures**
eg, list, queue,
hashtable, etc

----- software memory model -----

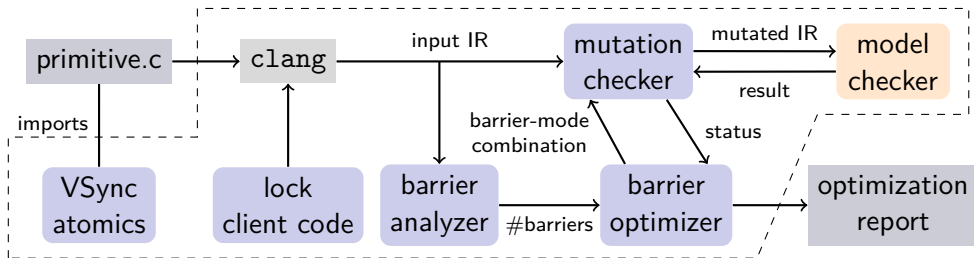
Atomic Interface
eg, atomic_{xchg, get_add}

----- hardware memory model -----

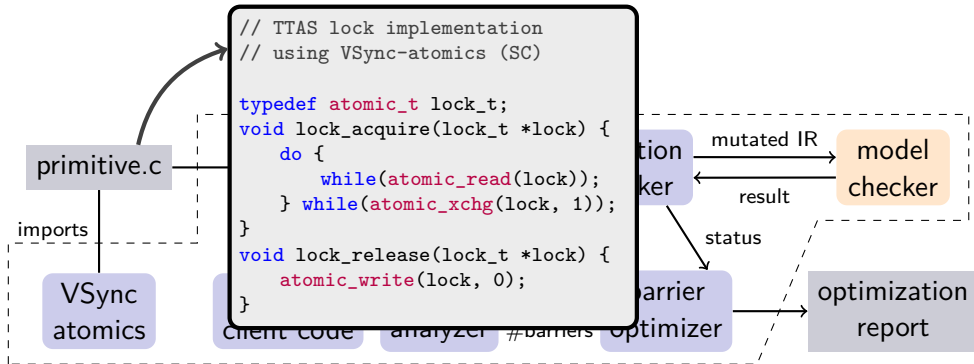
Multicore CPU
eg, ARMv8, RISC-V

VSync: Push-button Verification/Optimization on RMMs

- VSync: <https://dl.acm.org/doi/abs/10.1145/3445814.3446748> — **Best paper** @ ASPLOS'21
by Oberhauser, Chehab, Behrens, Fu, Paolillo, Oberhauser, Bhat, Wen, Chen, Kim, Vafeiadis
- Making relaxed memory models fair: <https://dl.acm.org/doi/abs/10.1145/3485475> — **Best paper** @ OOPSLA'21
by Lahav, Namakonov, Oberhauser, Podkopaev, Vafeiadis



VSync: Push-button Verification/Optimization on RMMs

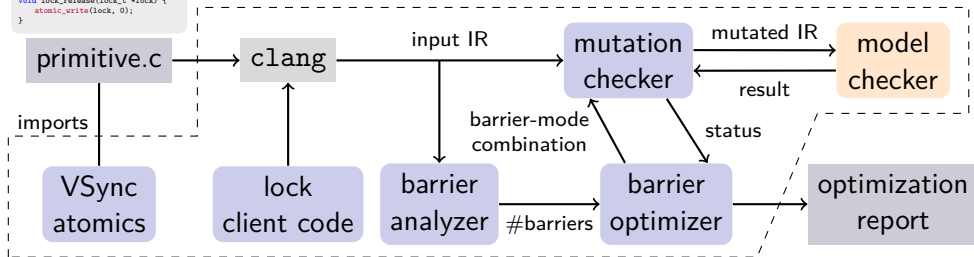


VSync: Push-button Verification/Optimization on RMMs



```
// TTAS lock implementation
// using VSync-atomics (SC)

typedef atomic_t lock_t;
void lock_acquire(lock_t *lock) {
    do {
        while(atomic_read(lock));
    } while(atomic_xchg(lock, 1));
}
void lock_release(lock_t *lock) {
    atomic_write(lock, 0);
}
```



VSync: Push-button Verification/Optimization on RMMs

```
// TTAS lock implementation
// using VSync-atomics (SC)

typedef atomic_t lock_t;
void lock_acquire(lock_t *lock) {
    do {
        while(atomic_read(lock));
    } while(atomic_xchg(lock, 1));
}
void lock_release(lock_t *lock) {
    atomic_write(lock, 0);
}
```

primitive.c

imports

VSync
atomics

lock
client code

barrier
analyzer

#barriers

barrier
optimizer

optimization
report

VSYNC-atomics

atomic_xchg
atomic_xchg_rel
atomic_xchg_acq
atomic_xchg_rlx
atomic_read
atomic_read_acq
atomic_read_rlx
atomic_write
atomic_write_rel
atomic_write_rlx
atomic_fence
atomic_fence_acq
atomic_fence_rel
atomic_fence_rlx

RISC-V

amoswap.w.aq.rl.sc
amoswap.w.rl
amoswap.w.aq
amoswap.w
fence [rw,rw]; lw; fence [r,rw]
lw; fence [r,rw]
lw
fence [rw,w]; sw; fence [rw,rw]
fence [rw,w]; sw
sw
fence [rw,rw]
fence [r,rw]
fence [rw,w]
nop

ARMv8

LDAXR;STLXR
LDXR;STLXR
LDAXR;STXR
LDXR;STXR
LDAR
LDAR
LDR
STLR
STLR
STR
DMB ISH
DMB ISHL
DMB ISH
NOP

x86

XCHG
XCHG
XCHG
XCHG
MOV
MOV
MOV
MOV;MFENCE
MOV
MOV
MOV
MFENCE
NOP
NOP
NOP







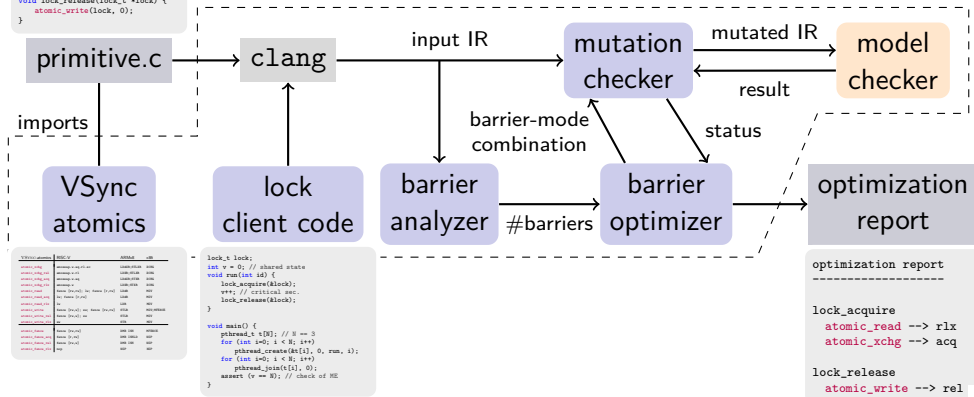


VSync: Push-button Verification/Optimization on RMMs

Problem: Some optimizations cause hangs on Arm CPUs! Why?

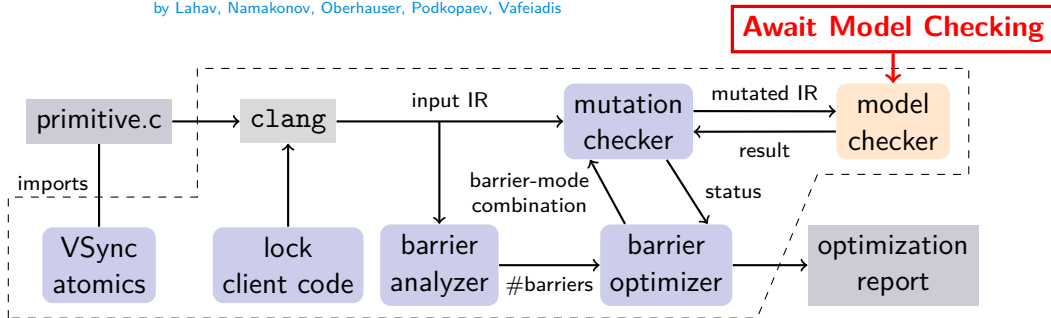
```
// TTAS lock implementation
// using VSync-atomics (SC)

typedef atomic_t lock_t;
void lock_acquire(lock_t *lock) {
    do {
        while(atomic_read(lock));
        while(atomic_xchg(lock, 1));
    }
}
void lock_release(lock_t *lock) {
    atomic_write(lock, 0);
}
```



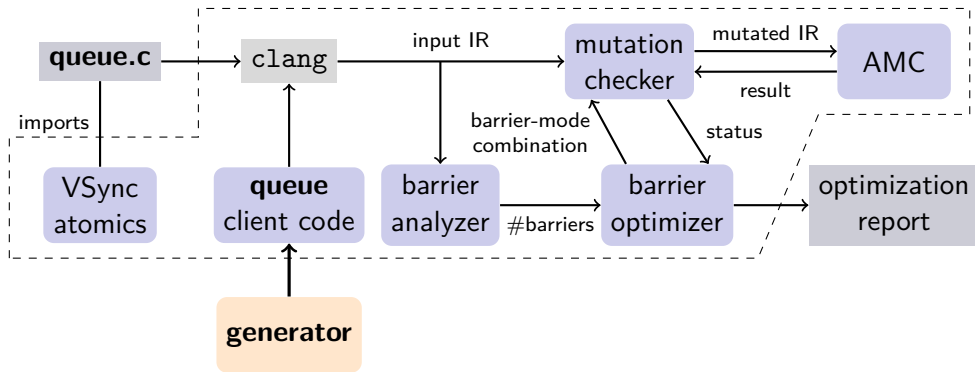
VSync: Push-button Verification/Optimization on RMMs

- VSync: <https://dl.acm.org/doi/abs/10.1145/3445814.3446748> — **Best paper @ ASPLOS'21**
by Oberhauser, Chehab, Behrens, Fu, Paolillo, Oberhauser, Bhat, Wen, Chen, Kim, Vafeiadis
- Making relaxed memory models fair: <https://dl.acm.org/doi/abs/10.1145/3485475> — **Best paper @ OOPSLA'21**
by Lahav, Namakonov, Oberhauser, Podkopaev, Vafeiadis



No scalable model checker for RMM was capable of checking termination of spinloops!

Current work: extend VSync to data structures
eg, queues, lists, stacks



- ▶ Modern Hardware and Concurrent Programs
- ▶ Enabling Performance with Practical Verification
- ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
- ▶ Wrap up and Outlook

Application**Database**

eg, OpenGauss DB

Operating System

eg, OpenHarmony, OpenEuler

**Synchronization
Primitives**eg, spinlock,
mutex, RCU**Lockless Data
Structures**eg, list, queue,
hashtable, etc

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

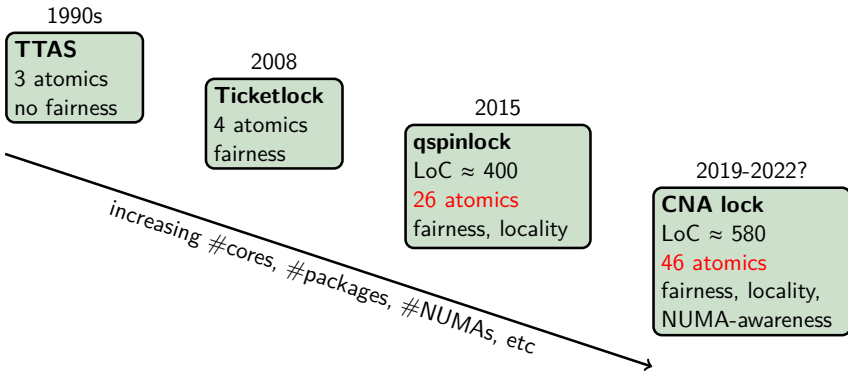
----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Verifying CNA-based Linux qspinlock on LKMM

LKMM = Linux Kernel Memory Model



"atomics" refers to racy accesses, ie, variables concurrently accessed by multiple cores

Verifying CNA-based Linux qspinlock on LKMM

Technical report @ arXiv Jul'22 — by Paolillo, Ponce de León, Haas, Behrens, Chehab, Fu, Meyer
<https://arxiv.org/abs/2111.15240>

- ▶ **Goal:**
 - ▶ Support merge of CNA patch (Oracle)
- ▶ **Verification:**
 - ▶ Linux Kernel Memory Model (LKMM)
 - ▶ Armv8 memory model
 - ▶ Power memory model
- ▶ **Verification time:**
 - ▶ With 4 threads, from 12h to 15h

Verifying CNA-based Linux qspinlock on LKMM

Technical report @ arXiv Jul'22 — by Paolillo, Ponce de León, Haas, Behrens, Chehab, Fu, Meyer
<https://arxiv.org/abs/2111.15240>

- ▶ **Goal:**
 - ▶ Support merge of CNA patch (Oracle)
- ▶ **Verification:**
 - ▶ Linux Kernel Memory Model (LKMM)
 - ▶ Armv8 memory model
 - ▶ Power memory model
- ▶ **Verification time:**
 - ▶ With 4 threads, from 12h to 15h
- ▶ **Linux qspinlock and CNA:**
 - ▶ **correct on Armv8**
 - ▶ correct on Power

Verifying CNA-based Linux qspinlock on LKMM

Technical report @ arXiv Jul'22 — by Paolillo, Ponce de León, Haas, Behrens, Chehab, Fu, Meyer
<https://arxiv.org/abs/2111.15240>

- ▶ **Goal:**
 - ▶ Support merge of CNA patch (Oracle)
- ▶ **Verification:**
 - ▶ Linux Kernel Memory Model (LKMM)
 - ▶ Armv8 memory model
 - ▶ Power memory model
- ▶ **Verification time:**
 - ▶ With 4 threads, from 12h to 15h
- ▶ **Linux qspinlock and CNA:**
 - ▶ **correct on Armv8**
 - ▶ correct on Power
 - ▶ **incorrect on LKMM!**
- ▶ **Lessons learned**
 - ▶ Was the model checker broken? No!
 - ▶ LKMM mismatches reality!
(**expert support**)

Verifying CNA-based Linux qspinlock on LKMM

Technical report @ arXiv Jul'22 — by Paolillo, Ponce de León, Haas, Behrens, Chehab, Fu, Meyer
<https://arxiv.org/abs/2111.15240>

- ▶ **Goal:**
 - ▶ Support merge of CNA patch (Oracle)
- ▶ **Verification:**
 - ▶ Linux Kernel Memory Model (LKMM)
 - ▶ Armv8 memory model
 - ▶ Power memory model
- ▶ **Verification time:**
 - ▶ With 4 threads, from 12h to 15h
- ▶ **Linux qspinlock and CNA:**
 - ▶ **correct on Armv8**
 - ▶ correct on Power
 - ▶ **incorrect on LKMM!**
- ▶ **Lessons learned**
 - ▶ Was the model checker broken? No!
 - ▶ LKMM mismatches reality!
(**expert support**)
 - ▶ Having 2 model checkers is helpful!

<https://github.com/hernanponcedeleon/Dat3M>
<https://github.com/MPI-SWS/genmc>

- ▶ Modern Hardware and Concurrent Programs
- ▶ Enabling Performance with Practical Verification
- ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
- ▶ Wrap up and Outlook

Application**Database**

eg, OpenGauss DB

Operating System

eg, OpenHarmony, OpenEuler

**Synchronization
Primitives**eg, spinlock,
mutex, RCU**Lockless Data
Structures**eg, list, queue,
hashtable, etc

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

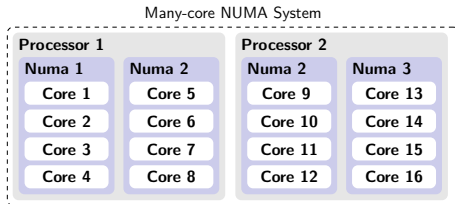
----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

Why would locks be too complex to model check?

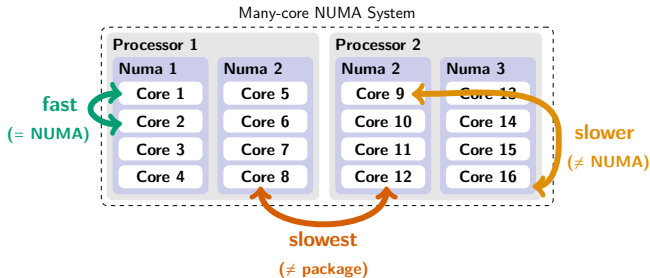
Modern lock-design challenges



Why would locks be too complex to model check?

Modern lock-design challenges

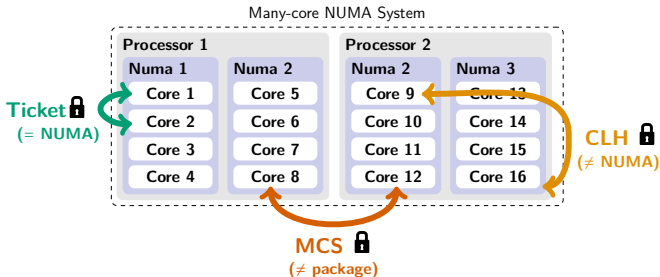
- **Deep NUMA hierarchies**
eg, packages, NUMA nodes, L3 cache partitions



Why would locks be too complex to model check?

Modern lock-design challenges

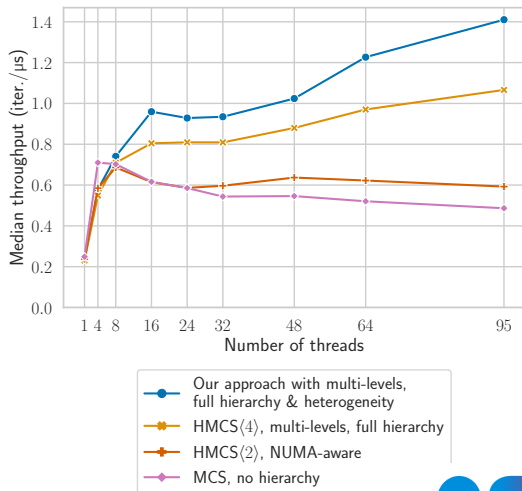
- **Deep NUMA hierarchies**
eg, packages, NUMA nodes, L3 cache partitions
- **Lock heterogeneity**
Different locks perform better in different contexts, eg, cores with shared cache or not



Why would locks be too complex to model check?

What are the potential improvements?

LevelDB benchmark, x86 server, 96 hyperthreads

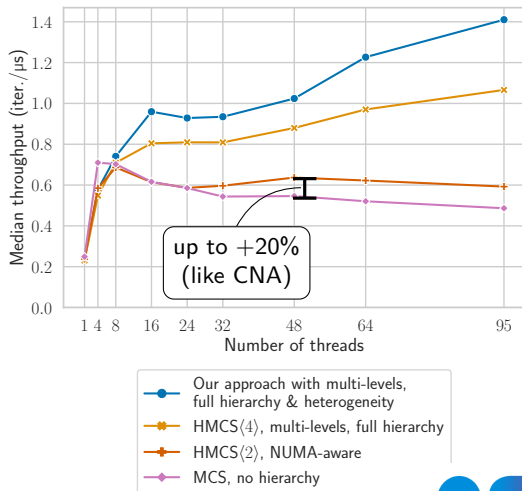


Why would locks be too complex to model check?

What are the potential improvements?

- 2-level is good

LevelDB benchmark, x86 server, 96 hyperthreads

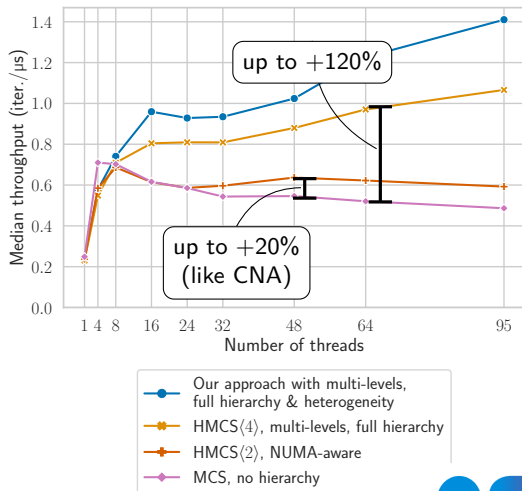


Why would locks be too complex to model check?

What are the potential improvements?

- ▶ 2-level is good
- ▶ Multi-level is great

LevelDB benchmark, x86 server, 96 hyperthreads

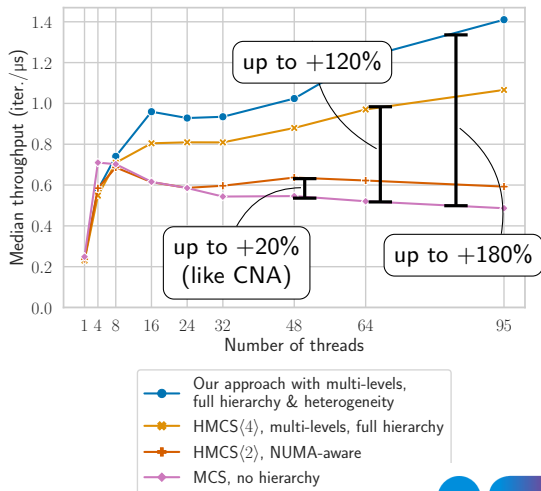


Why would locks be too complex to model check?

What are the potential improvements?

- ▶ 2-level is good
- ▶ Multi-level is great
- ▶ Multi-level+heterogeneity is awesome!

LevelDB benchmark, x86 server, 96 hyperthreads



Why would locks be too complex to model check?

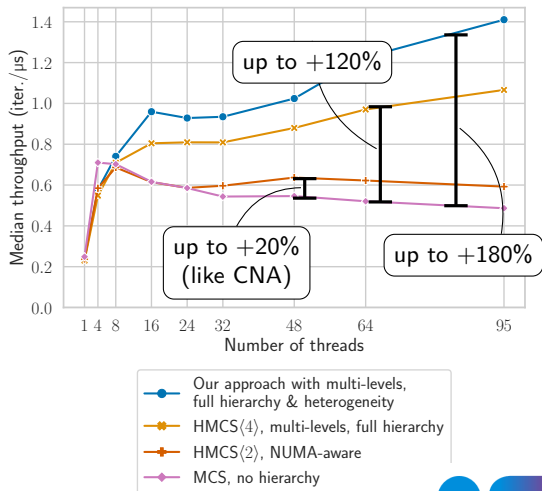
What are the potential improvements?

- ▶ 2-level is good
- ▶ Multi-level is great
- ▶ Multi-level+heterogeneity is awesome!

But model checkers can't handle that

- ▶ Example multi-level HMCS:

LevelDB benchmark, x86 server, 96 hyperthreads



Why would locks be too complex to model check?

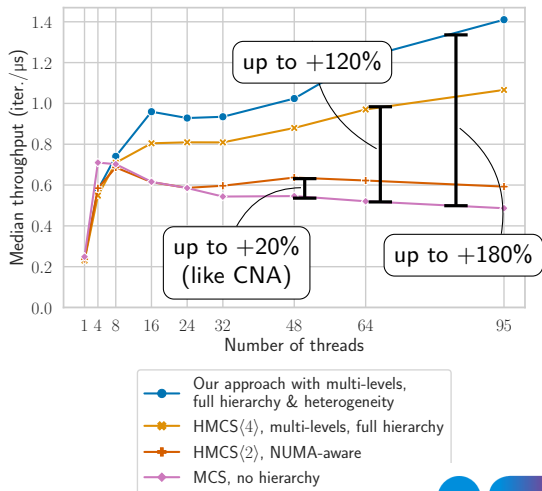
What are the potential improvements?

- ▶ 2-level is good
- ▶ Multi-level is great
- ▶ Multi-level+heterogeneity is awesome!

But model checkers can't handle that

- ▶ Example multi-level HMCS:
 - ▶ 2 levels: 2s
 - ▶ 3 levels: 10s
 - ▶ 4 levels: timeout after 24h

LevelDB benchmark, x86 server, 96 hyperthreads



Why would locks be too complex to model check?

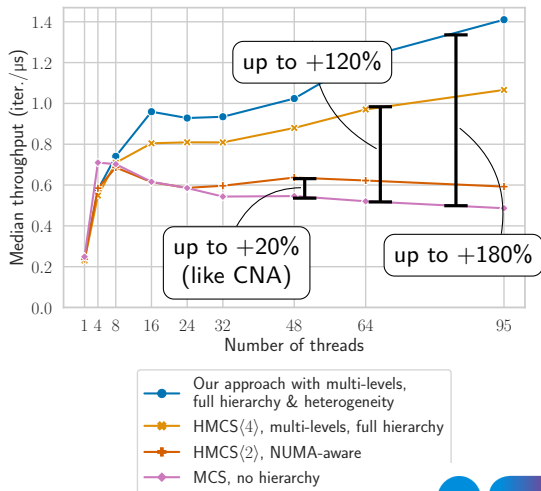
What are the potential improvements?

- ▶ 2-level is good
- ▶ Multi-level is great
- ▶ Multi-level+heterogeneity is awesome!

But model checkers can't handle that

- ▶ Example multi-level HMCS:
 - ▶ 2 levels: 2s
 - ▶ 3 levels: 10s
 - ▶ 4 levels: timeout after 24h
- ▶ **Enter CLoF!**

LevelDB benchmark, x86 server, 96 hyperthreads



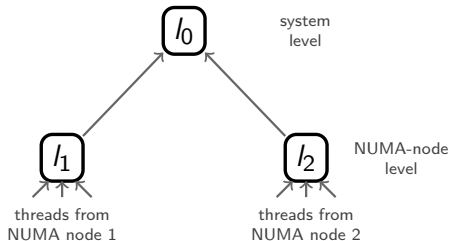
CLoF: Compositional Lock Framework

Published @ SOSP'21 — by Chehab, Paolillo, Behrens, Fu, Chen, Härtig
<https://dl.acm.org/doi/abs/10.1145/3477132.3483557>

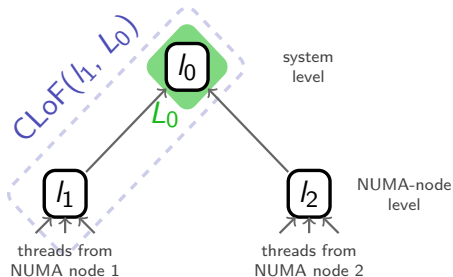
Two NUMA-node example



Two NUMA-node example



Two NUMA-node example

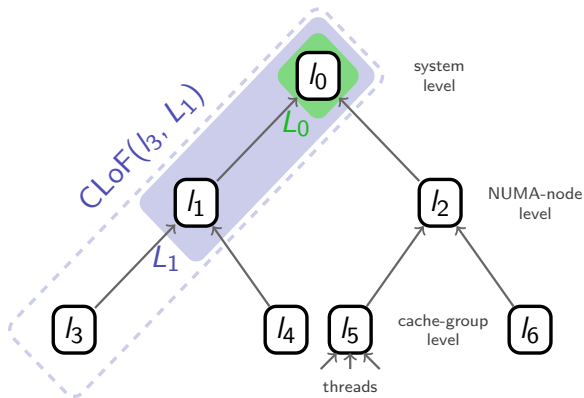


Simplified algorithm

```
CLoF( $l$ ,  $L$ )::acquire =
    acquire  $l$ ;
    if ( $\neg$ already has  $L$ )
        acquire  $L$ ;
```

```
CLoF( $l$ ,  $L$ )::release =
    if (others won't starve)
        release  $l$ ;
    else
        release  $L$ ;
        release  $l$ ;
```

Two NUMA-node example



Simplified algorithm

```
CLoF( $l$ ,  $L$ )::acquire =
  acquire  $l$ ;
  if ( $\neg$ already has  $L$ )
    acquire  $L$ ;
```

```
CLoF( $l$ ,  $L$ )::release =
  if (others won't starve)
    release  $l$ ;
  else
    release  $L$ ;
    release  $l$ ;
```

CLoF: Compositional Lock Framework

Published @ SOSP'21 — by Chehab, Paolillo, Behrens, Fu, Chen, Härtig
<https://dl.acm.org/doi/abs/10.1145/3477132.3483557>

Base Step

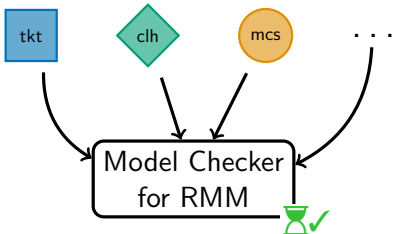
Induction Step



CLoF: Compositional Lock Framework

Published @ SOSP'21 — by Chehab, Paolillo, Behrens, Fu, Chen, Härtig
<https://dl.acm.org/doi/abs/10.1145/3477132.3483557>

Base Step



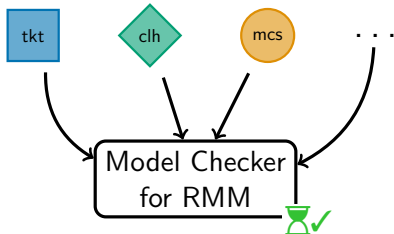
As in VSync paper

Induction Step

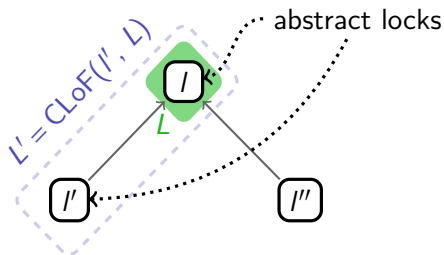
CLoF: Compositional Lock Framework

Published @ SOSP'21 — by Chehab, Paolillo, Behrens, Fu, Chen, Härtig
<https://dl.acm.org/doi/abs/10.1145/3477132.3483557>

Base Step



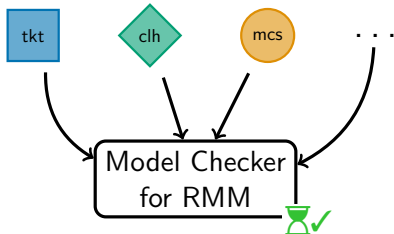
Induction Step



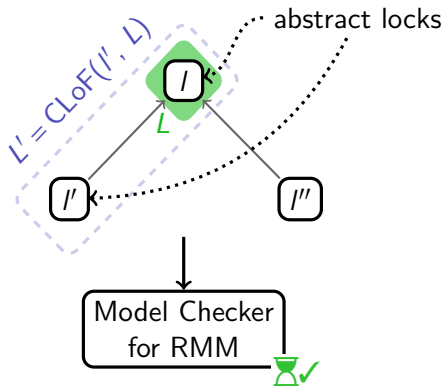
CLoF: Compositional Lock Framework

Published @ SOSP'21 — by Chehab, Paolillo, Behrens, Fu, Chen, Härtig
<https://dl.acm.org/doi/abs/10.1145/3477132.3483557>

Base Step



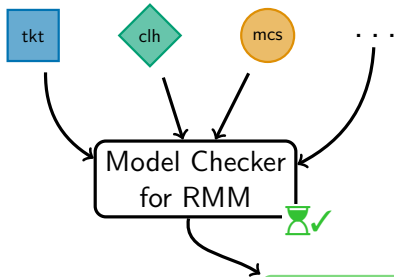
Induction Step



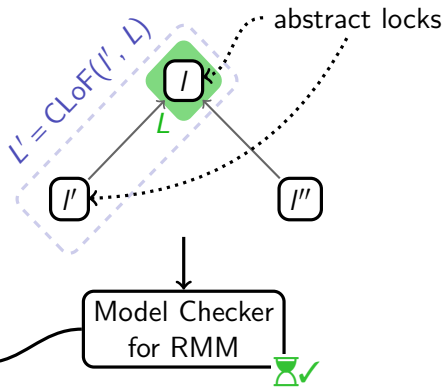
CLoF: Compositional Lock Framework

Published @ SOSP'21 — by Chehab, Paolillo, Behrens, Fu, Chen, Härtig
<https://dl.acm.org/doi/abs/10.1145/3477132.3483557>

Base Step



Induction Step



- ▶ **Concurrency must be smart and hardware tailored**
Otherwise we miss big opportunities
- ▶ **Modularity is essential for model checking**
Ideally by design, not in hindsight
- ▶ **Support proof sketches are faster, scale, and often OK**
Goodbye fully-automated verification 😞

- ▶ Modern Hardware and Concurrent Programs
- ▶ Enabling Performance with Practical Verification
- ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
- ▶ Wrap up and Outlook

Application**Database**

eg, OpenGauss DB

Operating System

eg, OpenHarmony, OpenEuler

**Synchronization
Primitives**eg, spinlock,
mutex, RCU**Lockless Data
Structures**eg, list, queue,
hashtable, etc

----- software memory model -----

Atomic Interface

eg, atomic_{xchg, get_add}

----- hardware memory model -----

Multicore CPU

eg, ARMv8, RISC-V

What about concurrent data structures?



What about concurrent data structures?



Ringbuffers are pervasive

What about concurrent data structures?

Ringbuffers are pervasive



What about concurrent data structures?

Ringbuffers are pervasive



What about concurrent data structures?

Ringbuffers are pervasive



BBQ: Block-based Bounded Queue

In 6 days @ USENIX ATC'22 — by Wang, Behrens, Fu, Oberhauser, Oberhauser, Lei, Chen, Härtig, Chen

► Conventional designs:

- most **favor simplicity**
- performance \sim interference enq/deq



BBQ: Block-based Bounded Queue

In 6 days @ USENIX ATC'22 — by Wang, Behrens, Fu, Oberhauser, Oberhauser, Lei, Chen, Härtig, Chen

► Conventional designs:

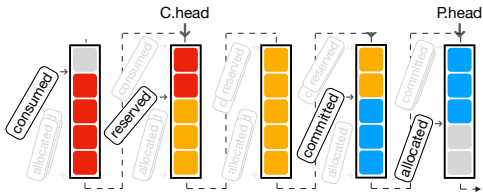
- most **favor simplicity**
- performance \sim interference enq/deq

► Block Approach:

- Split ringbuffer and metadata in blocks
- Drastically reduce enq/deq interference
- 1.5x to 50x higher throughput
- **Verified with VSync**

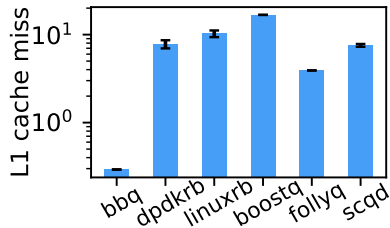
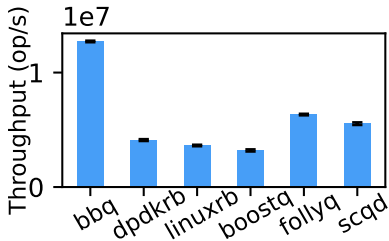


BBQ:



BBQ: Performance Experiments

SPSC Microbenchmarks



2x throughput of Meta's FollyQ

See paper for DPDK, Disruptor, io_uring, multiple modes, etc.

The cost of increased performance

DPDK-like algorithm

```
1 enqueue(data){
2   again:
3   ph = LOAD(P.head);
4   pn = ph + 1;
5   if (pn > LOAD(C.tail) + SZ)
6     return FULL;
7   if (!CAS(P.head, ph, pn))
8     goto again;
9   entry[pn % SZ] = data;
10  while(LOAD(P.tail) != pn);
11  STORE(P.tail, pn);
12  return OK;
13 }

14 dequeue(){
15  again:
16  ch = LOAD(C.head);
17  cn = ch + 1;
18  if (cn > LOAD(P.tail))
19    return EMPTY;
20  if (!CAS(C.head, ch, cn))
21    goto again;
22  data = entry[cn % SZ];
23  while(LOAD(C.tail) != cn);
24  STORE(C.tail, cn);
25  return data;
26 }
```

≈ 10 atomics

The cost of increased performance

Part of BBQ

DPDK-like algorithm

```

1 enqueue(data){
2   again:
3   ph = LOAD(P.head);
4   pn = ph + 1;
5   if (pn > LOAD(C.tail) + SZ)
6     return FULL;
7   if (!CAS(P.head, ph, pn))
8     goto again;
9   entry[pn % SZ] = data;
10  while(LOAD(P.tail) != pn){
11    STORE(P.tail, pn);
12    return OK;
13 }

14 dequeue(){
15   again:
16   ch = LOAD(C.head);
17   cn = ch + 1;
18   if (cn > LOAD(P.tail))
19     return EMPTY;
20   if (!CAS(C.head, ch, cn))
21     goto again;
22   data = entry[cn % SZ];
23   while(LOAD(C.tail) != ch){
24     STORE(C.tail, cn);
25     return data;
26 }

```

≈ 10 atomics

increased complexity

```

1 <Head, Block> BBQ<T>::get_phead_and_block(){
2   ph = LOAD(phead);
3   return (ph, blocks[ph.idx]);
4 }
5 state BBQ<T>::allocate_entry(Block blk){
6   if (LOAD(blk.allocated).off >= BLOCK_SIZE)
7     return BLOCK_DONE;
8   old = TBA(blk.allocated, 1).off;
9   if (old >= BLOCK_SIZE)
10    return BLOCK_DONE;
11   return ALLOCATED(EntryDesc{.block=blk, .offset=old});
12 }
13 void BBQ<T>::commit_entry(EntryDesc e, T data){
14   e.block.entries[e.offset] = data;
15   ADD(e.block.committed, 1);
16 }
17 state BBQ<T>::advance_phead(Head ph) {
18   nbhk = blocks[(ph.idx + 1) % BLOCK_NUM];
19   cons = LOAD(nbhk.consumed);
20   if (cons.van < ph.van ||
21       (cons.van == ph.van && cons.off != BLOCK_SIZE)) {
22     reserved = LOAD(nbhk.reserved);
23     if (reserved.off == cons.off) return NO_ENTRY;
24     else return NOT_AVAILABLE;
25 }
26 cmd = LOAD(nbhk.committed);
27 if (cmd.van == ph.van && cmd.off != BLOCK_SIZE)
28   return NOT_AVAILABLE;
29 MAX(nbhk.committed, Cursor{.van=ph.van + 1});
30 MAX(nbhk.allocated, Cursor{.van=ph.van + 1});
31 return SUCCESS;
32 }
33 class BBQ<T> {
34   shared<Head> phead, chead;
35   Block<T>[] blocks;
36 }
37 class Block<T> {
38   shared<Cursor> allocated, committed;
39   shared<Cursor> reserved, consumed;
40   T[] entries;
41 }
42 class EntryDesc {
43   Block block; Offset offset; Version version; }

```

```

44 <Head, Block> BBQ<T>::get_thead_and_block(){
45   ch = LOAD(thead);
46   return (ch, blocks[ch.idx]);
47 }
48 state BBQ<T>::reserve_entry(Block blk){
49   again:
50   reserved = LOAD(blk.reserved);
51   if (reserved.off < BLOCK_SIZE) {
52     committed = LOAD(blk.committed);
53     if (reserved.off == committed.off)
54       return NO_ENTRY;
55     if (committed.off != BLOCK_SIZE) {
56       allocated = LOAD(blk.allocated);
57       if (allocated.off != committed.off)
58         return NOT_AVAILABLE;
59     }
60     if (MAX(blk.reserved, reserved + 1) == reserved)
61       return RESERVED(EntryDesc{.block=blk,
62                                   .offset=reserved.off,
63                                   .version=reserved.van});
64     else goto again;
65 }
66 return BLOCK_DONE(reserved.van);
67 }
68 T BBQ<T>::consume_entry(EntryDesc e){
69   data = e.block.entries[e.offset];
70   ADD(e.block.consumed, 1);
71   allocated = LOAD(e.block.allocated);
72   if (allocated.van != e.version) return NULL;
73   return data;
74 }
75 bool BBQ<T>::advance_thead(Head ch, Version vsn){
76   nbhk = blocks[(ch.idx + 1) % BLOCK_NUM];
77   committed = LOAD(nbhk.committed);
78   if (committed.van != ch.van + 1)
79     return false;
80   MAX(nbhk.consumed, Cursor{.van=ch.van + 1});
81   MAX(nbhk.reserved, Cursor{.van=ch.van + 1});
82   if (committed.van < vsn + (ch.idx == 0))
83     return false;
84   MAX(nbhk.reserved, Cursor{.van=committed.van});
85   MAX(thead, ch + 1);
86   return true;
87 }

```

retry-new mode | drop-old mode

≈ 20 atomics

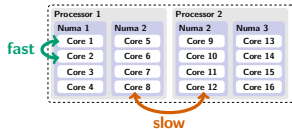
How to ensure correctness?

- ▶ **Long stress testing**
by engineers
- ▶ **Identification of corner cases**
by experts and engineers
- ▶ **Model checking of corner cases**
by engineers with expert support
- ▶ **Only a few corner cases necessary**
queue full/empty, FIFO, wrap-around
- ▶ **3 bugs found model checking them**
Not found while stress testing
- ▶ **Reproducible on real hardware**
Test cases were built in retrospect

- ▶ Modern Hardware and Concurrent Programs
- ▶ Enabling Performance with Practical Verification
- ▶ Practical Verification in Practice
 - ▶ VSync: dealing with Relaxed Memory Models (RMMs)
 - ▶ CNA on RMM: verifying next-gen Linux spinlock
 - ▶ CLoF: dealing with NUMA hierarchies and heterogeneity
 - ▶ BBQ: building highly-efficient ringbuffers
- ▶ **Wrap up and Outlook**

Modern hardware features

- ▶ Relaxed Memory Models, eg, Arm, RISC-V
- ▶ Deep NUMA hierarchies



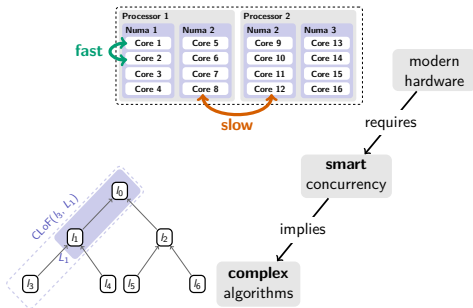
modern
hardware

Modern hardware features

- ▶ Relaxed Memory Models, eg, Arm, RISC-V
- ▶ Deep NUMA hierarchies

Consequences to concurrency

- ▶ Must be **smarter** to boost performance



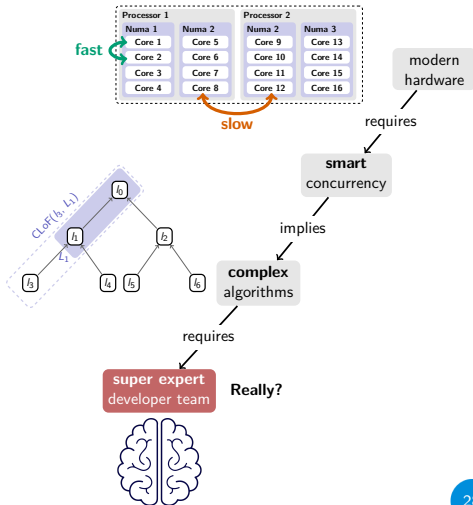
Enabling Performance with Practical Verification

Modern hardware features

- ▶ Relaxed Memory Models, eg, Arm, RISC-V
- ▶ Deep NUMA hierarchies

Consequences to concurrency

- ▶ Must be **smarter** to boost performance
- ▶ But **complexity gets out of control!**



Enabling Performance with Practical Verification

Modern hardware features

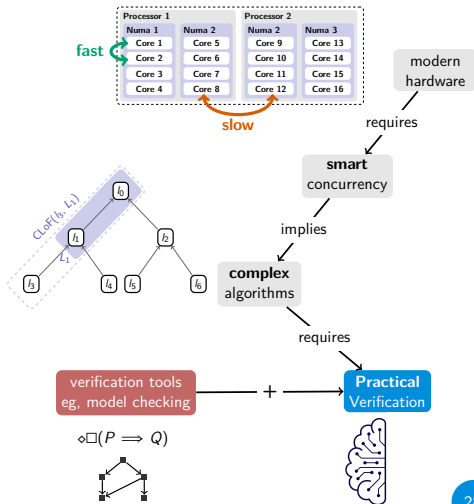
- ▶ Relaxed Memory Models, eg, Arm, RISC-V
- ▶ Deep NUMA hierarchies

Consequences to concurrency

- ▶ Must be **smarter** to boost performance
- ▶ But **complexity gets out of control!**

Practical verification

- ▶ **Formal verification tools**
 - confidence on code correctness



Enabling Performance with Practical Verification

Modern hardware features

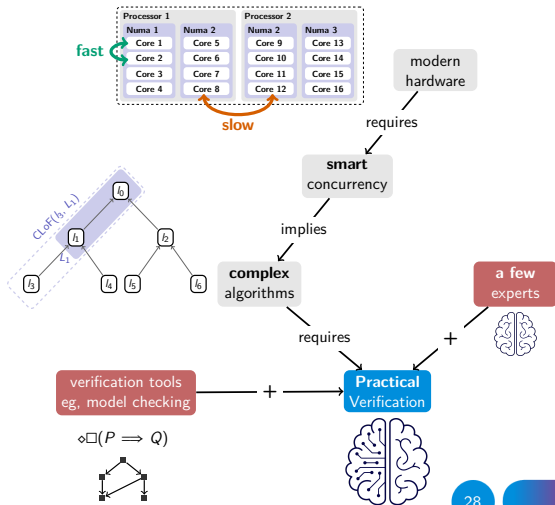
- ▶ Relaxed Memory Models, eg, Arm, RISC-V
- ▶ Deep NUMA hierarchies

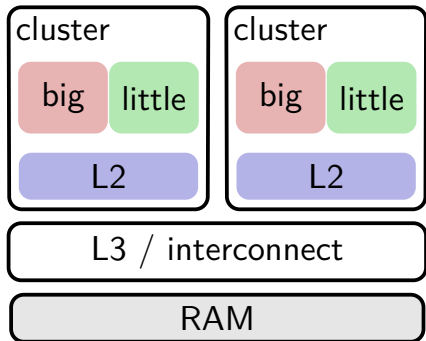
Consequences to concurrency

- ▶ Must be **smarter** to boost performance
- ▶ But **complexity gets out of control!**

Practical verification

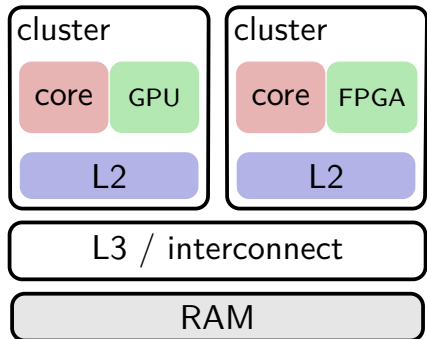
- ▶ **Formal verification tools**
 - confidence on code correctness
- ▶ **A few experts**
 - scalability and coverage of tools





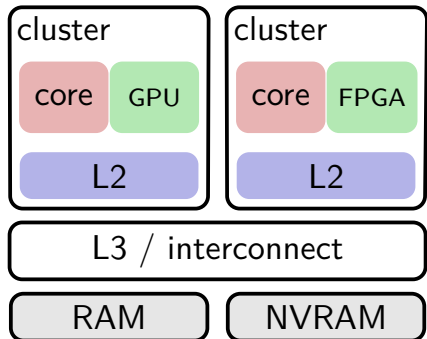
Deeper and more complex hierarchies

- ▶ Heterogeneous processing power
eg, Arm big.LITTLE, Intel Alder Lake



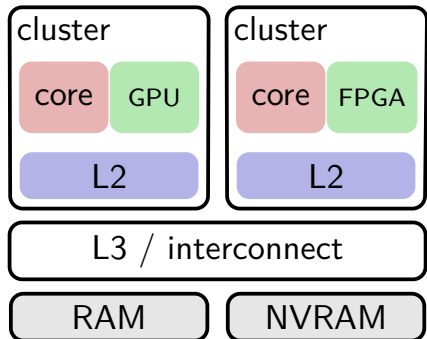
Deeper and more complex hierarchies

- ▶ Heterogeneous processing power
eg, Arm big.LITTLE, Intel Alder Lake
- ▶ Accelerators on shared memory
eg, GPUs, NPUs, FPGAs



Deeper and more complex hierarchies

- ▶ Heterogeneous processing power
eg, Arm big.LITTLE, Intel Alder Lake
- ▶ Accelerators on shared memory
eg, GPUs, NPUs, FPGAs
- ▶ Non-volatile memories



Deeper and more complex hierarchies

- ▶ Heterogeneous processing power
eg, Arm big.LITTLE, Intel Alder Lake
- ▶ Accelerators on shared memory
eg, GPUs, NPU's, FPGAs
- ▶ Non-volatile memories

How to consider everything together?

- ▶ Practical verification FTW!
- ▶ Great potential of HW-SW collaboration!



THANK YOU

非常感谢你

Copyright © 2022 Huawei Technologies Co., Ltd. All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.