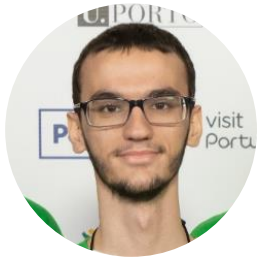


# CLoF: A Compositional Lock Framework for Multi-level NUMA Systems

**Rafael Chehab**



Antonio Paolillo



Diogo Behrens



Ming Fu



Hermann Härtig



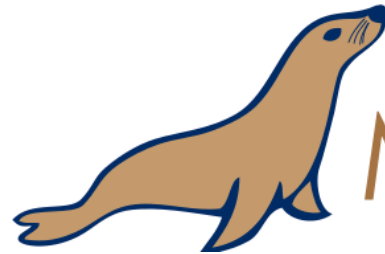
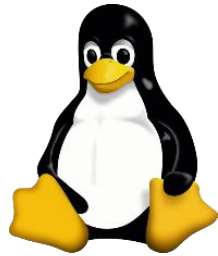
Haibo Chen



March 18<sup>th</sup>, 2022

# Concurrency is Everywhere

Modern operating systems, databases & applications resort to **multi-core concurrency** to achieve high performance.



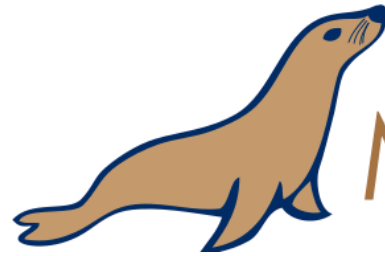
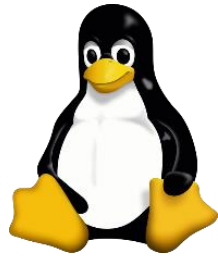
MariaDB



PostgreSQL

# Concurrency is Everywhere

Modern operating systems, databases & applications resort to **multi-core concurrency** to achieve high performance.



MariaDB

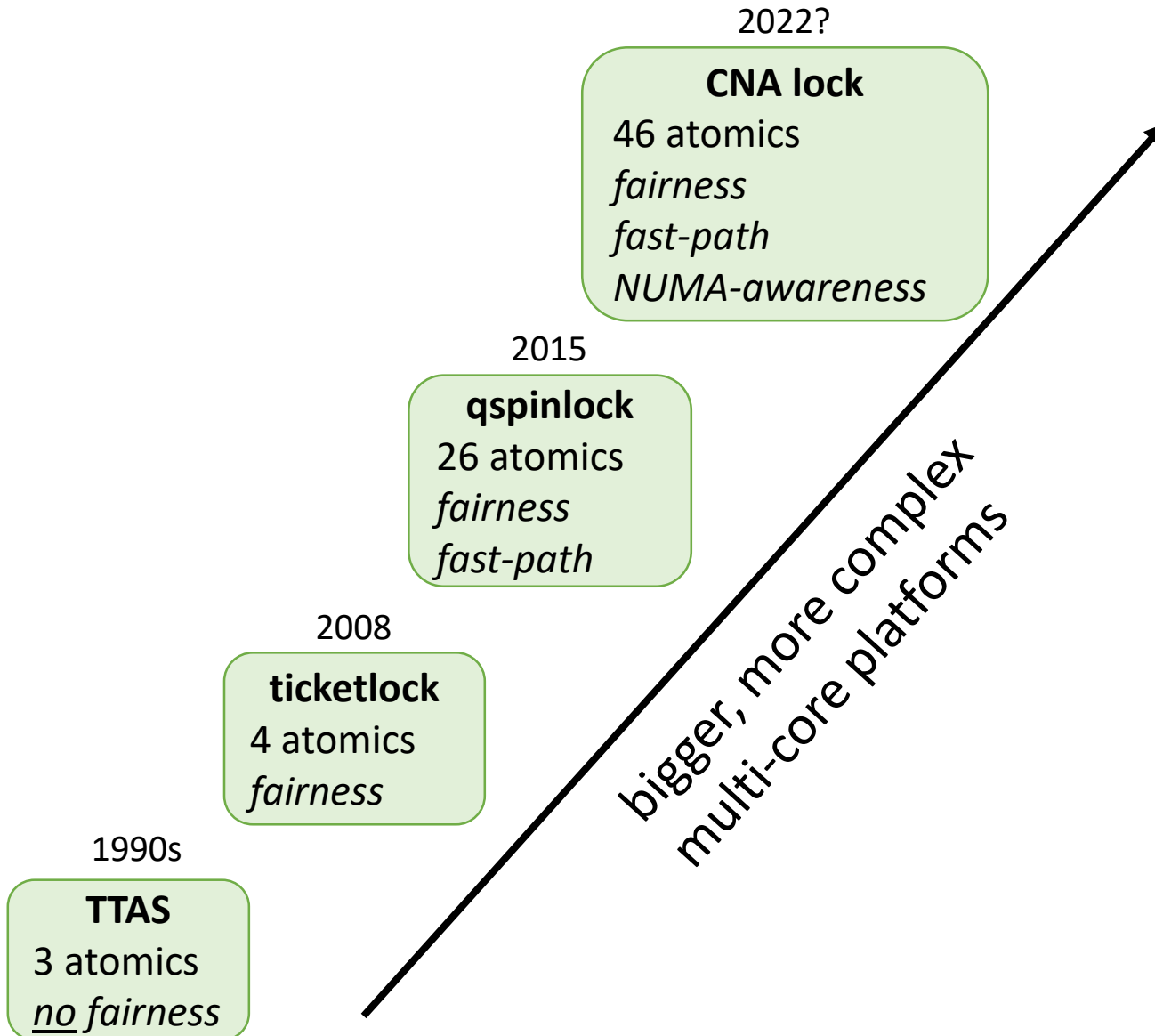


PostgreSQL

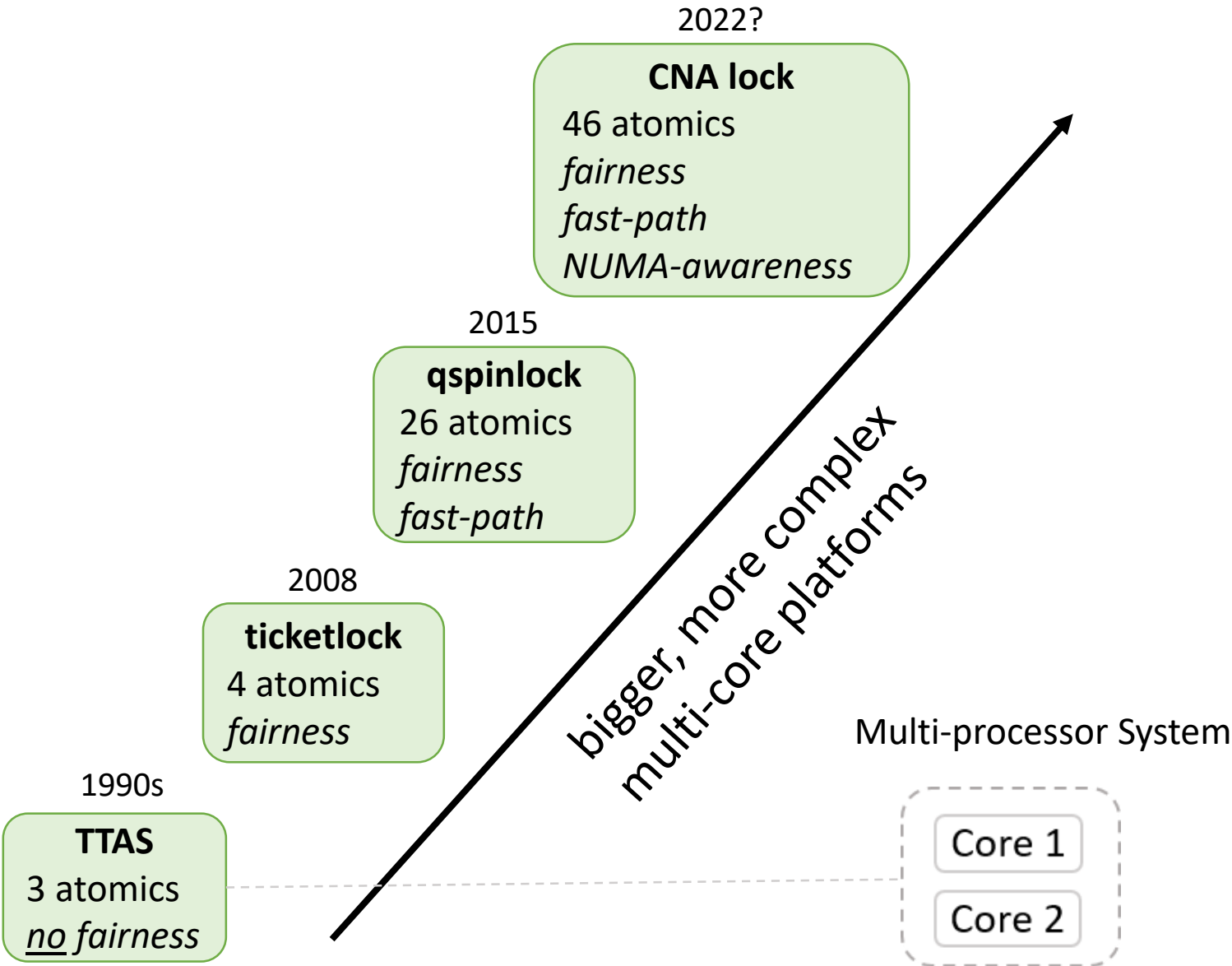
## **Multi-core concurrency:**

One of its most important tasks is to *synchronize* access to shared variables

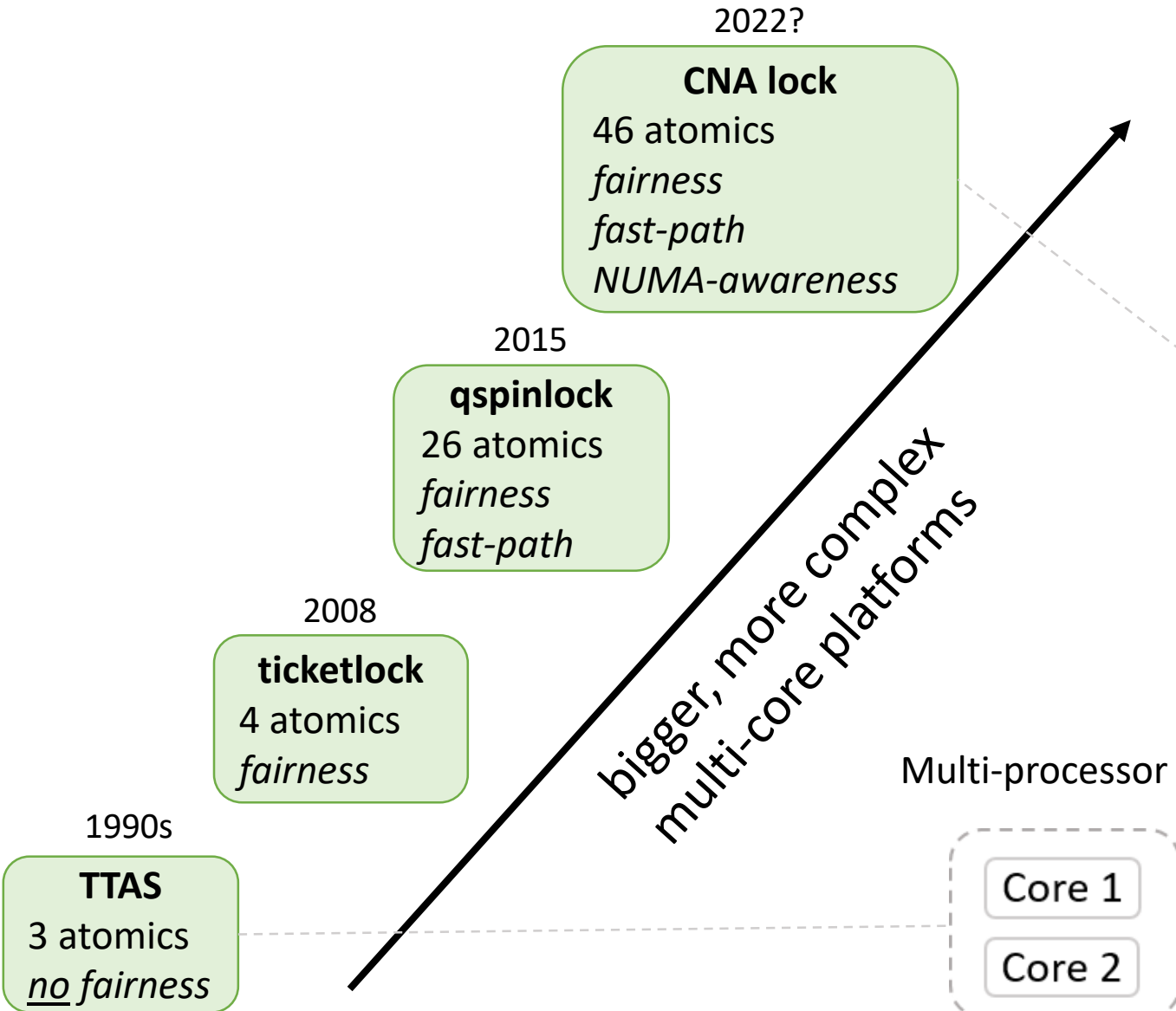
# Linux spinlock evolution



# Linux spinlock evolution

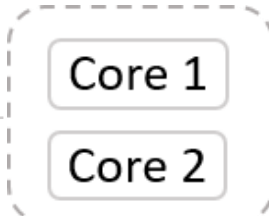


# Linux spinlock evolution

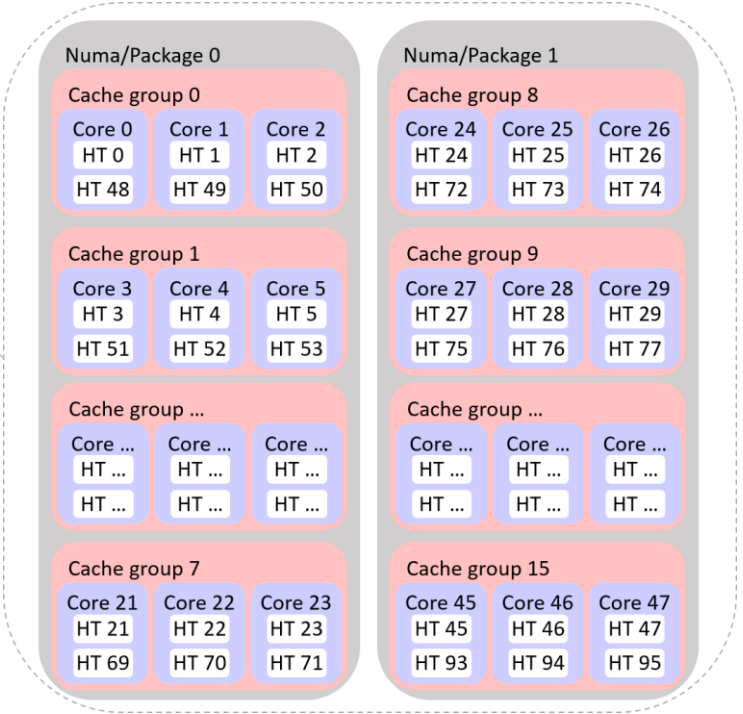


bigger, more complex multi-core platforms

Multi-processor System

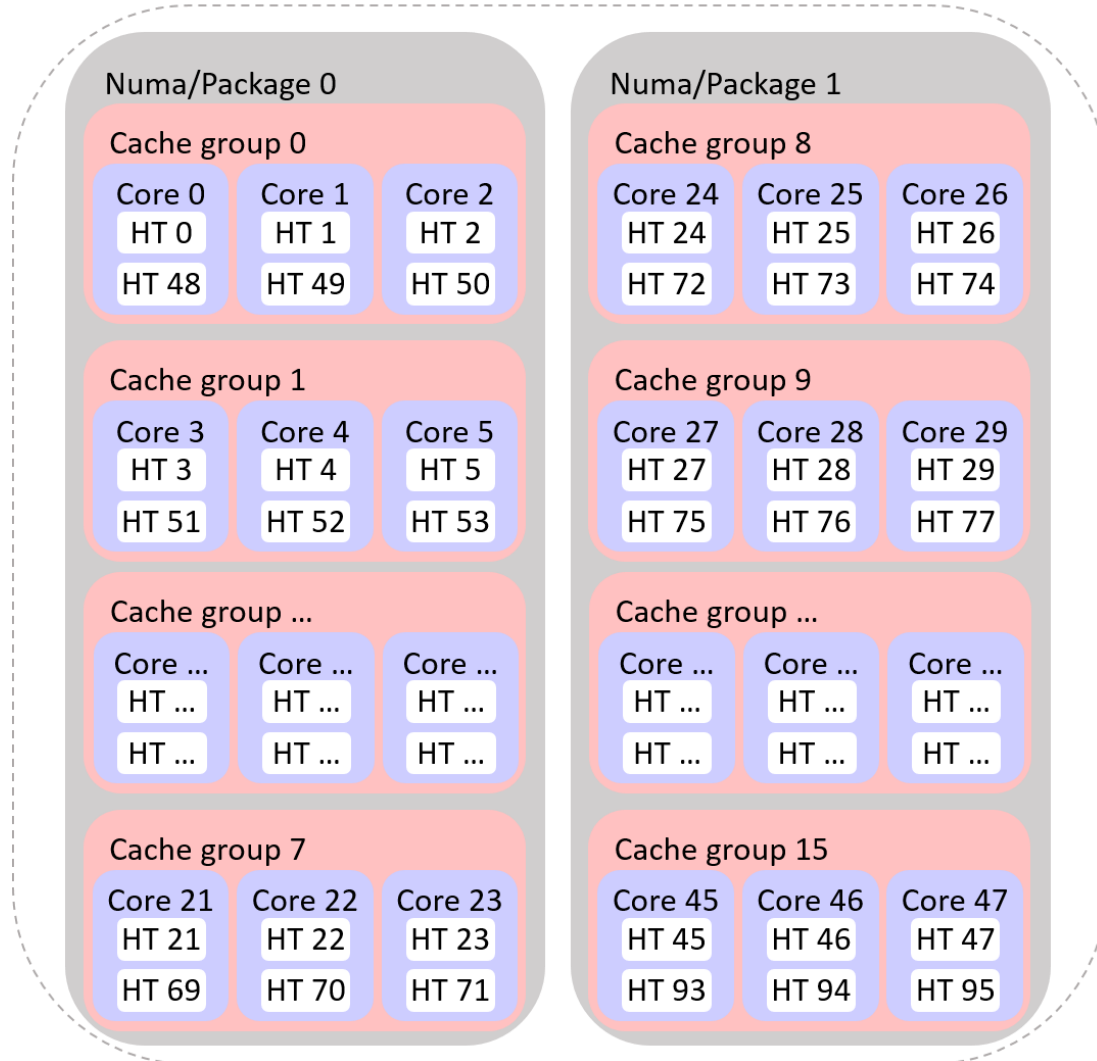


Many-core NUMA system



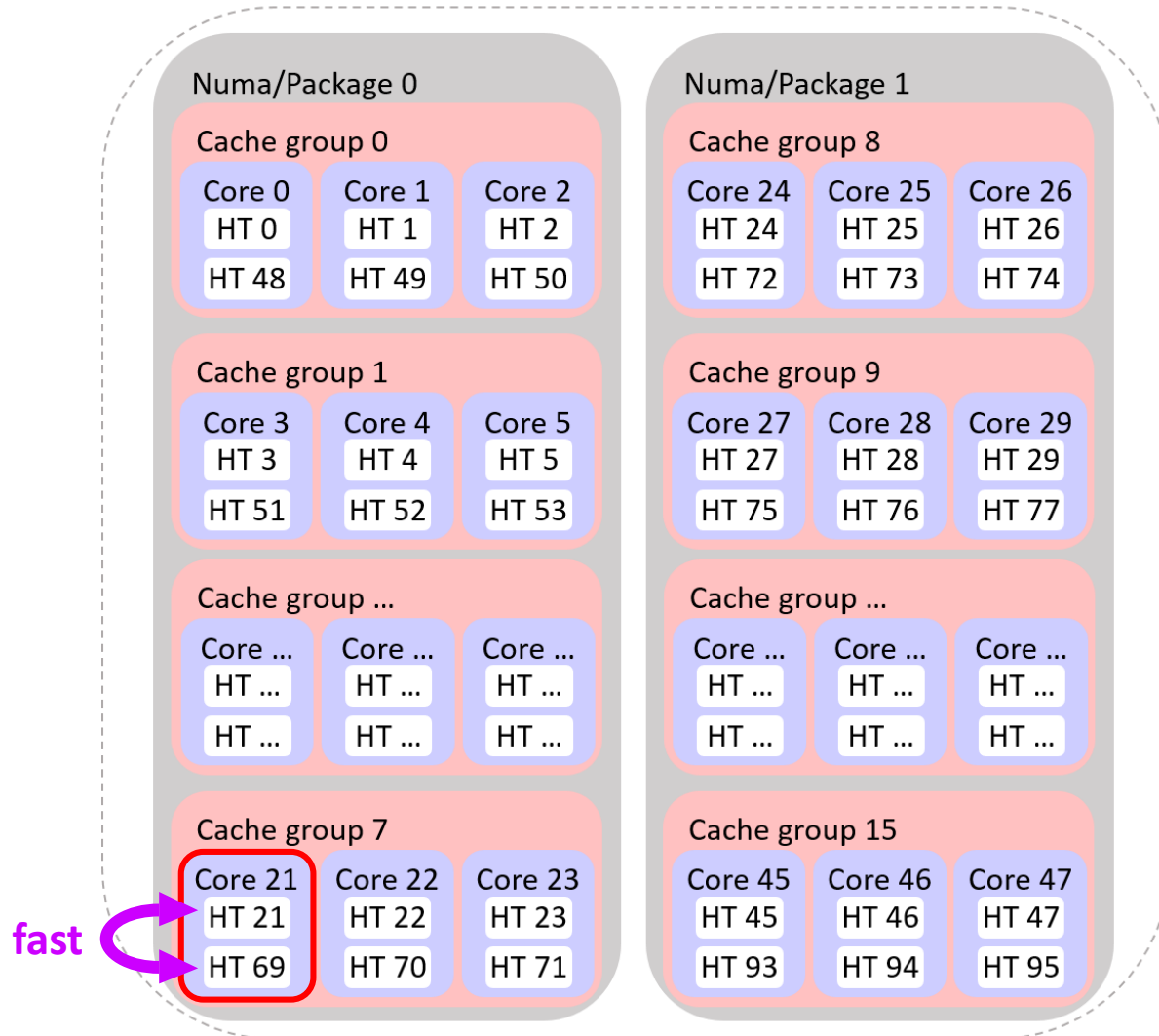
# Deep-hierarchy NUMA systems

Many-core NUMA system



# Deep-hierarchy NUMA systems

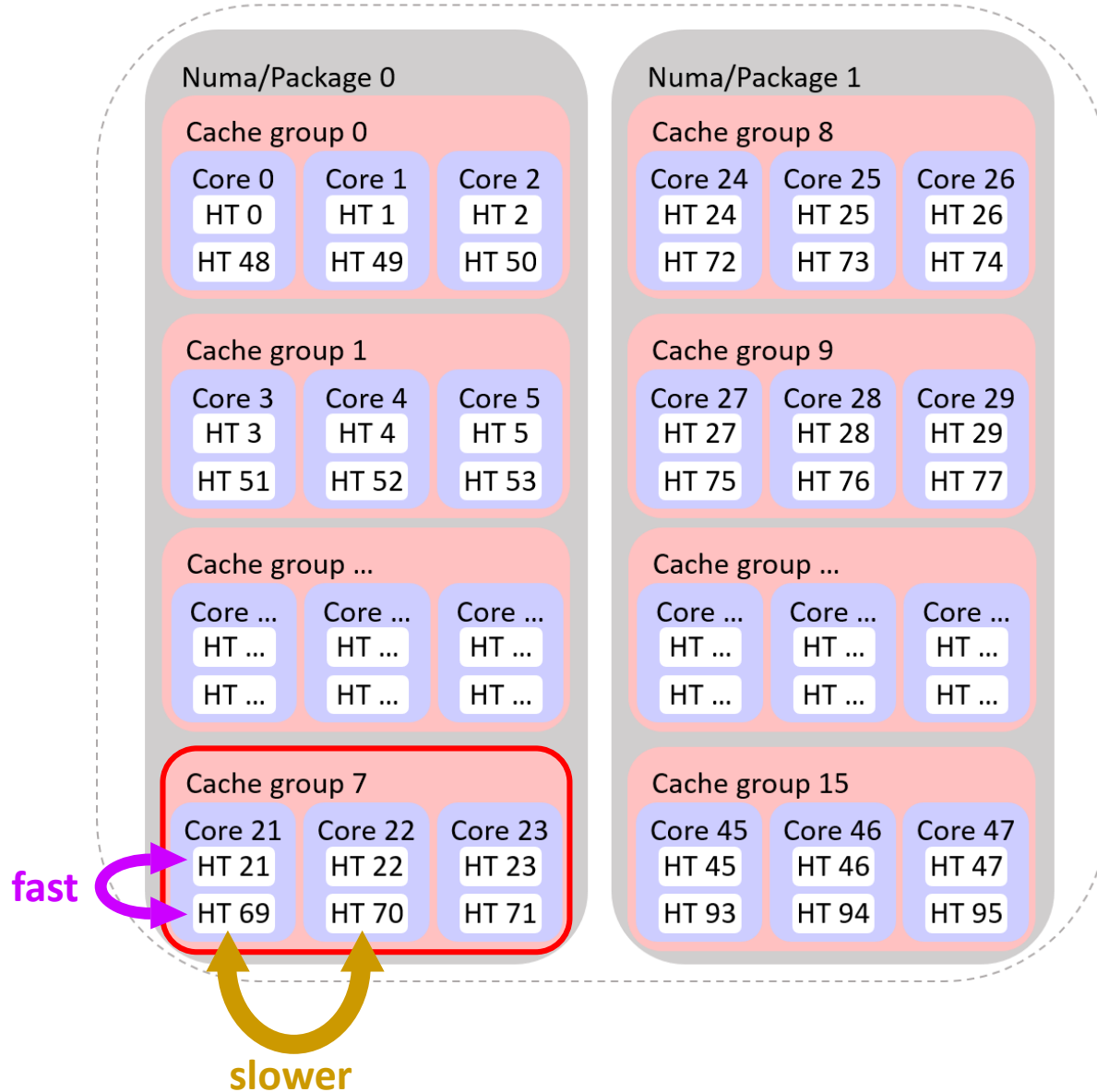
Many-core NUMA system





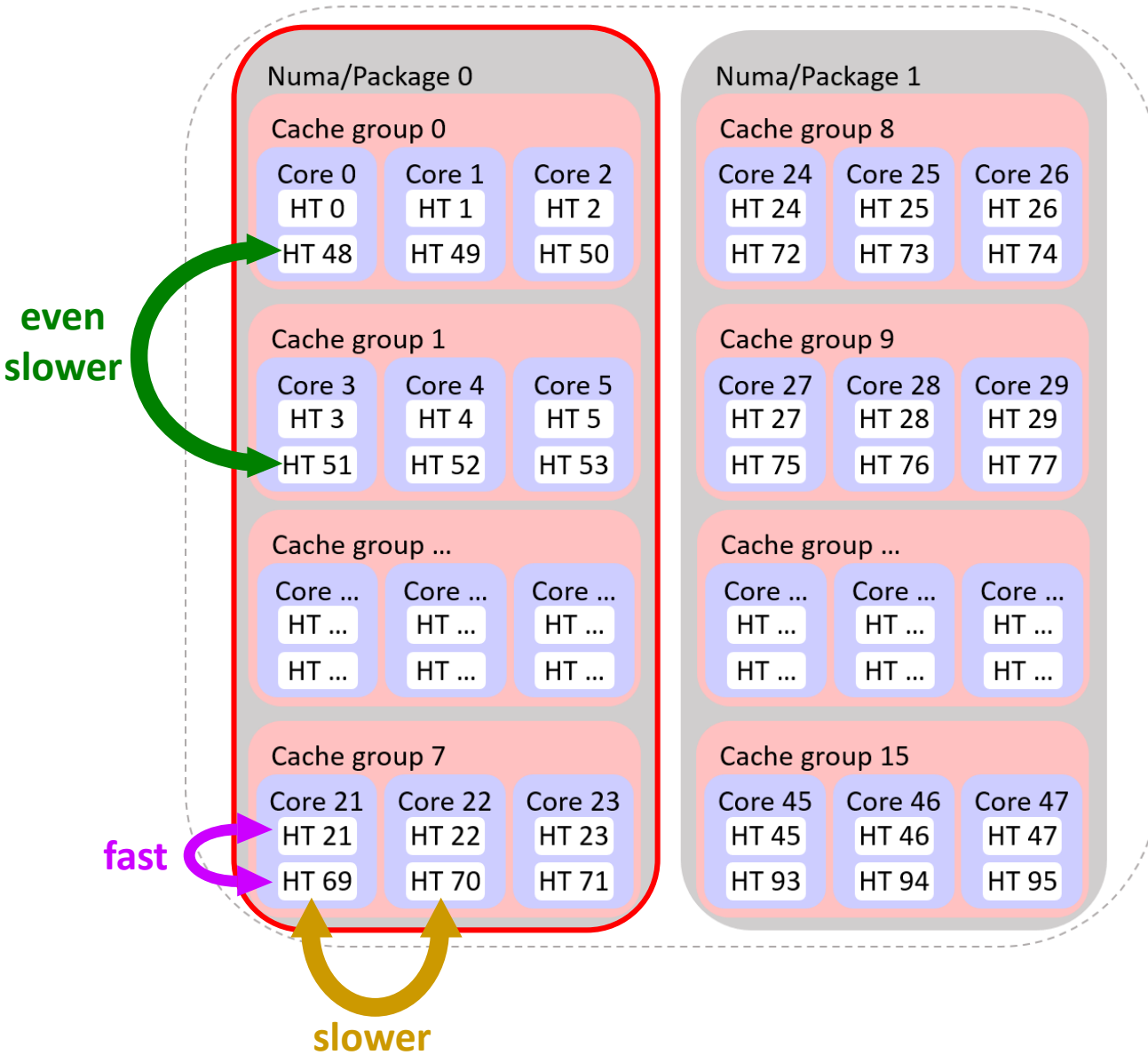
# Deep-hierarchy NUMA systems

Many-core NUMA system



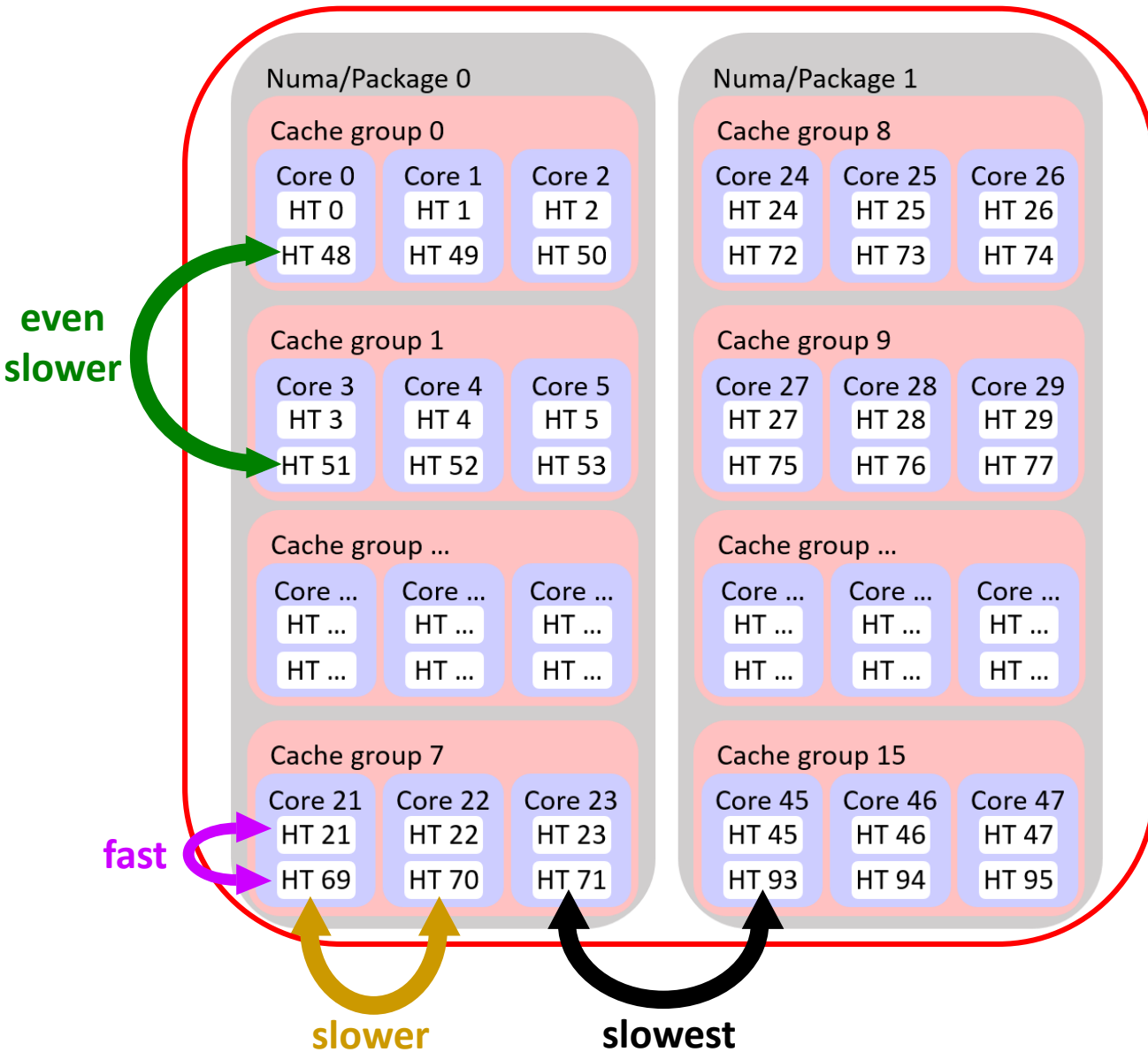
# Deep-hierarchy NUMA systems

Many-core NUMA system



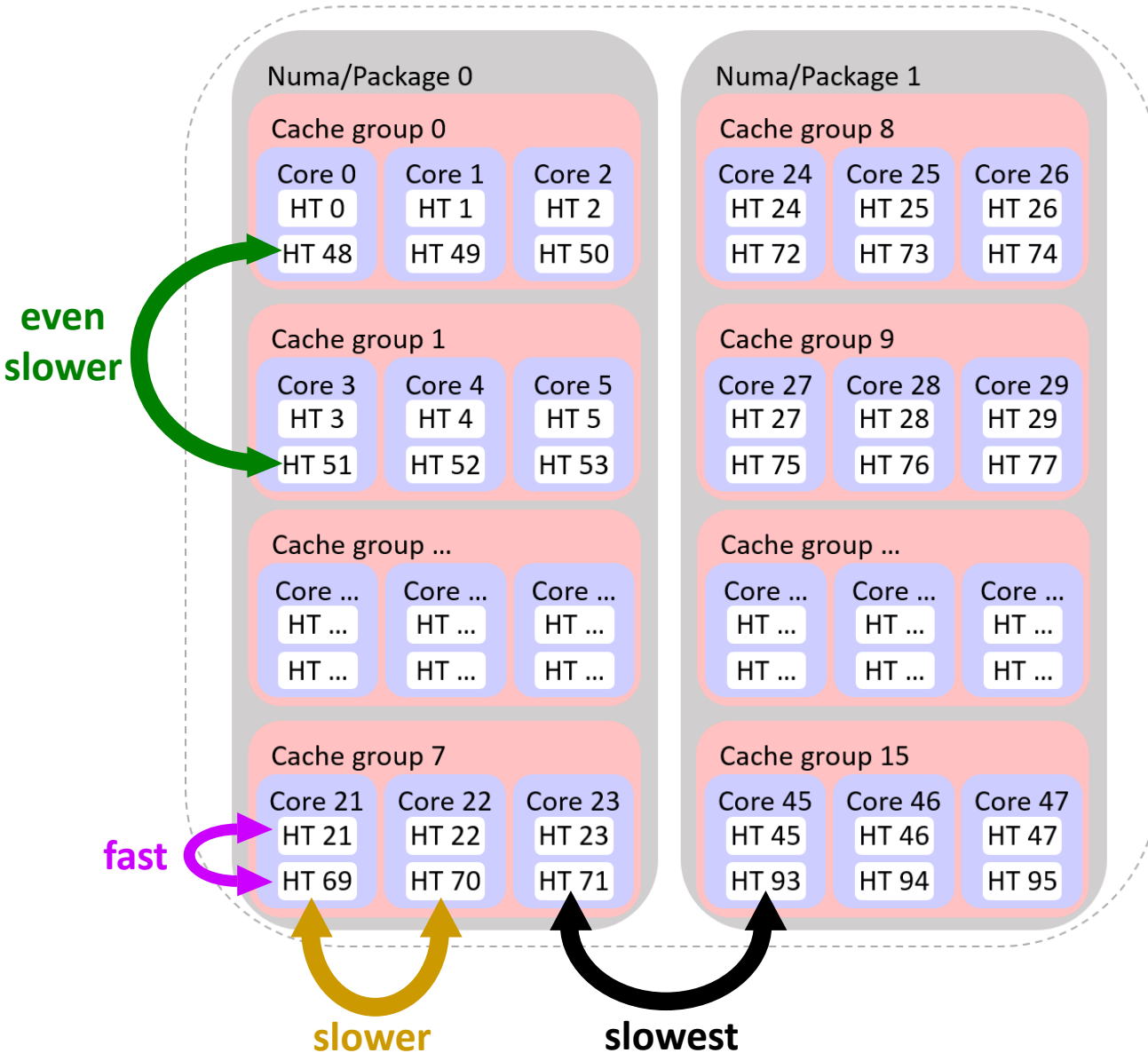
# Deep-hierarchy NUMA systems

Many-core NUMA system



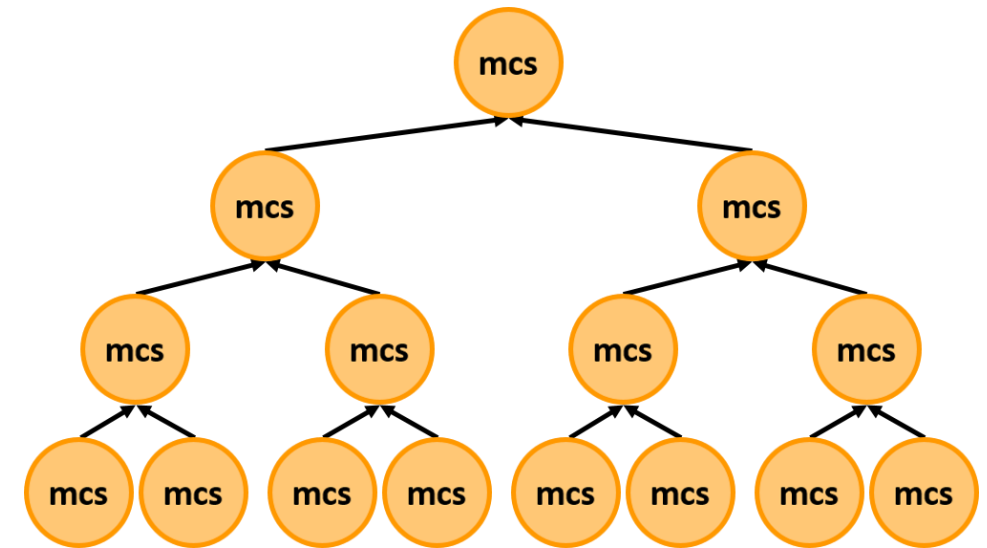
# Deep-hierarchy NUMA systems

Many-core NUMA system



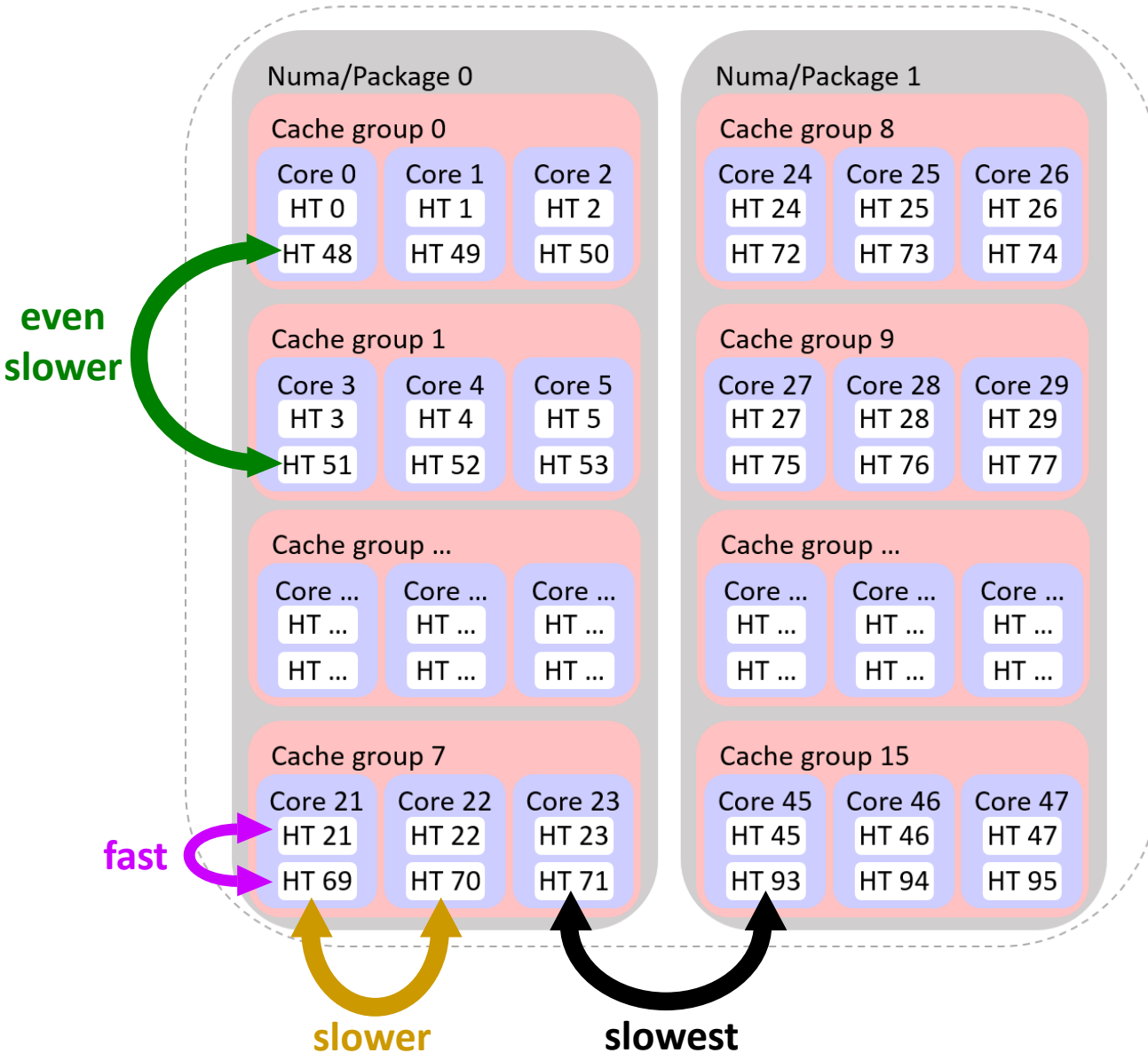
# Multi-level NUMA-aware Locks

- HMCS creates a hierarchy of MCS locks
- Arbitrary number of levels

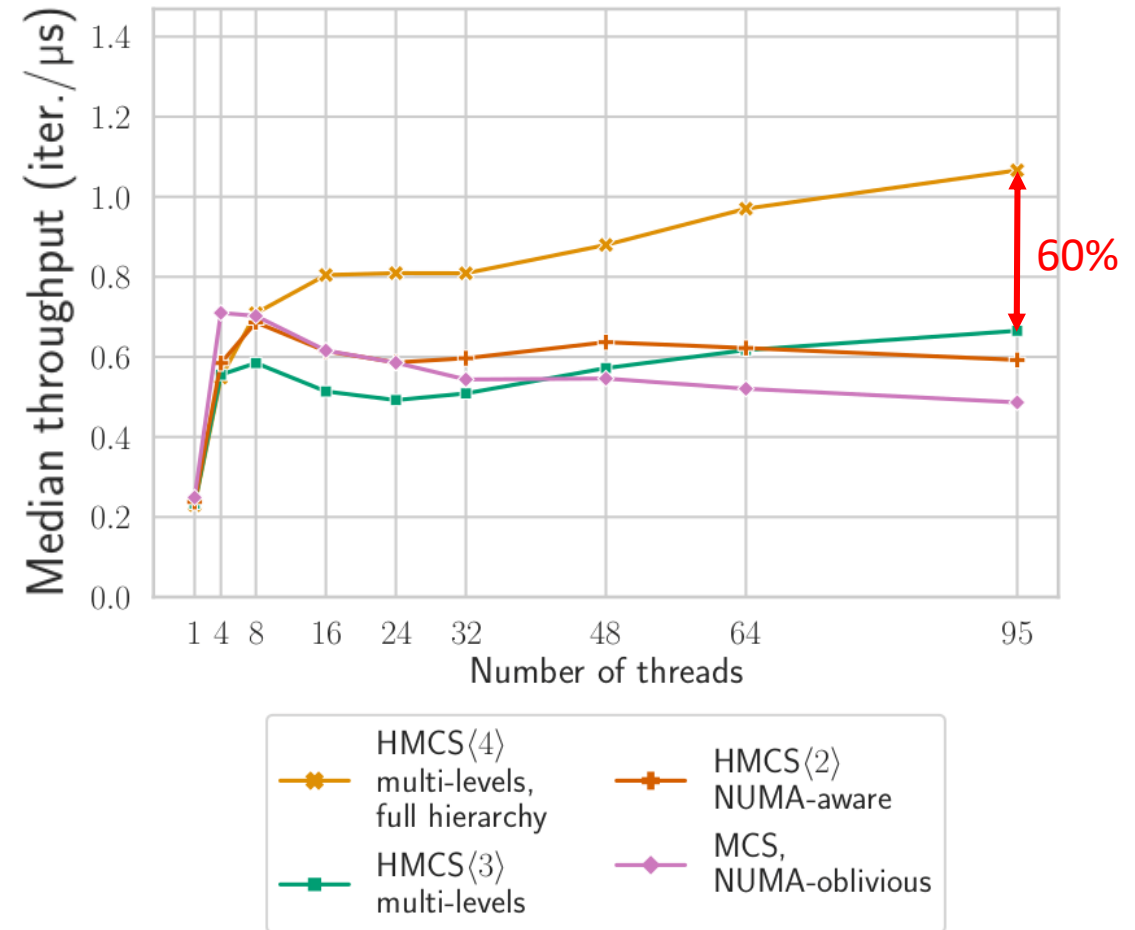


# Deep-hierarchy NUMA systems

Many-core NUMA system

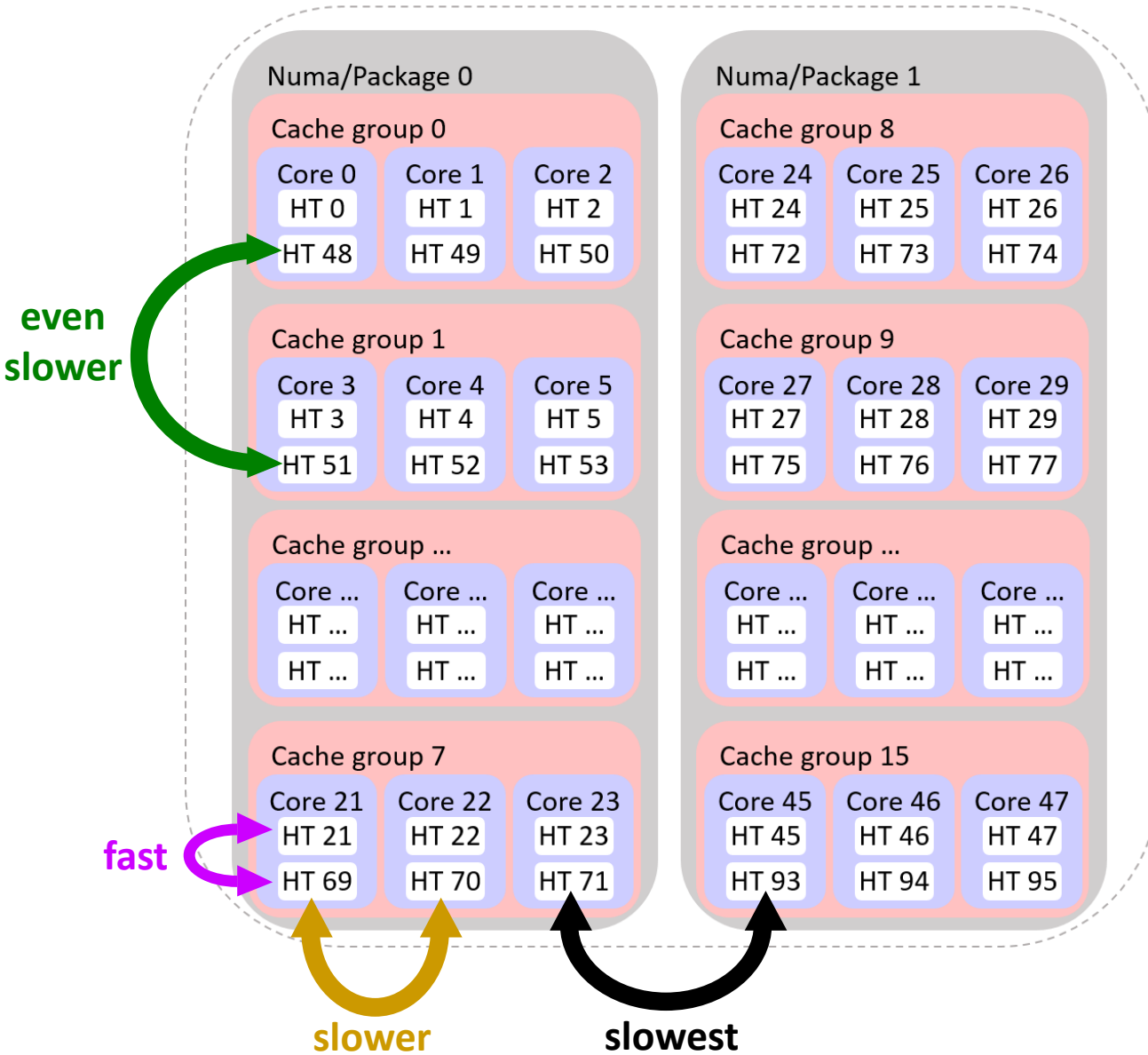


# Multi-level NUMA-aware Locks

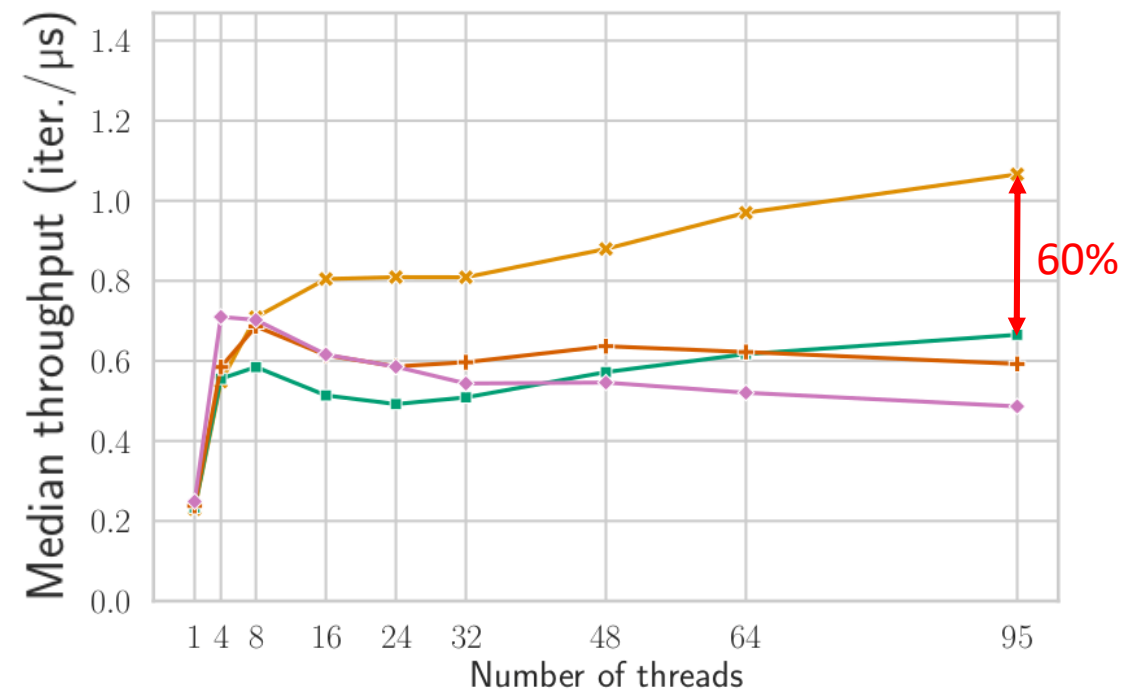


# Deep-hierarchy NUMA systems

Many-core NUMA system



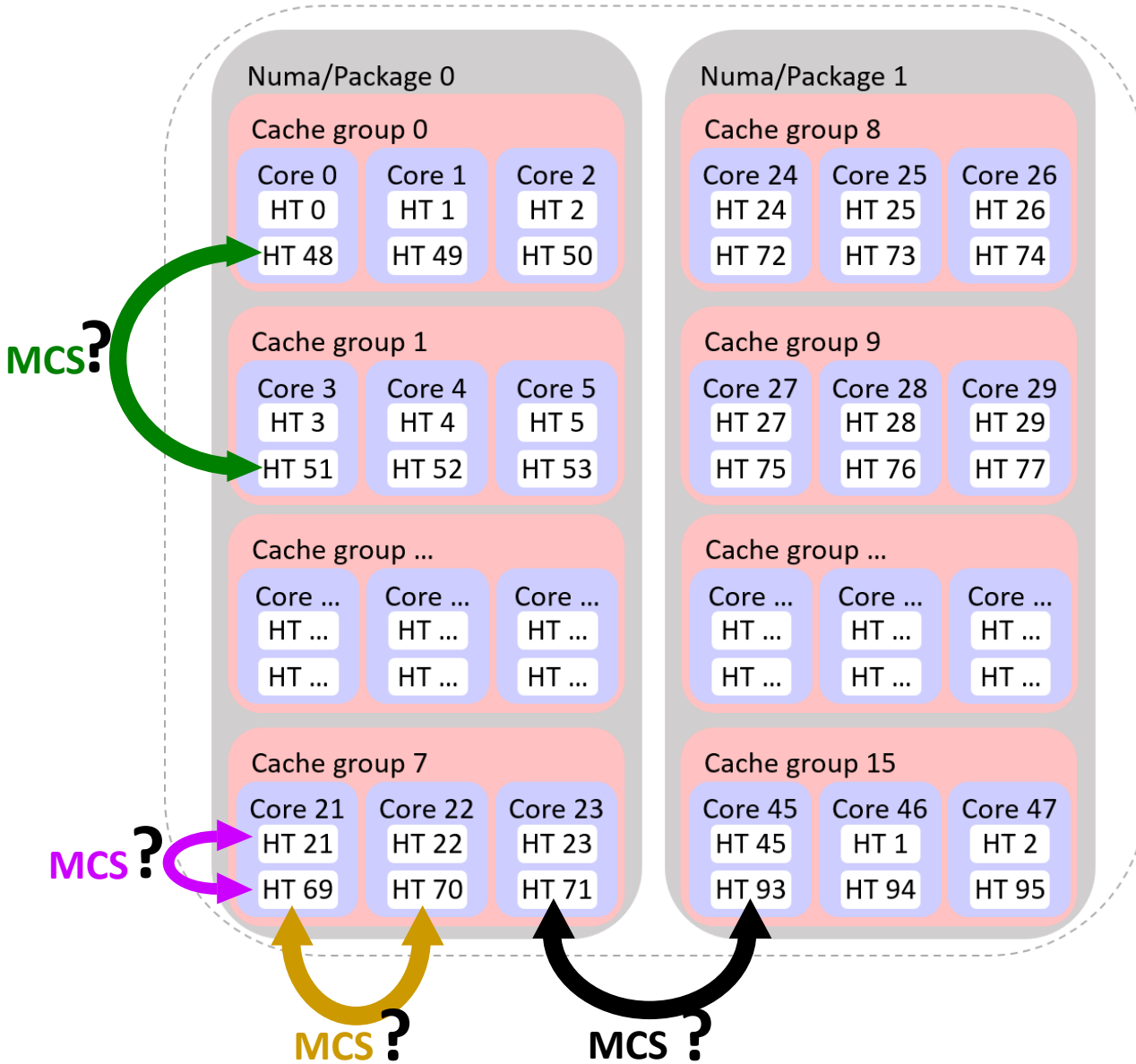
# Multi-level NUMA-aware Locks



*Multi-Level:*  
Utilizing the full deep-hierarchy in a lock improves performance

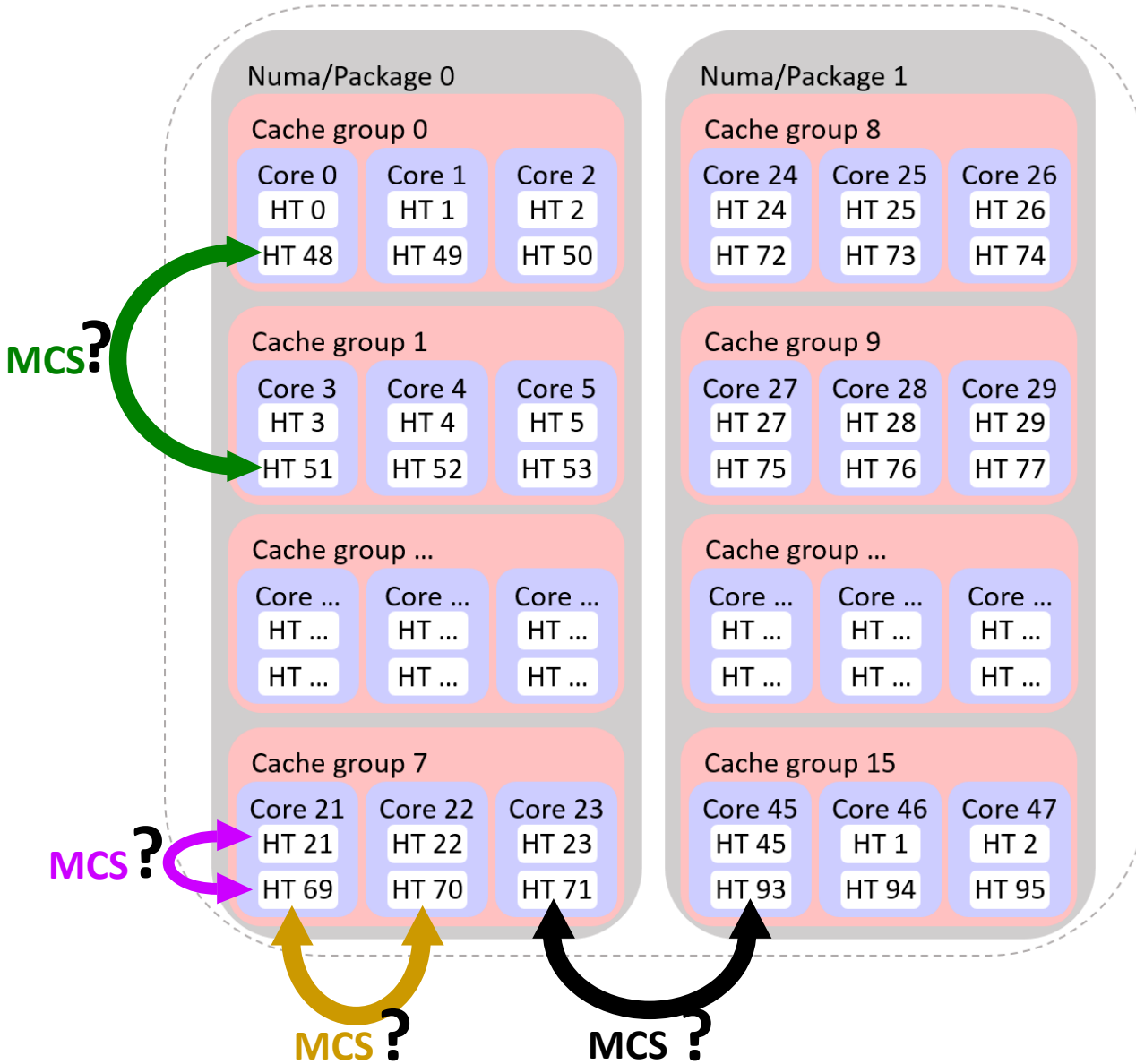
# Level-heterogeneity

Many-core NUMA system



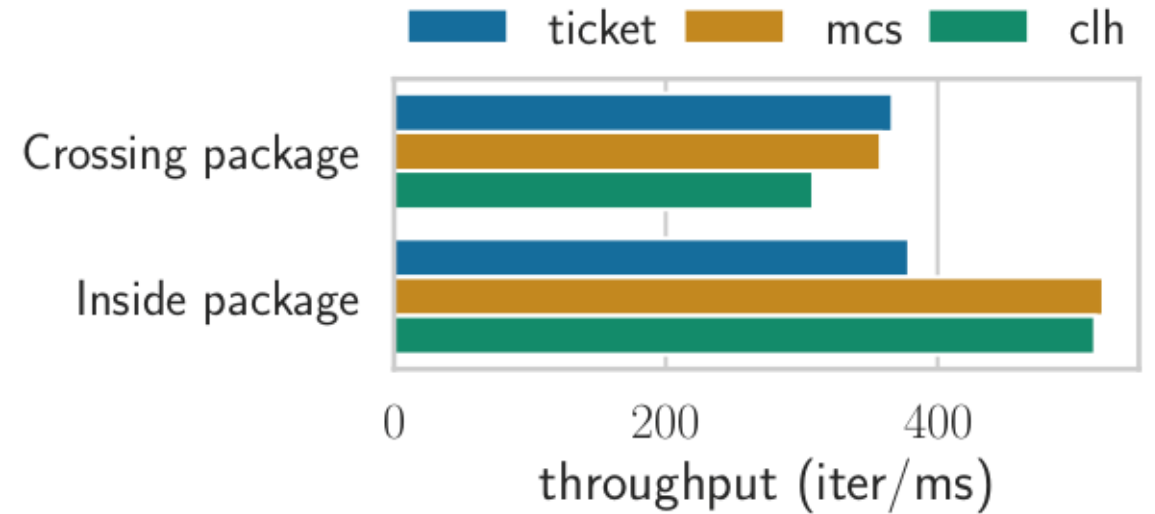
# Level-heterogeneity

Many-core NUMA system



# Experiment - Heterogeneity

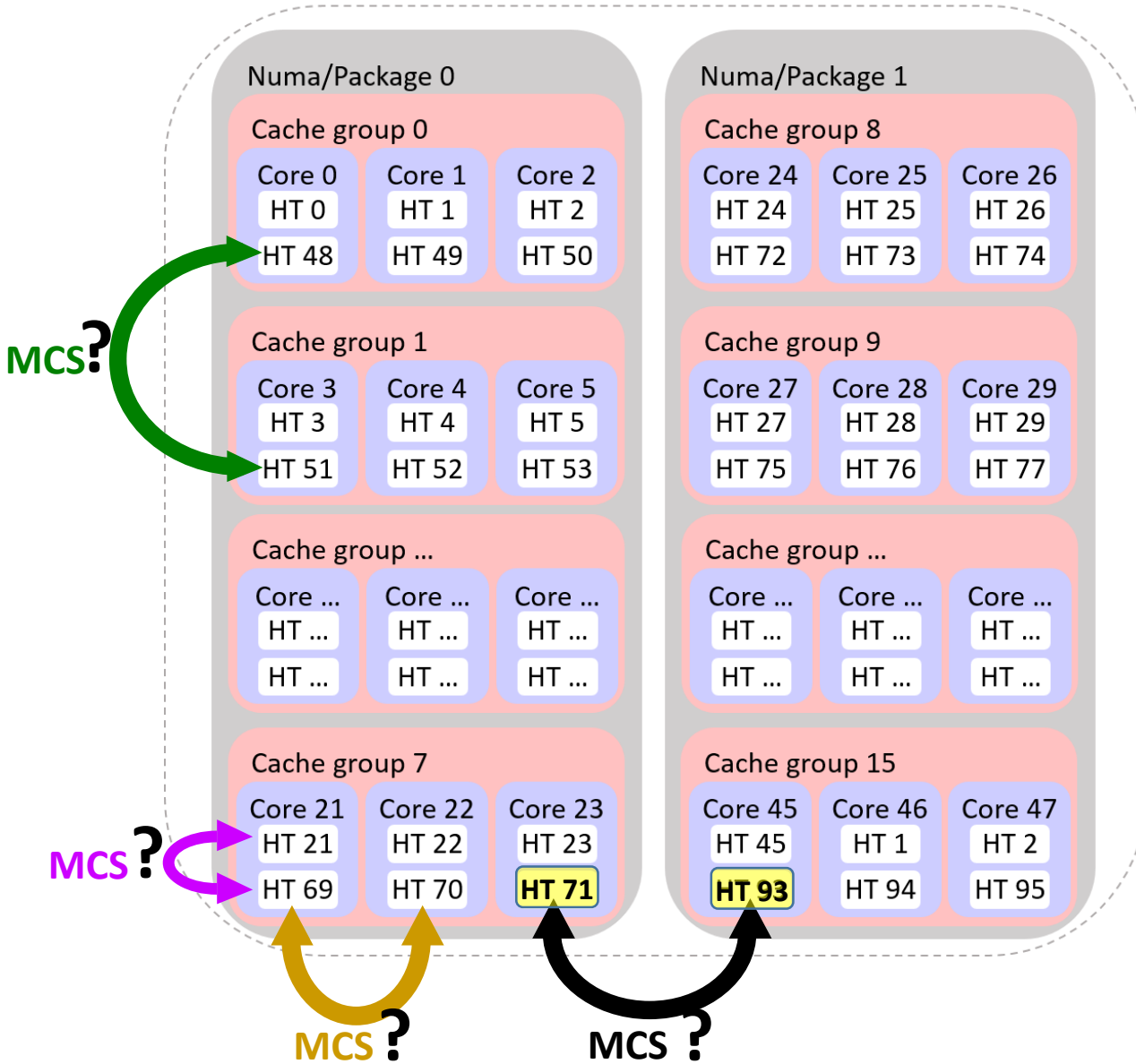
x86 server – execution of classical locks in isolation





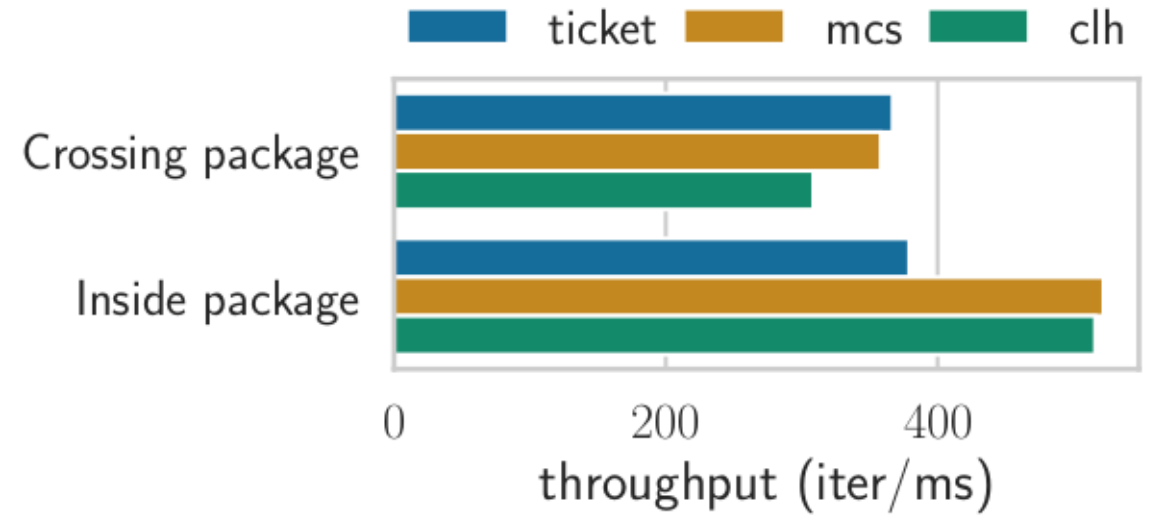
# Level-heterogeneity

Many-core NUMA system



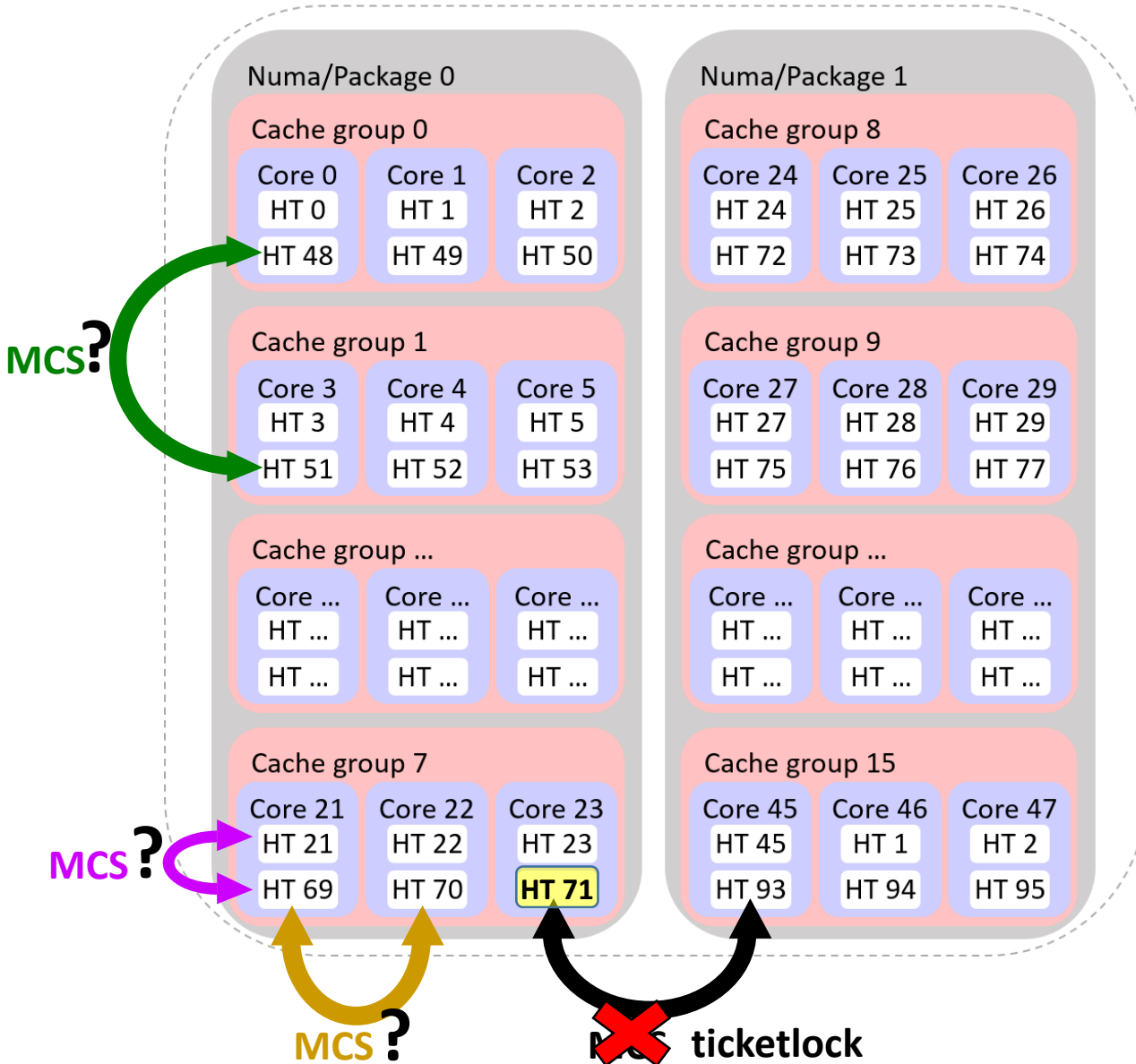
# Experiment - Heterogeneity

x86 server – execution of classical locks in isolation



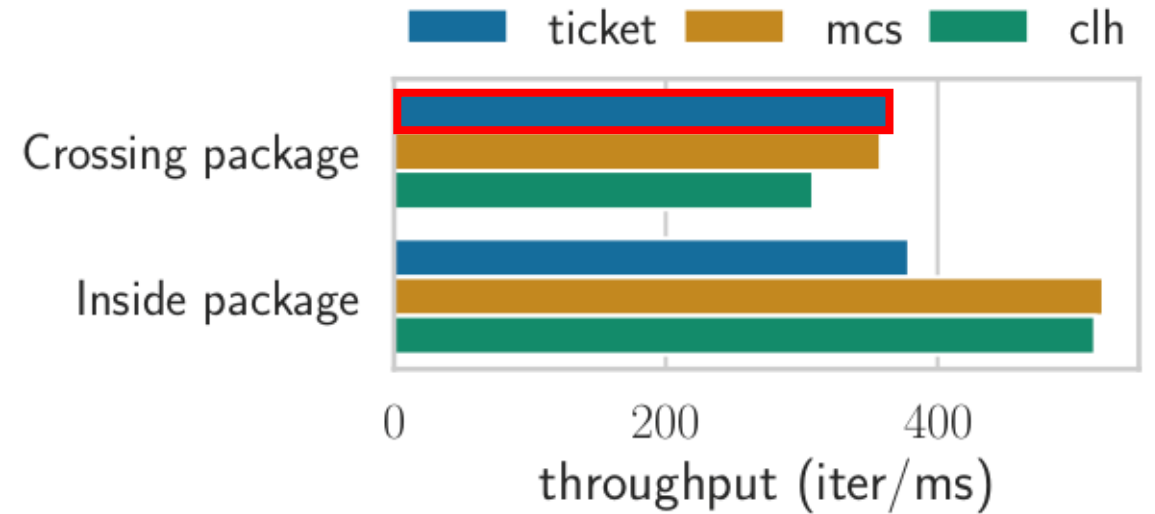
# Level-heterogeneity

Many-core NUMA system



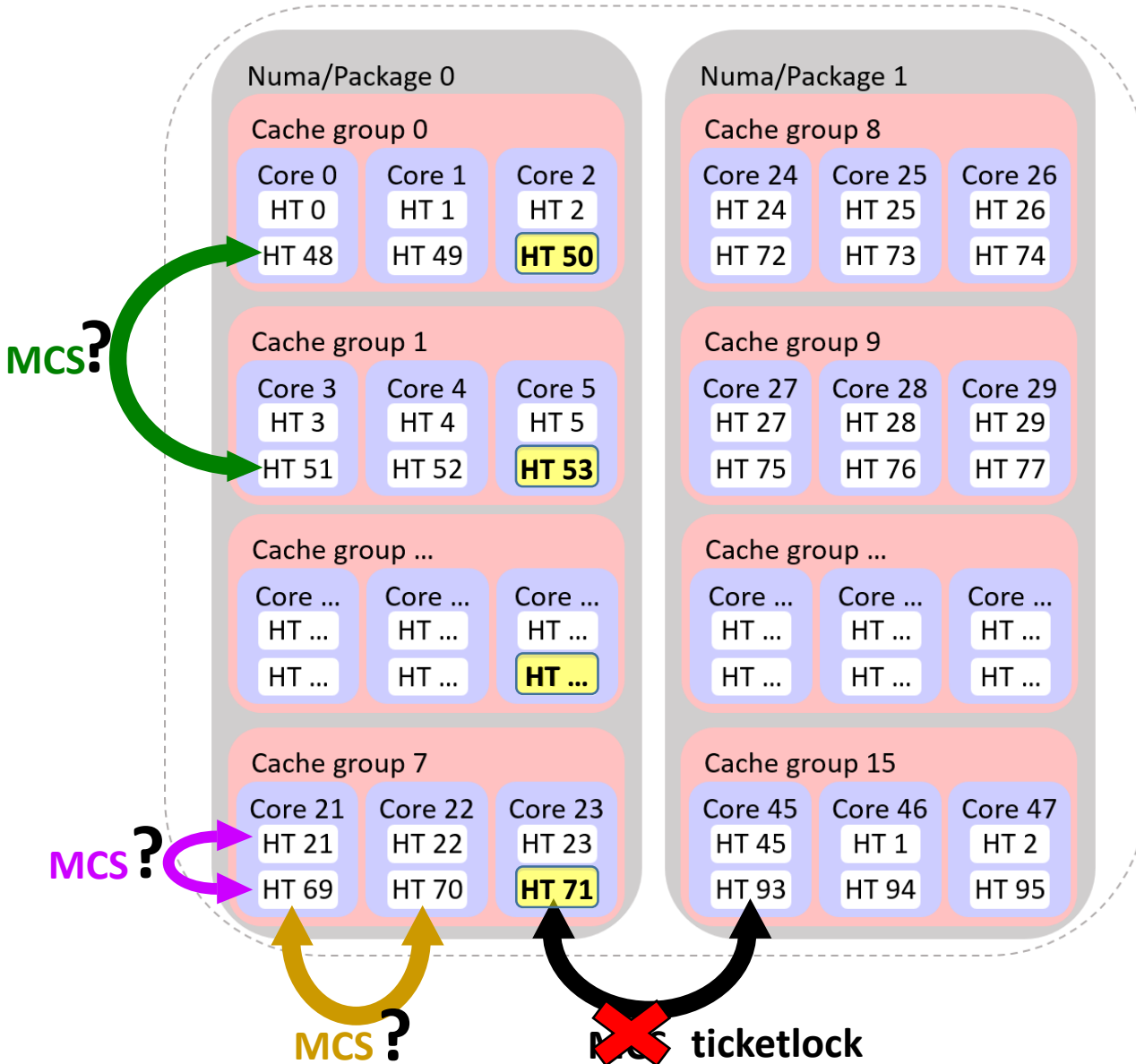
# Experiment - Heterogeneity

x86 server – execution of classical locks in isolation



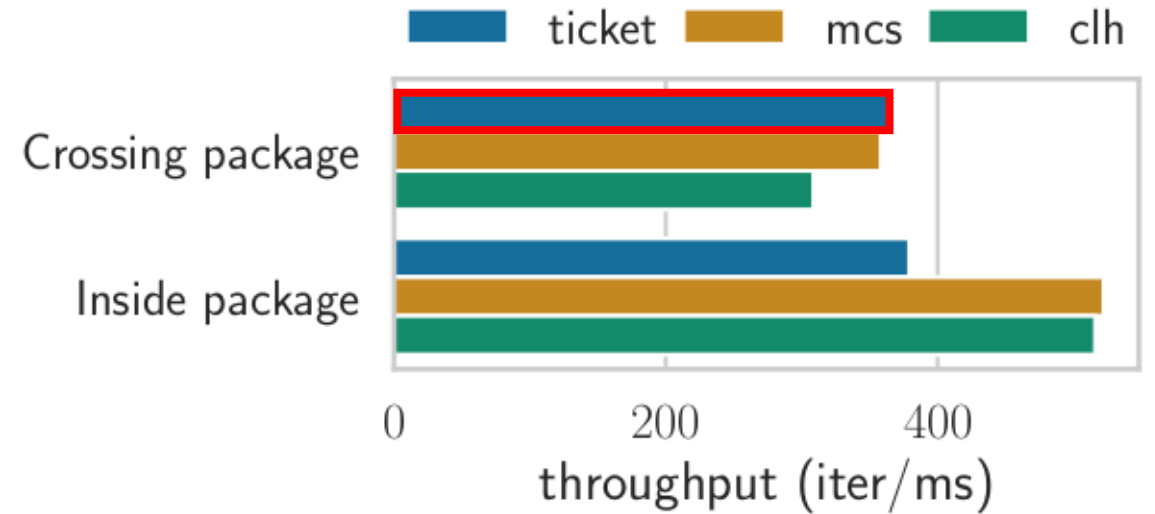
# Level-heterogeneity

Many-core NUMA system



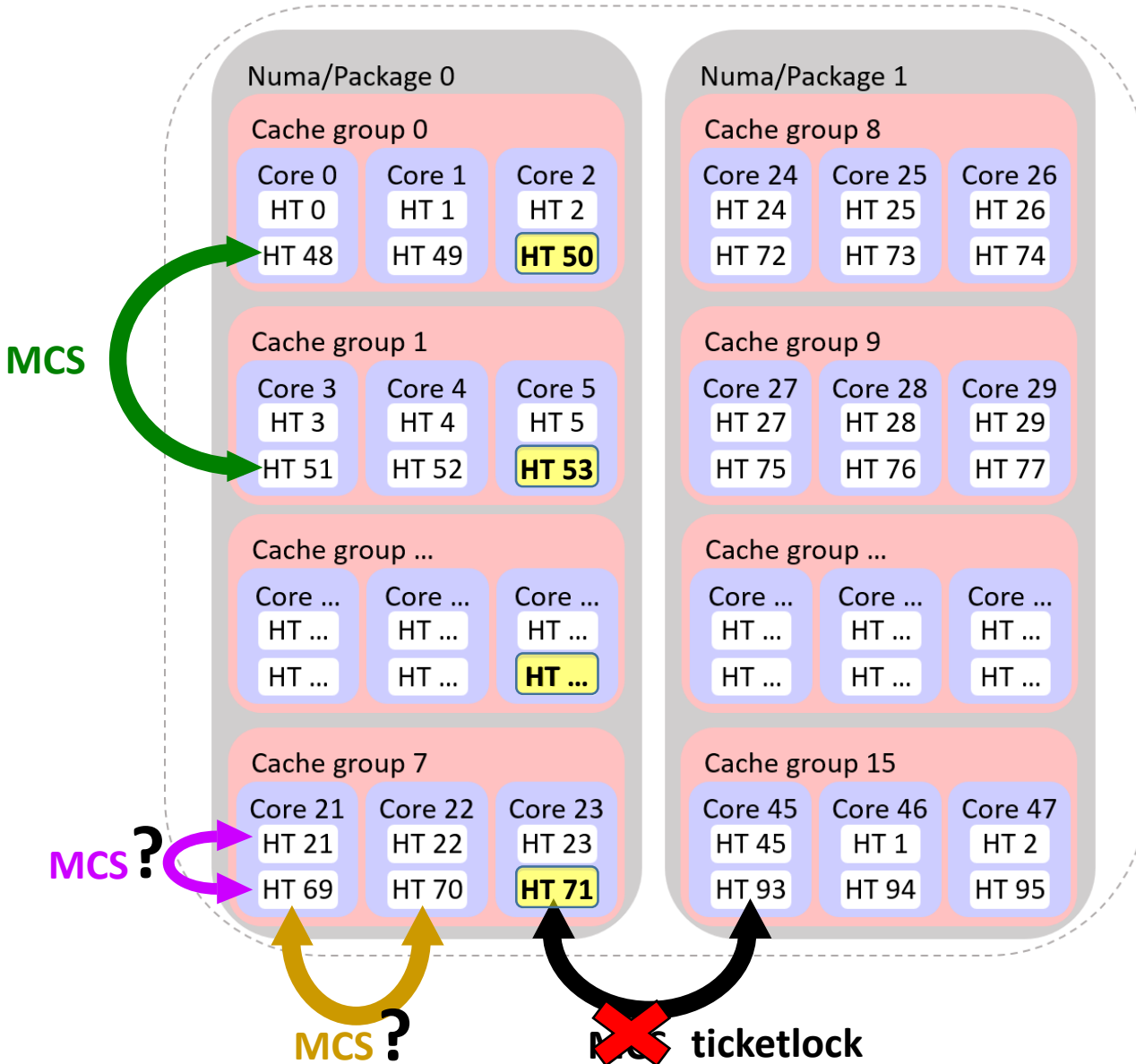
# Experiment - Heterogeneity

x86 server – execution of classical locks in isolation



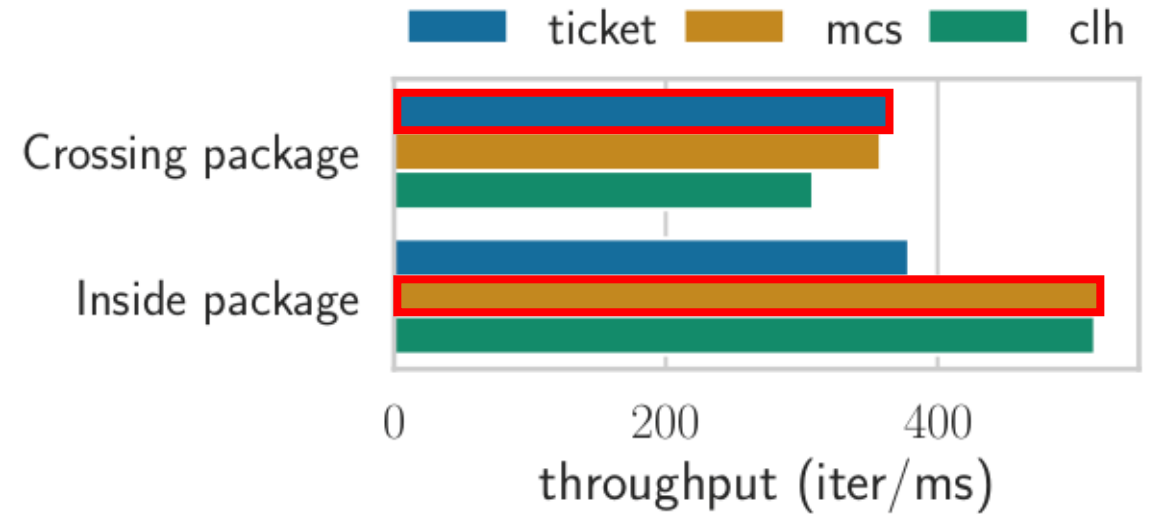
# Level-heterogeneity

Many-core NUMA system



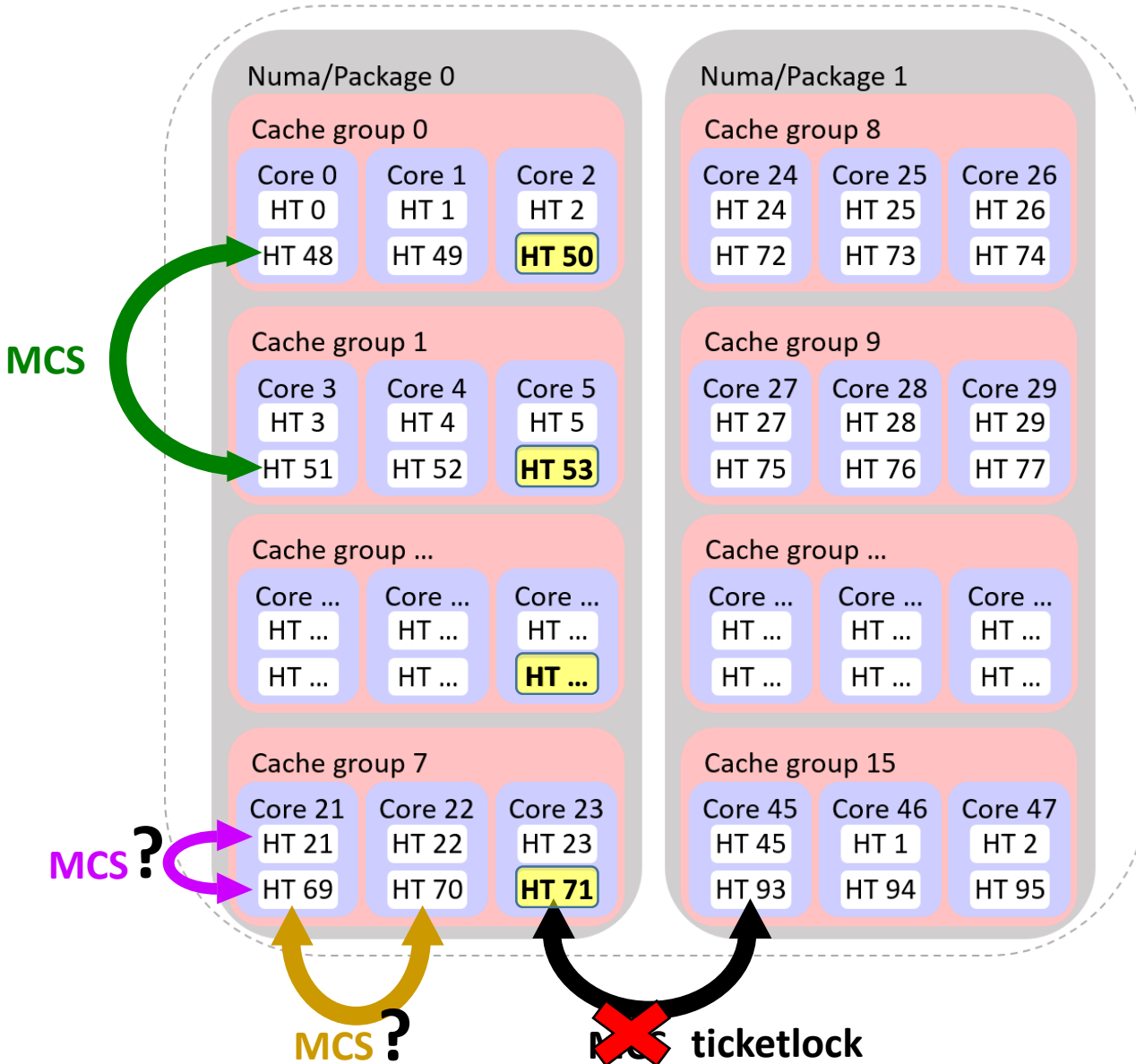
# Experiment - Heterogeneity

x86 server – execution of classical locks in isolation



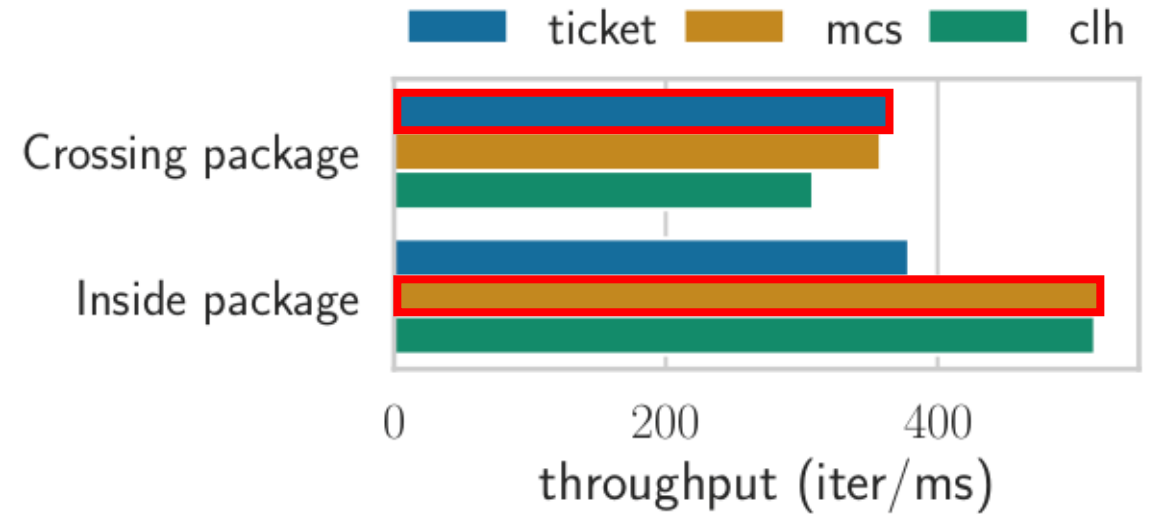
# Level-heterogeneity

Many-core NUMA system



# Experiment - Heterogeneity

x86 server – execution of classical locks in isolation

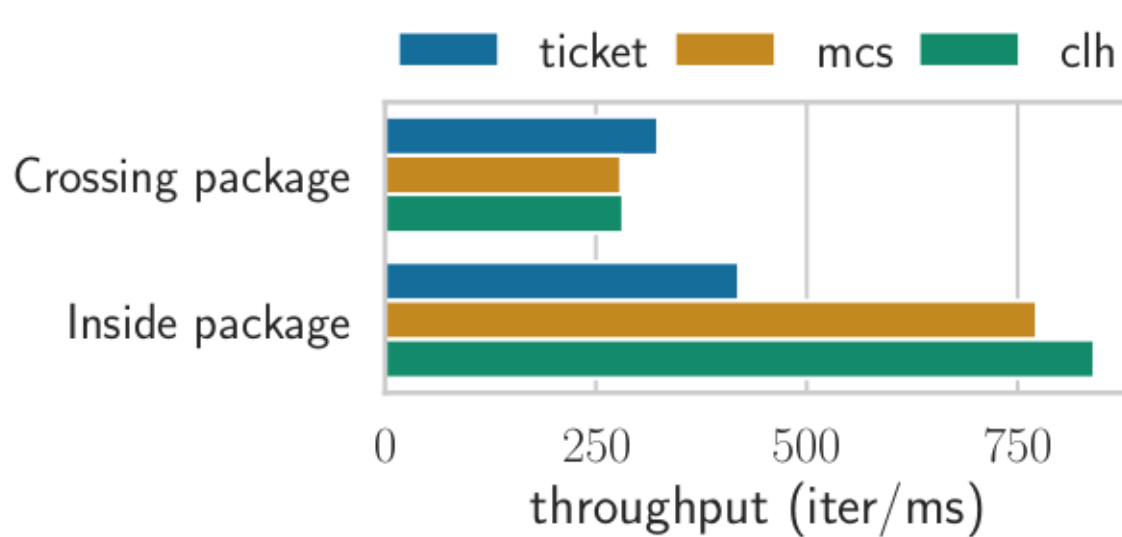


*Level-heterogeneity:*  
For different levels,  
the best lock may differ

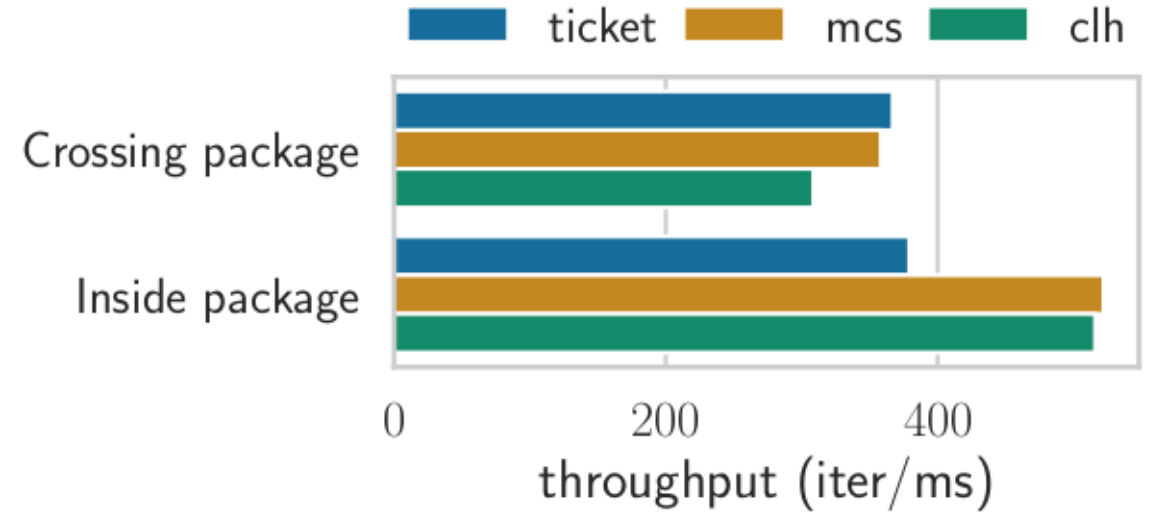
# Experiment – Platform Optimization

# Experiment - Heterogeneity

**Arm server** – execution of classical locks in isolation



**x86 server** – execution of classical locks in isolation

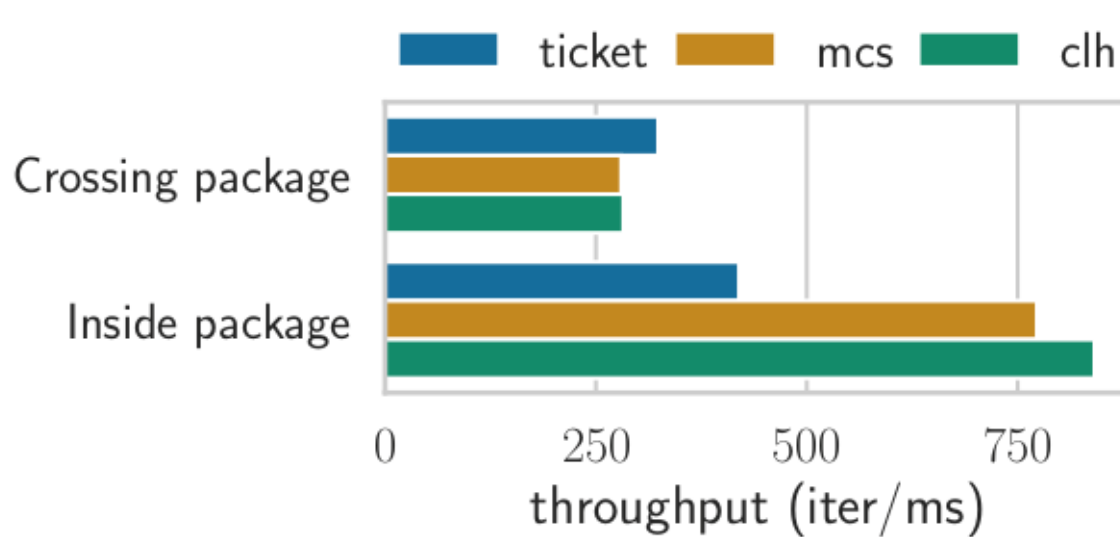


*Level-heterogeneity:*  
For different levels,  
the best lock may differ

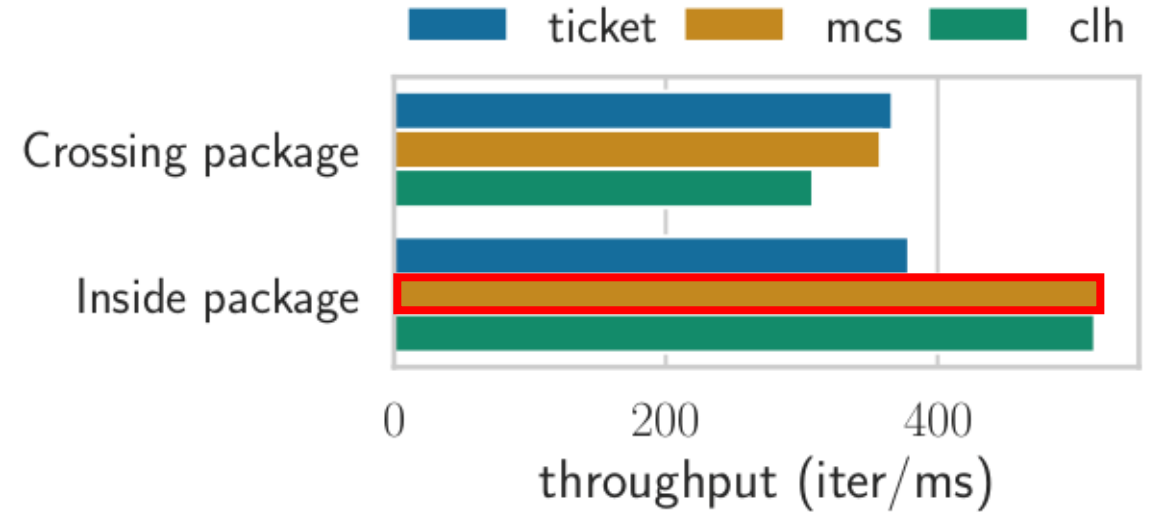
# Experiment – Platform Optimization

# Experiment - Heterogeneity

**Arm server** – execution of classical locks in isolation



**x86 server** – execution of classical locks in isolation

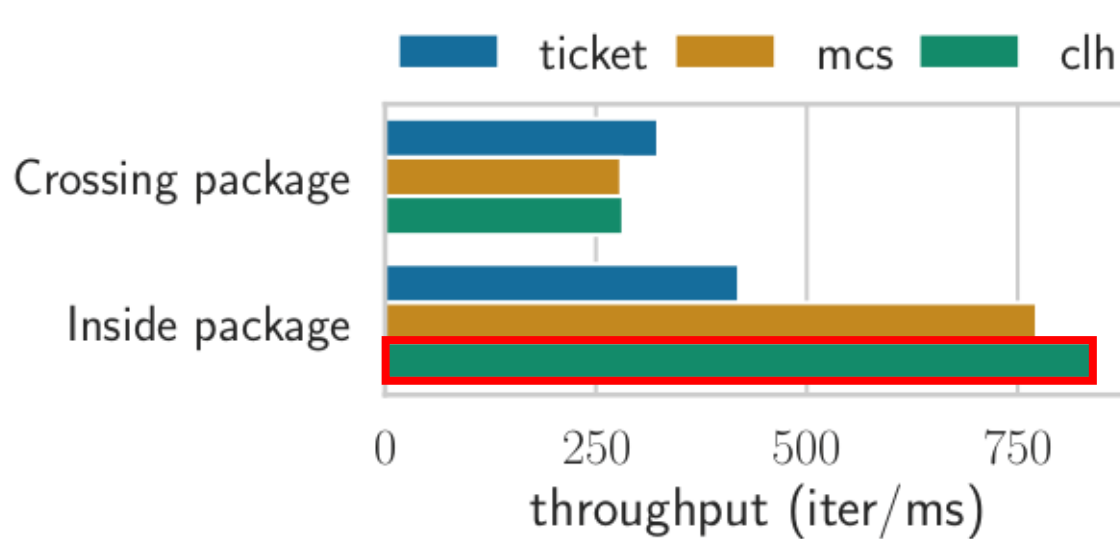


*Level-heterogeneity:*  
For different levels,  
the best lock may differ

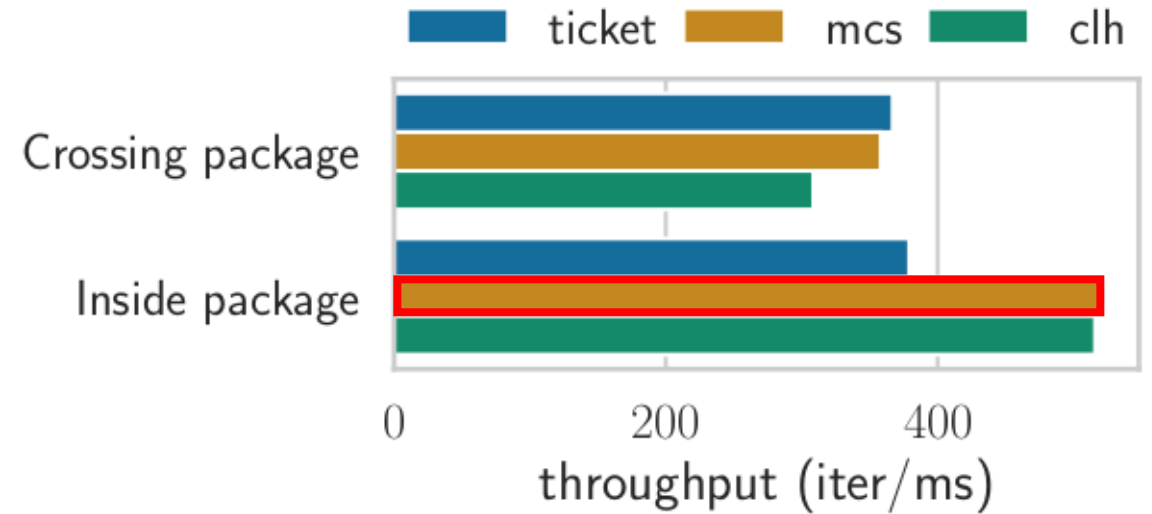
# Experiment – Platform Optimization

# Experiment - Heterogeneity

**Arm server** – execution of classical locks in isolation



**x86 server** – execution of classical locks in isolation



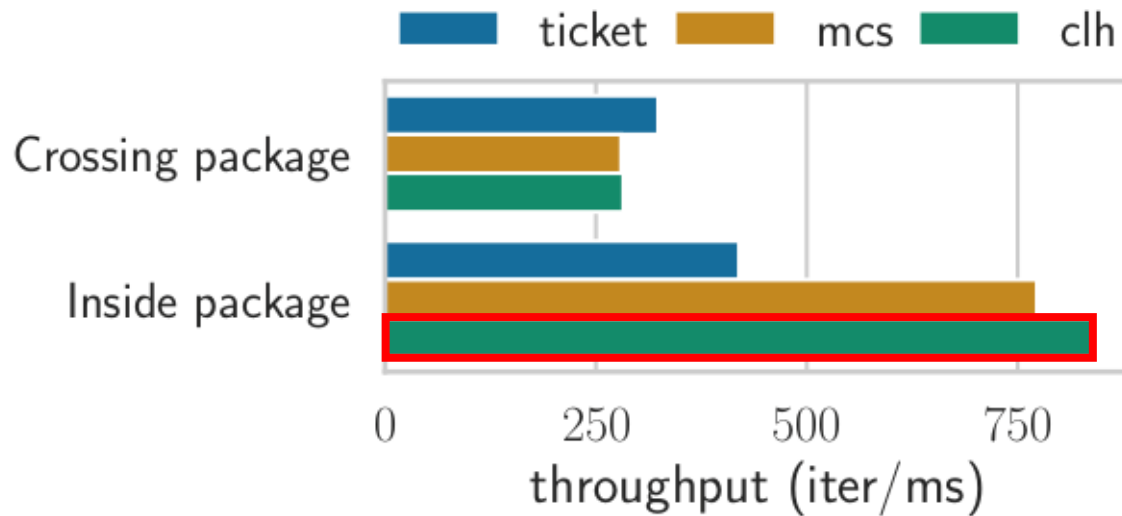
*Level-heterogeneity:*  
For different levels,  
the best lock may differ



# Experiment – Platform Optimization

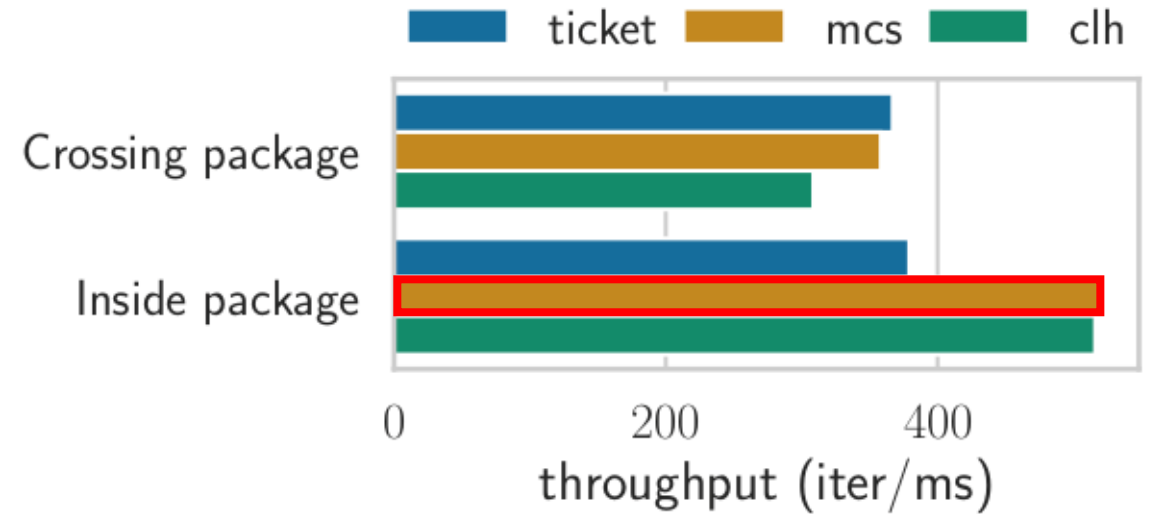
# Experiment - Heterogeneity

**Arm server** – execution of classical locks in isolation



*Platform Optimization:*  
For different platforms,  
the best lock for a level may differ

**x86 server** – execution of classical locks in isolation

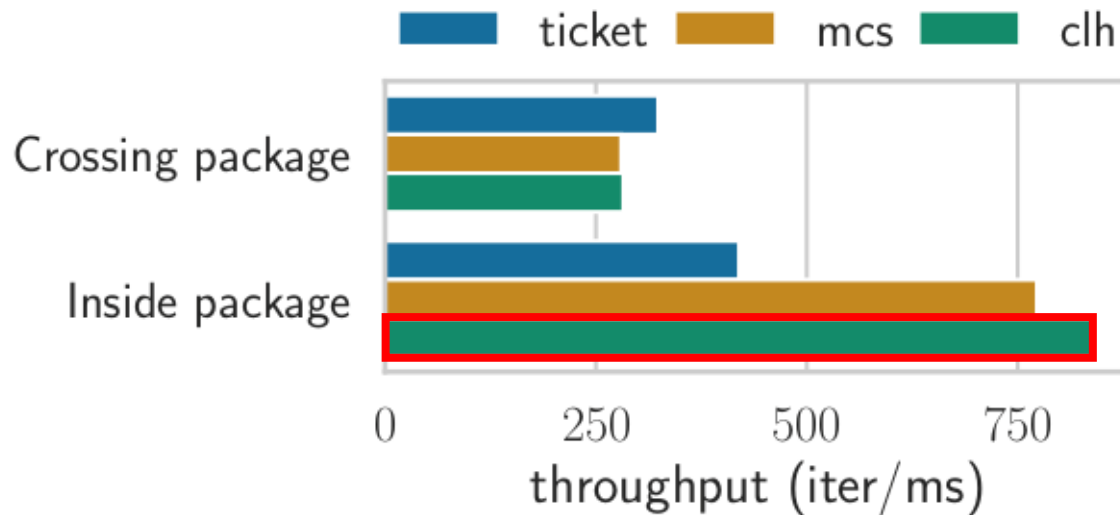


*Level-heterogeneity:*  
For different levels,  
the best lock may differ

# Experiment – Platform Optimization

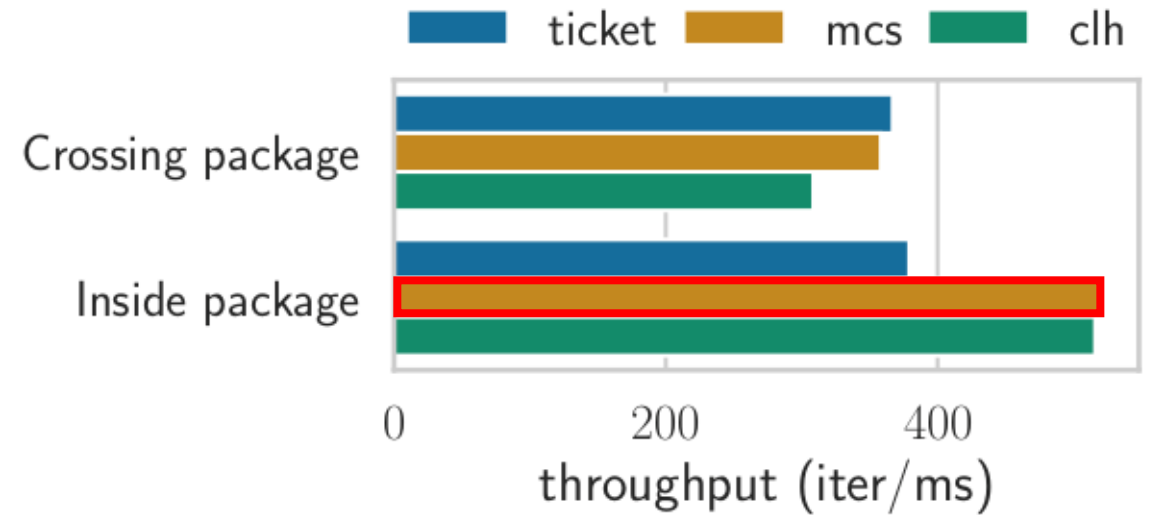
# Experiment - Heterogeneity

**Arm server** – execution of classical locks in isolation



*Platform Optimization:*  
For different platforms,  
the best lock for a level may differ

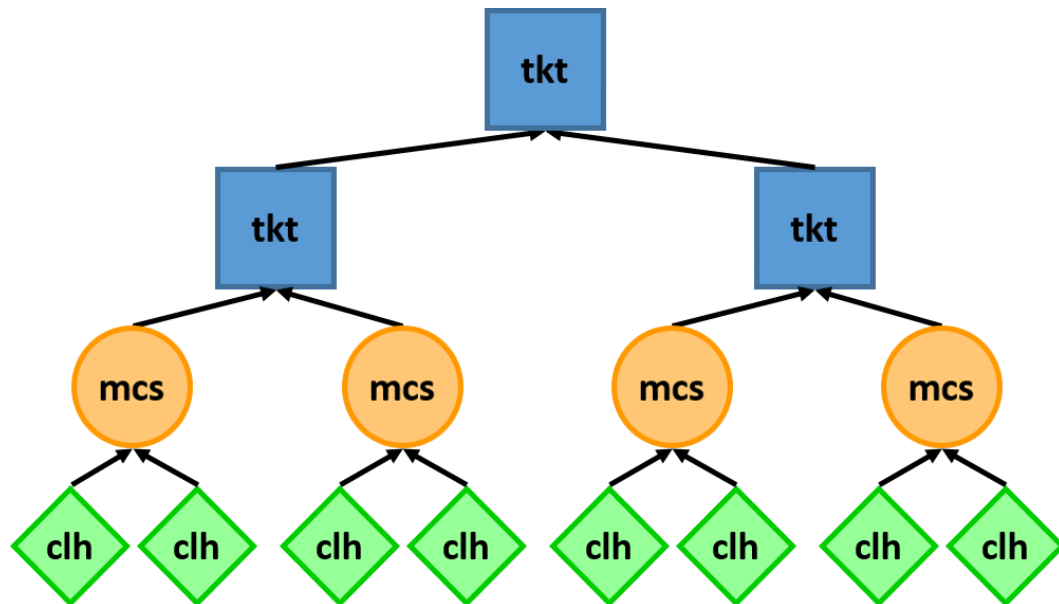
**x86 server** – execution of classical locks in isolation



*Level-heterogeneity:*  
For different levels,  
the best lock may differ

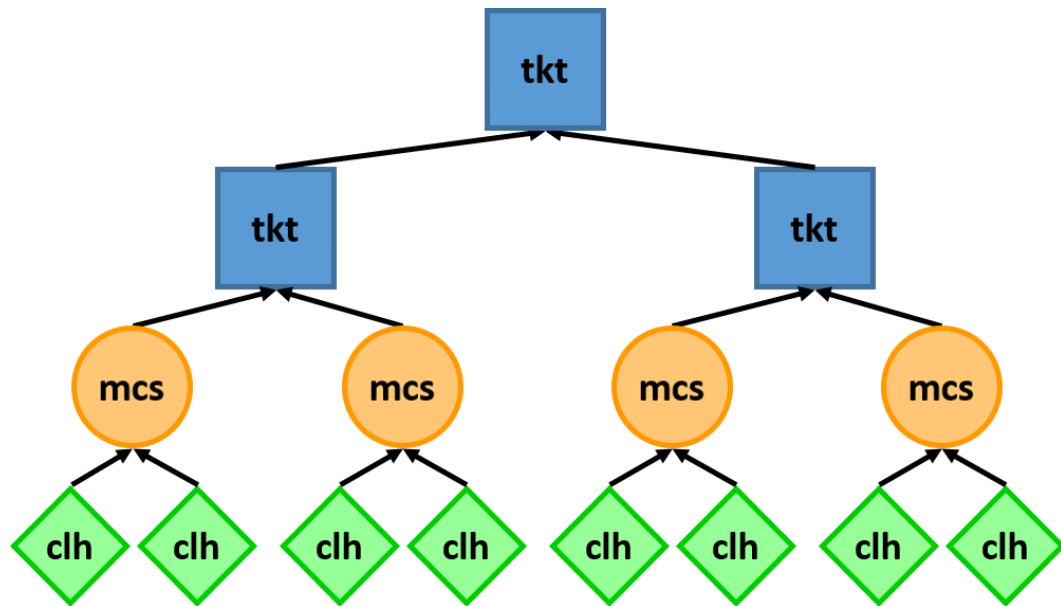
So, how does this affects our lock design?

# Our Desired NUMA-aware lock



- Multi-Level
- Level-Heterogeneous
- Configurability for Platform Optimization

# Our Desired NUMA-aware lock



- Multi-Level
- Level-Heterogeneous
- Configurability for Platform Optimization
- Showing lock correctness is challenging
  - Weak Memory Models (WMMs) make it even more complicated

## Our contribution: CLoF

We propose CLoF, a framework to generate locks for a target platform:

- that support an arbitrary hierarchy;
- for each level, the lock implementation may be different;
- that are *correct-by-construction* on Weak Memory Models.

# Our contribution: CLoF

We propose CLoF, a framework to generate locks for a target platform:

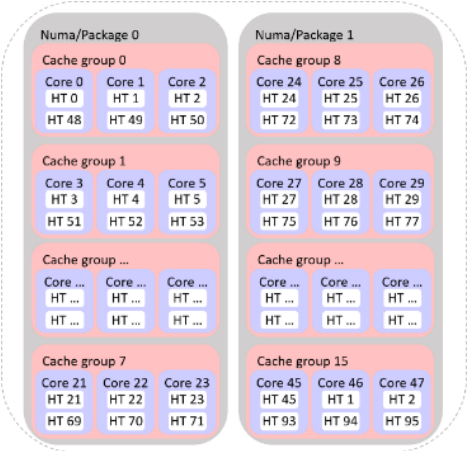
- that support an arbitrary hierarchy;
- for each level, the lock implementation may be different;
- that are *correct-by-construction* on Weak Memory Models.

NUMA-aware locks		Correctness on WMMs	Level heterogeneity & Architecture optimization	Multi-Level
lock cohorting	PPPoPP'12	✗	✓	✗
HMCS	PPoPP'15	✗ <sup>1</sup>	✗	✓
CNA lock	EuroSys'19	✗	✗	✗
ShflLock	SOSP'19	✗	✗	✗
CLoF	SOSP'21	✓	✓	✓

<sup>1</sup>Insufficient barriers, fixed in Oberhauser *et al.*, Verifying and Optimizing the HMCS Lock for Arm Servers, NETYS'2021

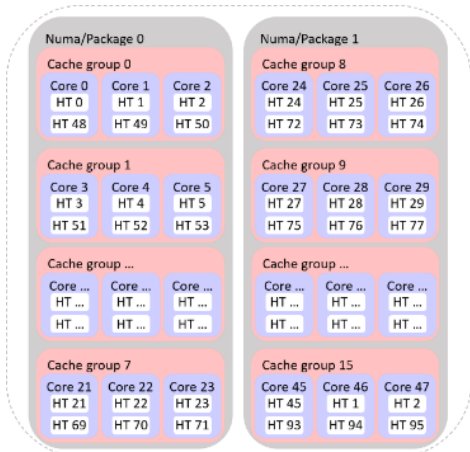
# CLoF Workflow

## Discover Memory Hierarchy



# CLoF Workflow

## Discover Memory Hierarchy



## Verify correctness on WMMs

tkl

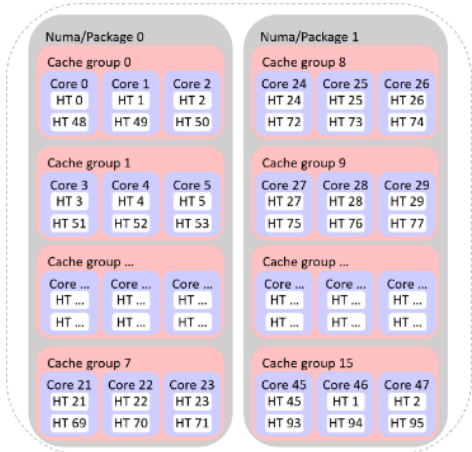
clh

mcs

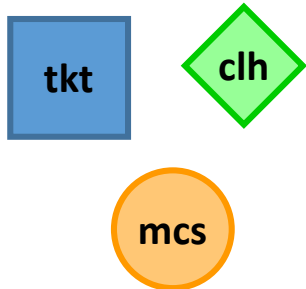


# CLoF Workflow

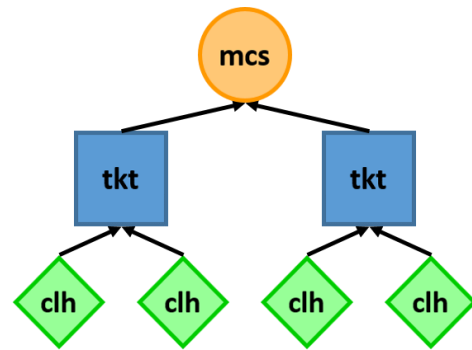
Discover Memory Hierarchy



Verify correctness on WMMs

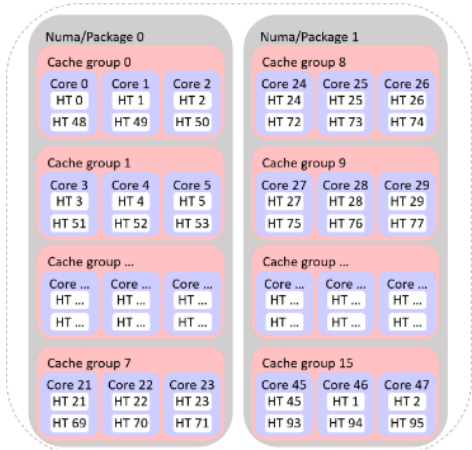


CLoF Lock Generator

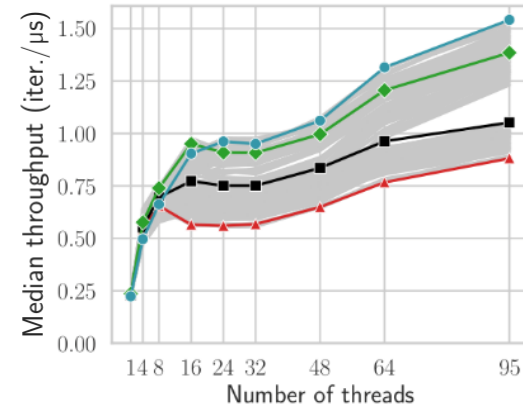
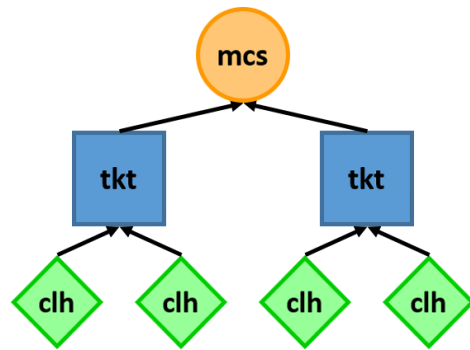
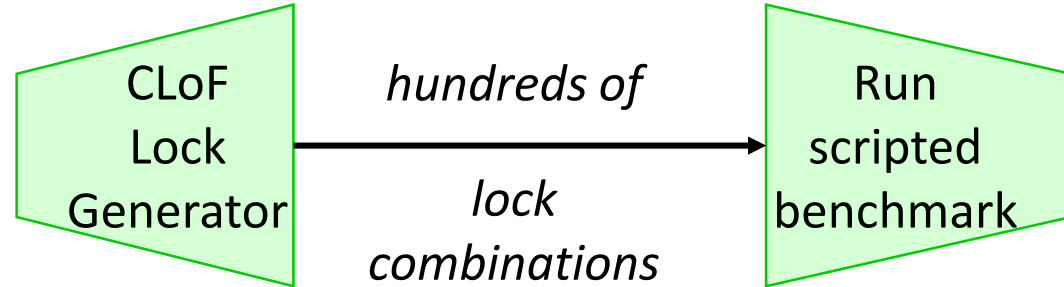
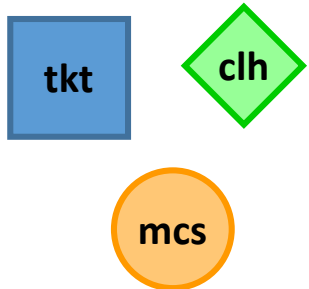


# CLoF Workflow

Discover Memory Hierarchy

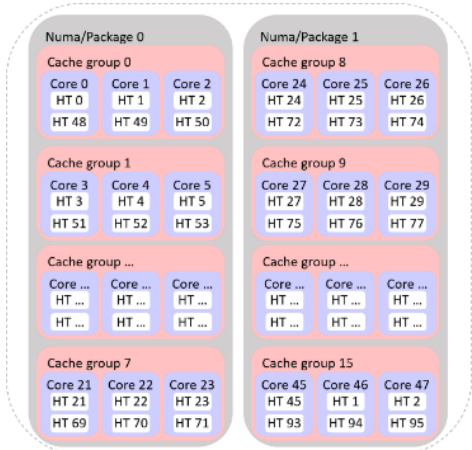


Verify correctness on WMMs

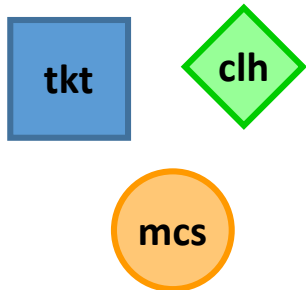


# CLoF Workflow

Discover Memory Hierarchy



Verify correctness on WMMs

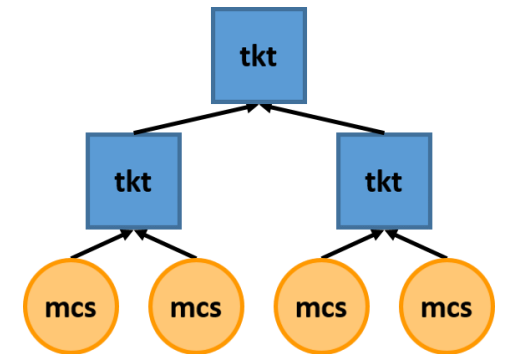
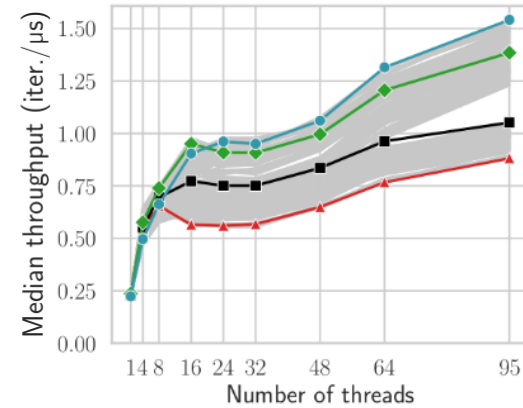
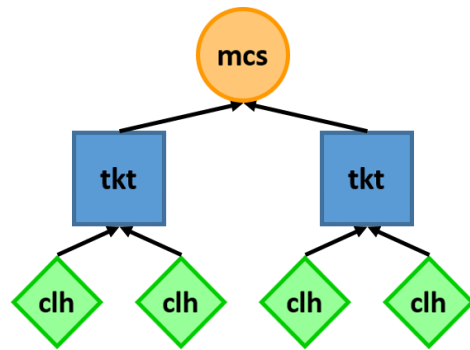


CLoF Lock Generator

*hundreds of lock combinations*

Run scripted benchmark

Select Best Lock



# Agenda

- How to figure out the hierarchy we need to use?
- How does our CLoF Lock Generator works?
- How do we know it is correct?
- How do we pick the best lock for the target platform?

# Agenda

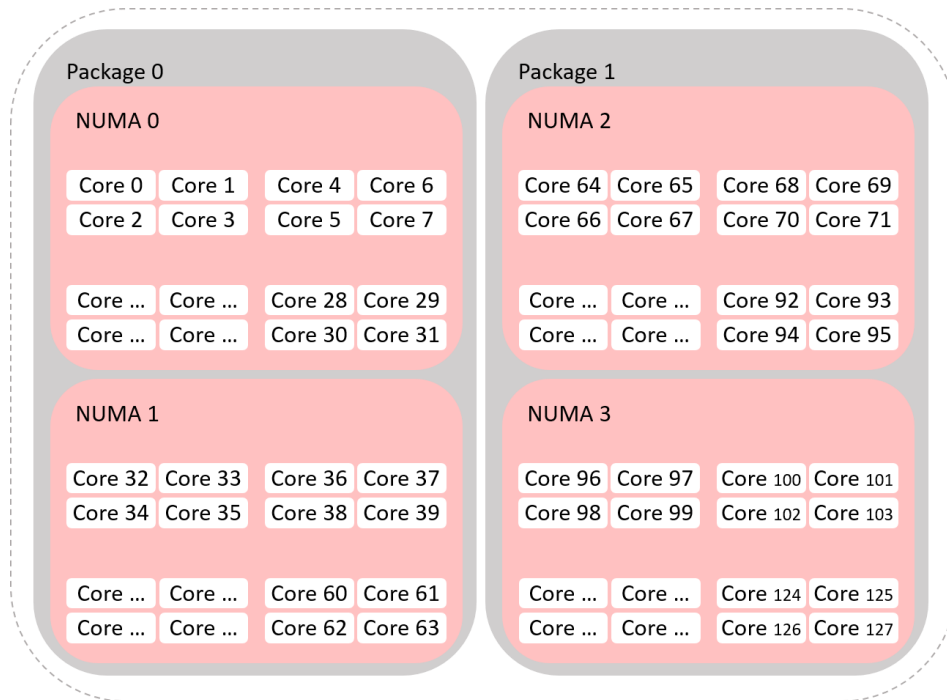
- How to figure out the hierarchy we need to use?
- How does our CLoF Lock Generator works?
- How do we know it is correct?
- How do we pick the best lock for the target platform?

# Discovering the Memory Hierarchy

Operating systems know the hierarchy

- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8

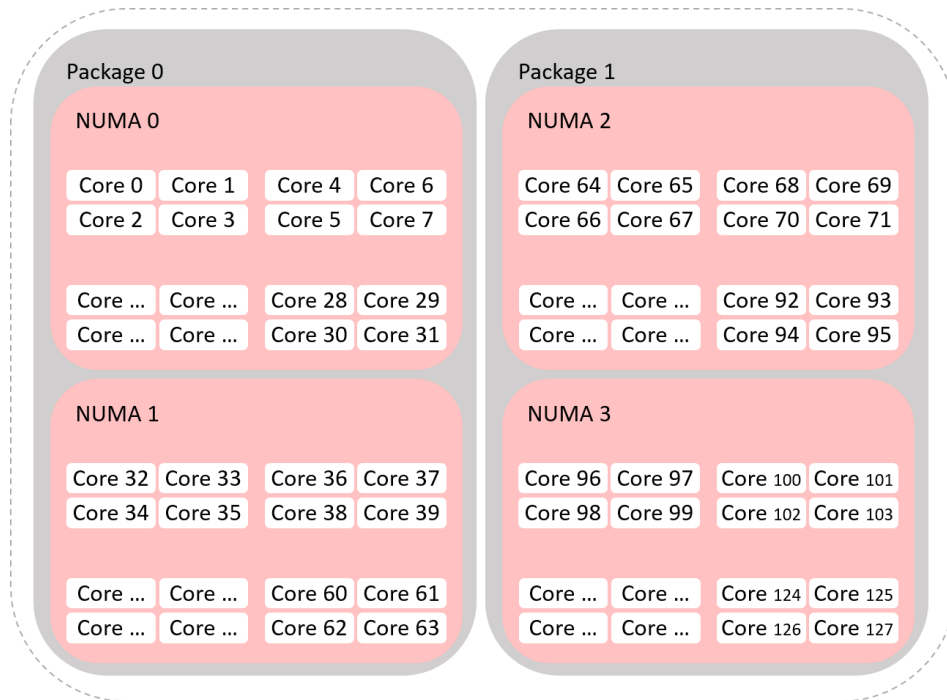


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



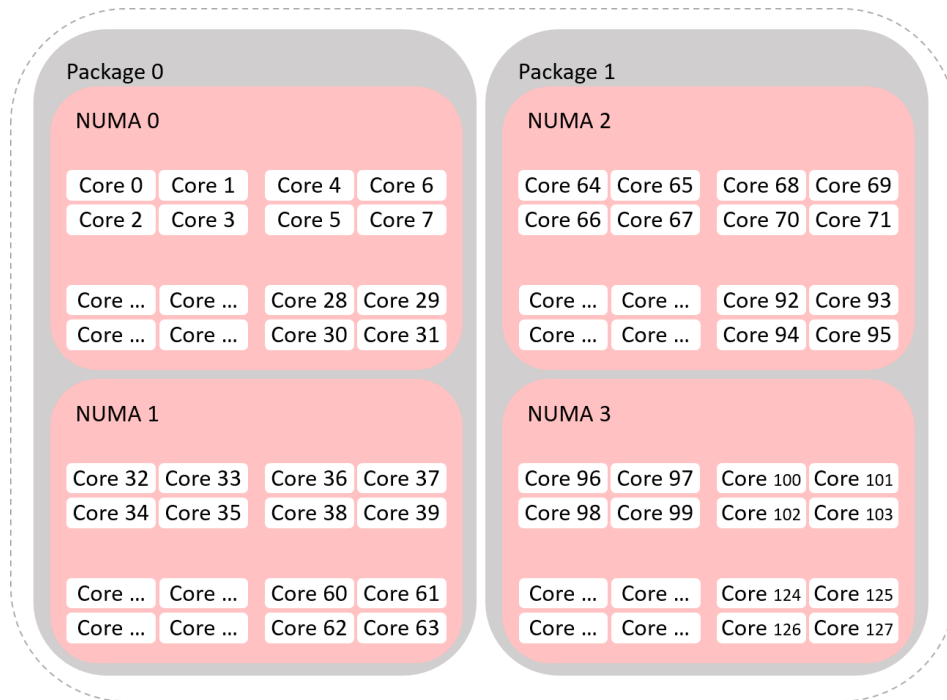
Is this the **full** hierarchy?

# Discovering the Memory Hierarchy

Operating systems know the hierarchy

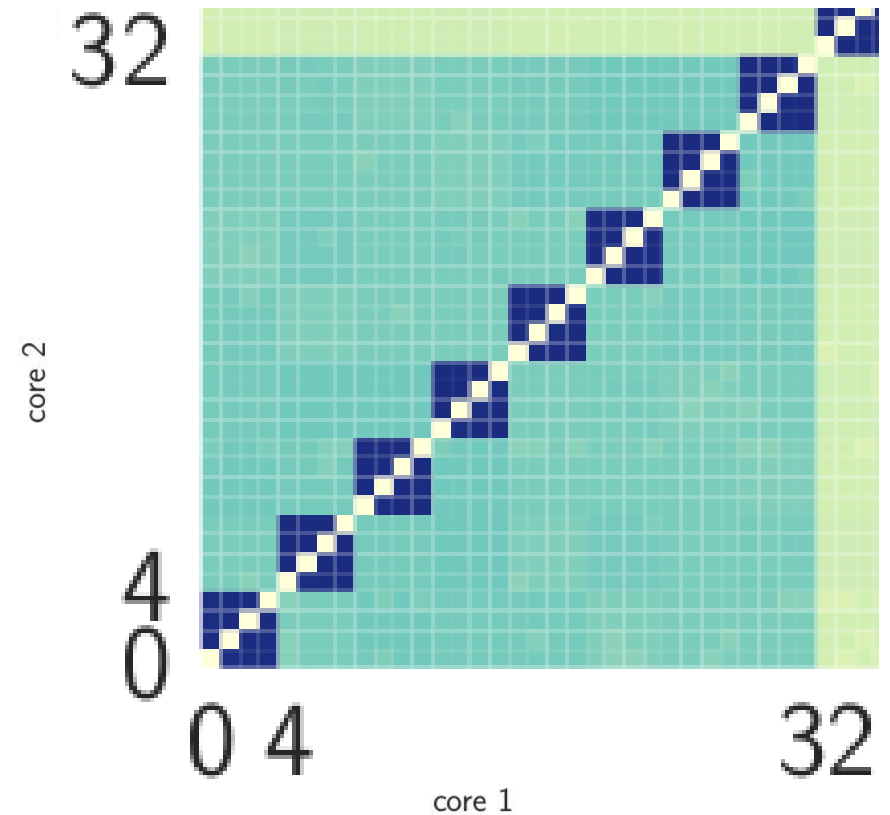
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput



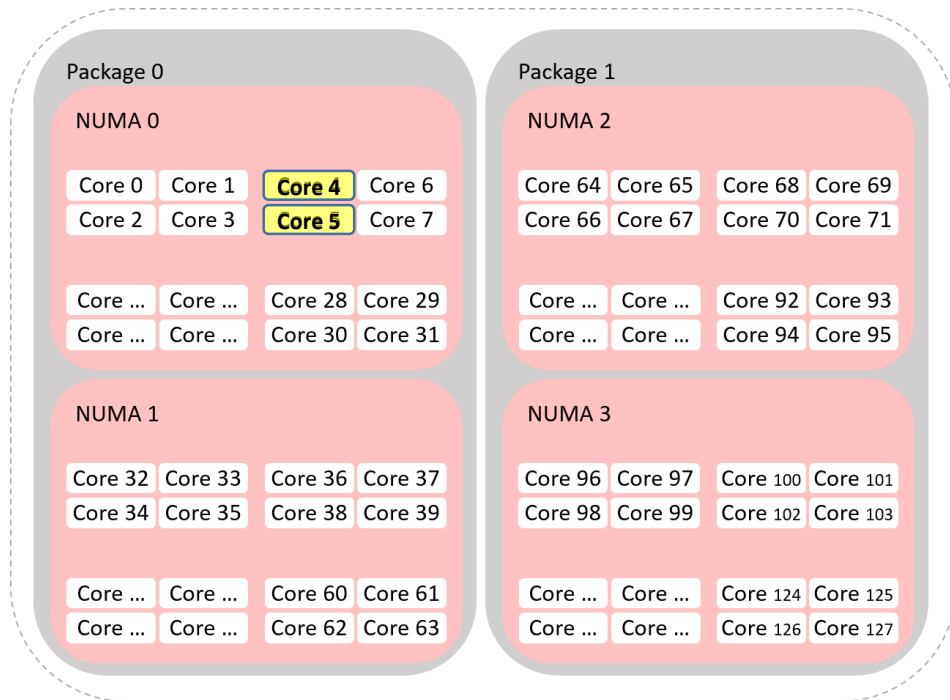


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

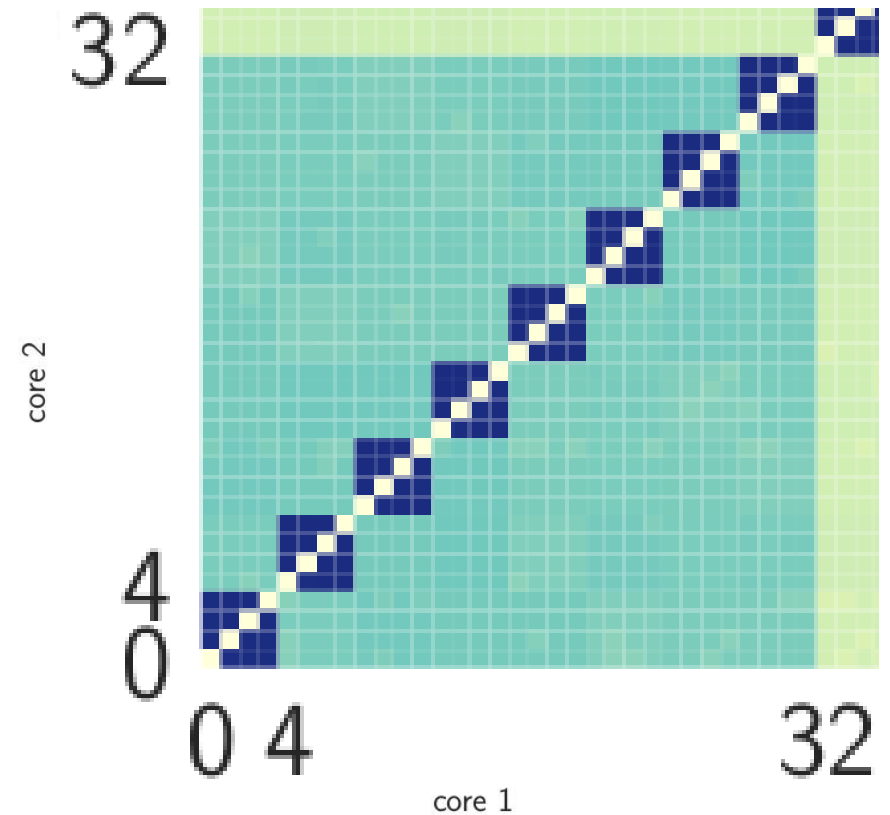
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

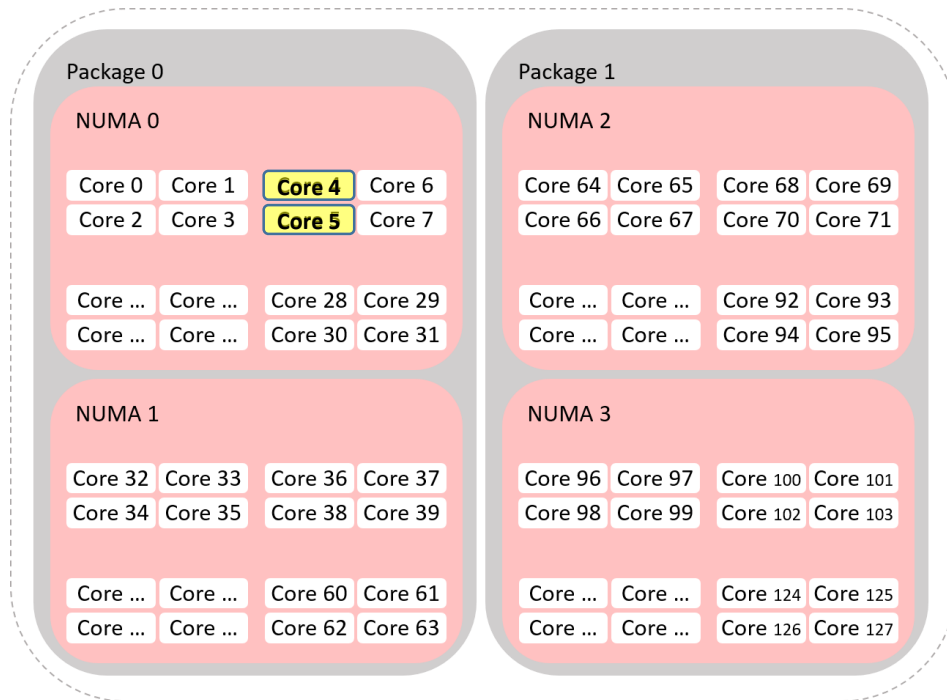


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

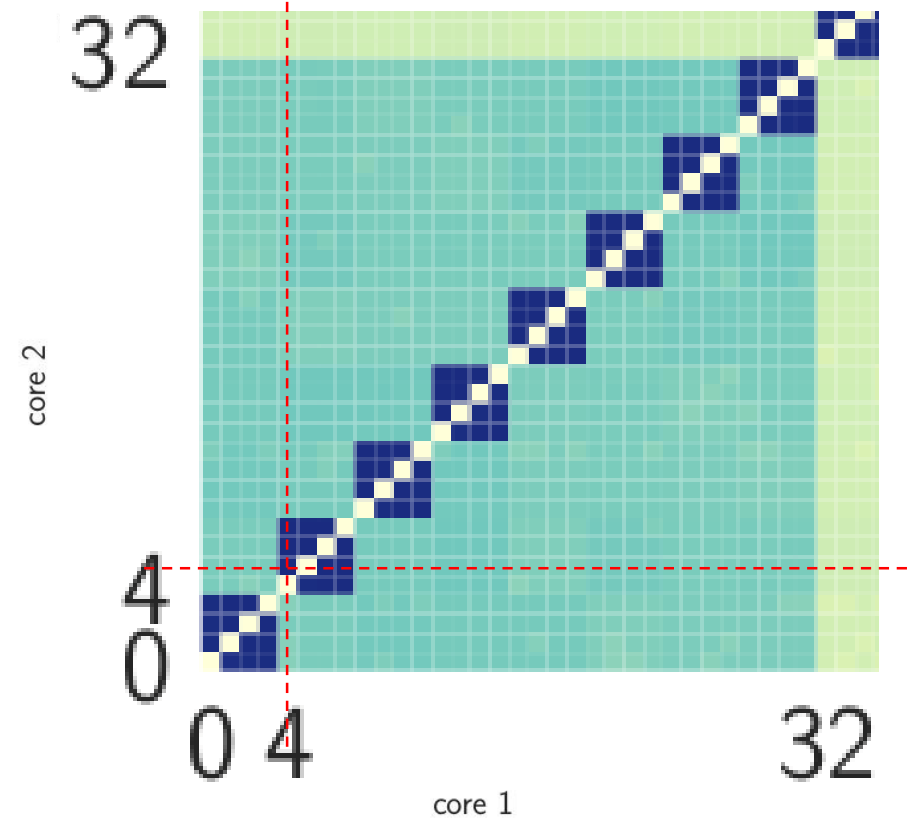
- Visible with Linux's lscpu and lstopo

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

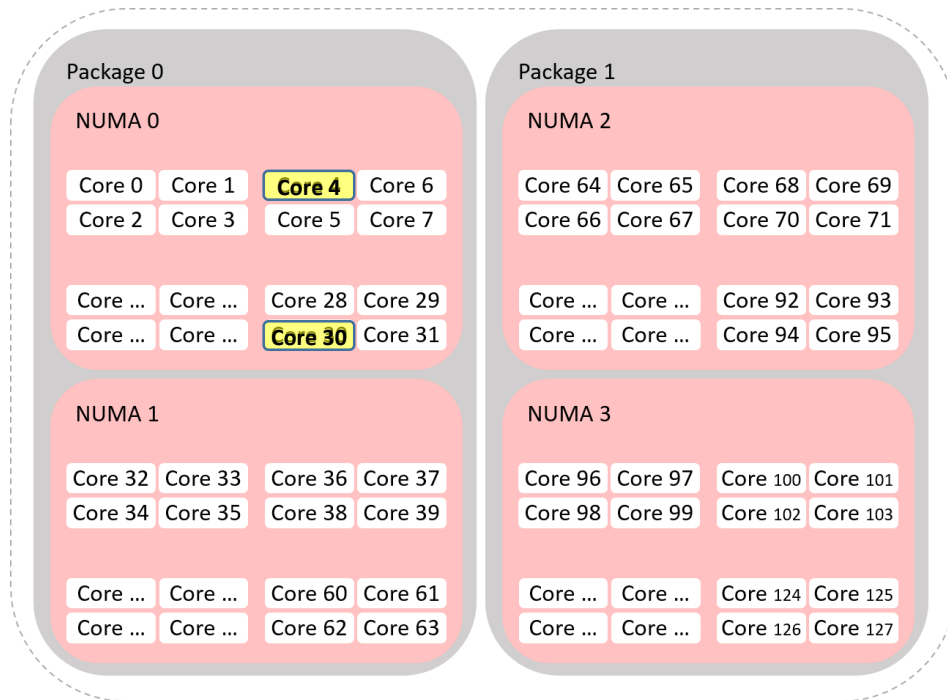


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

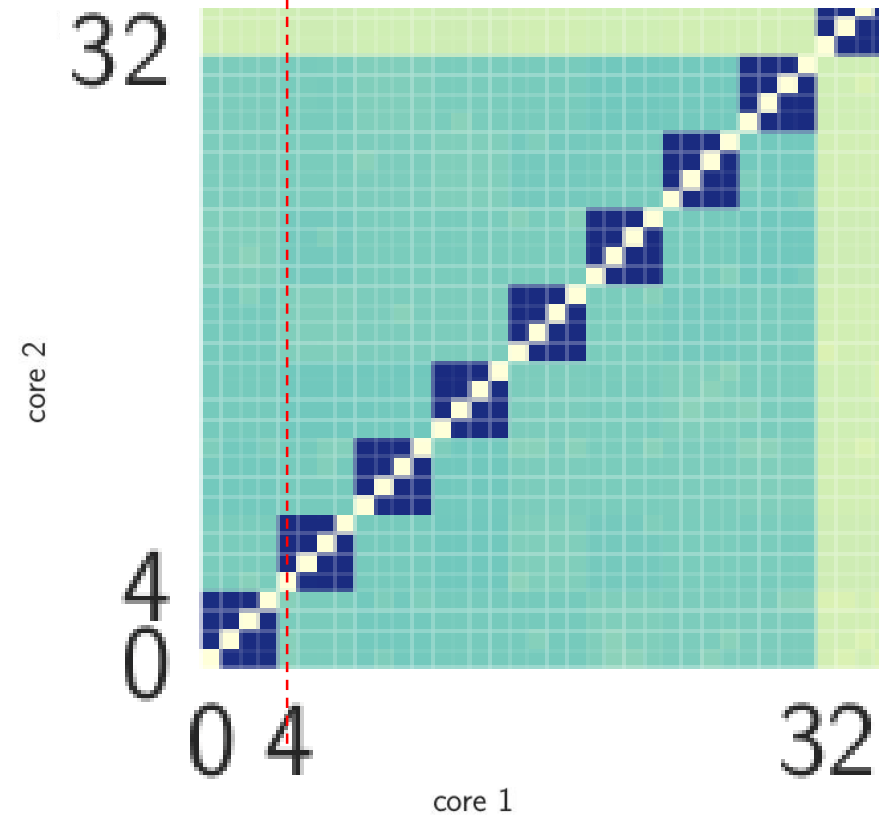
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

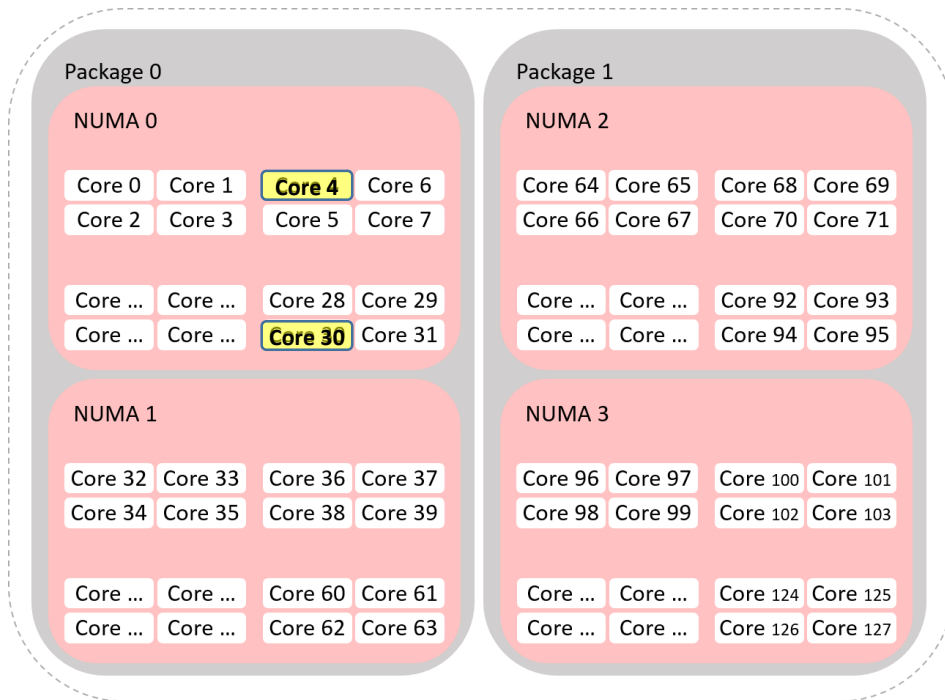


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

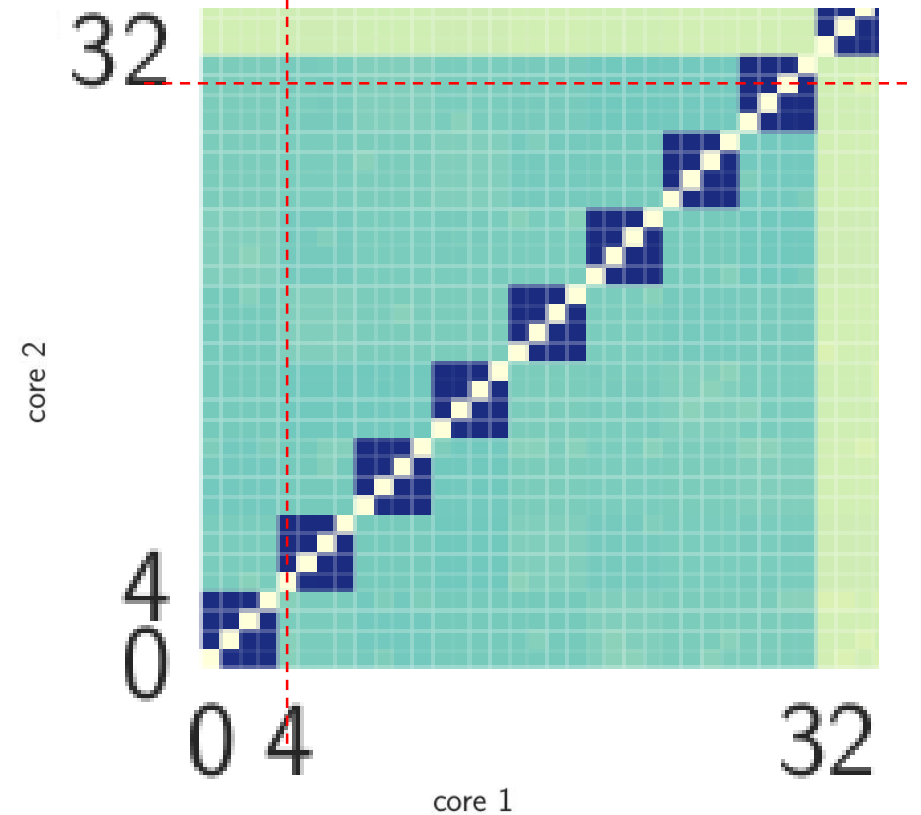
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

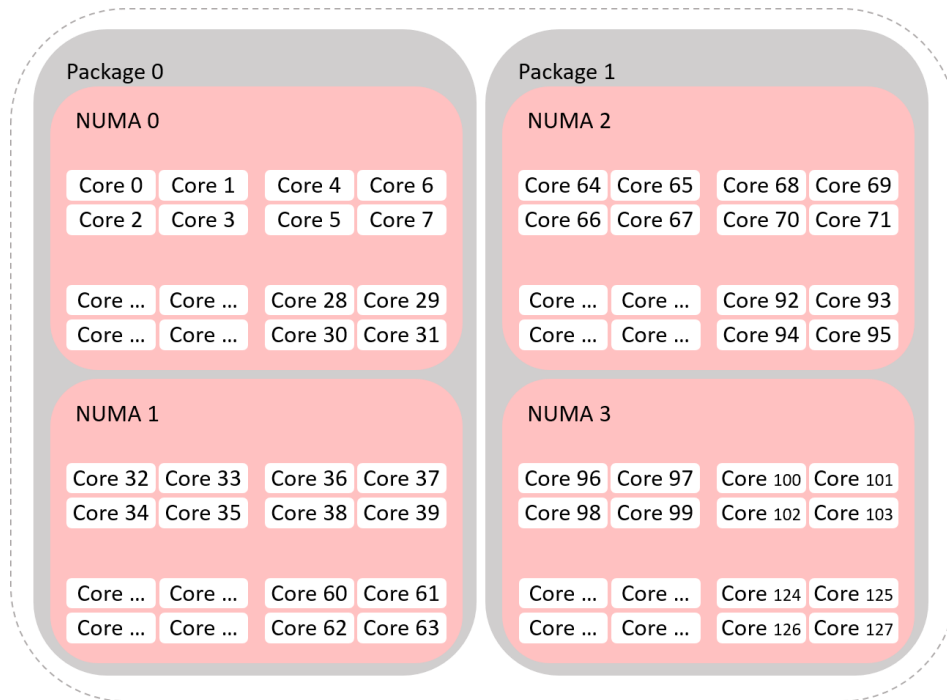


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

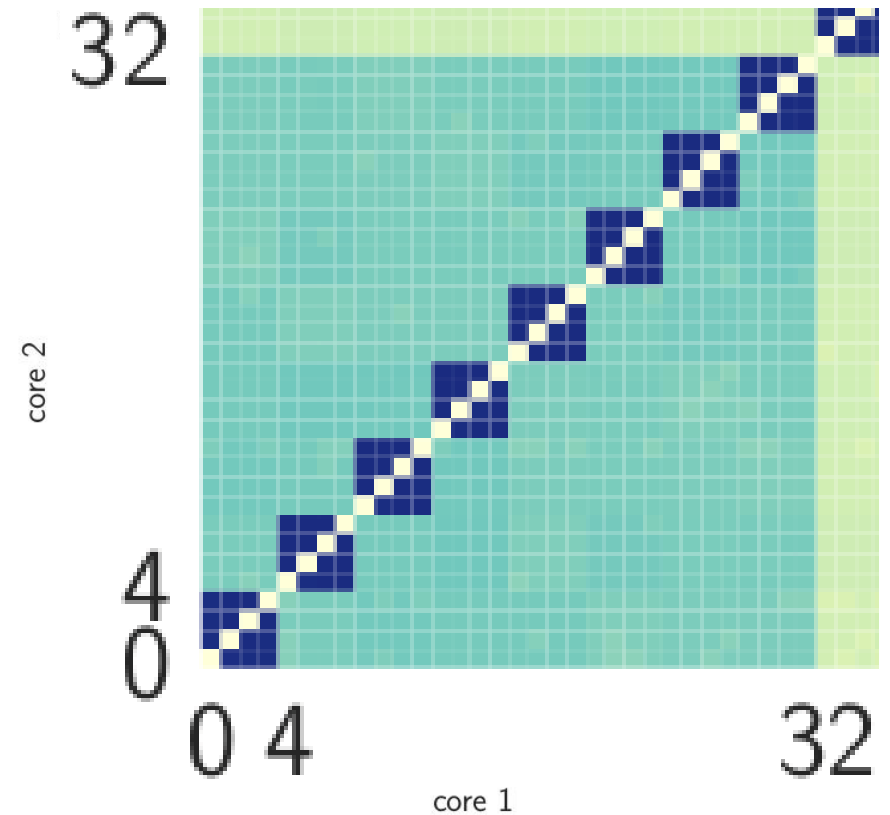
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

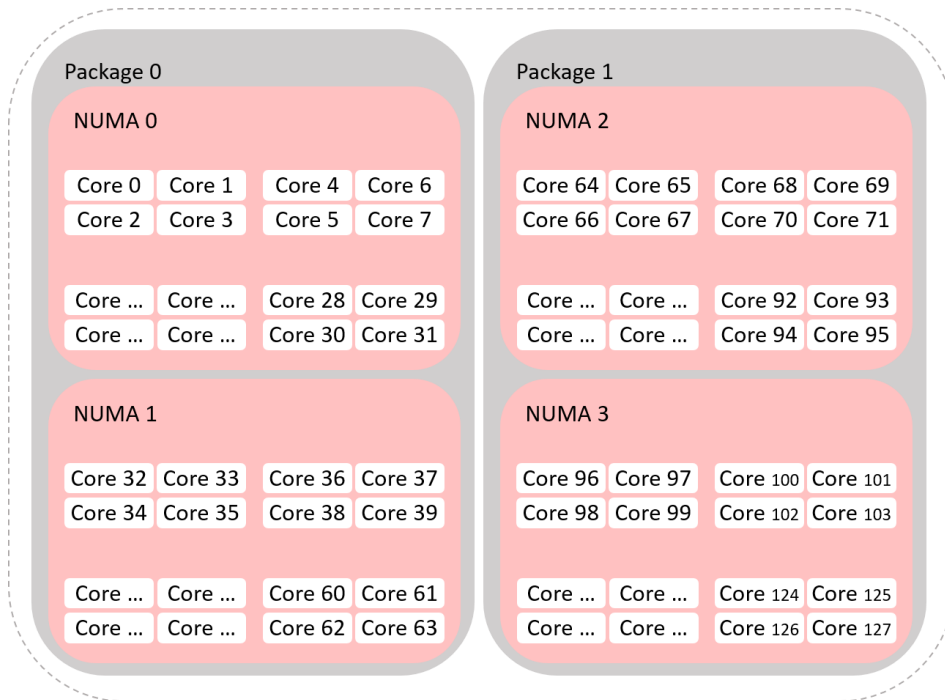


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

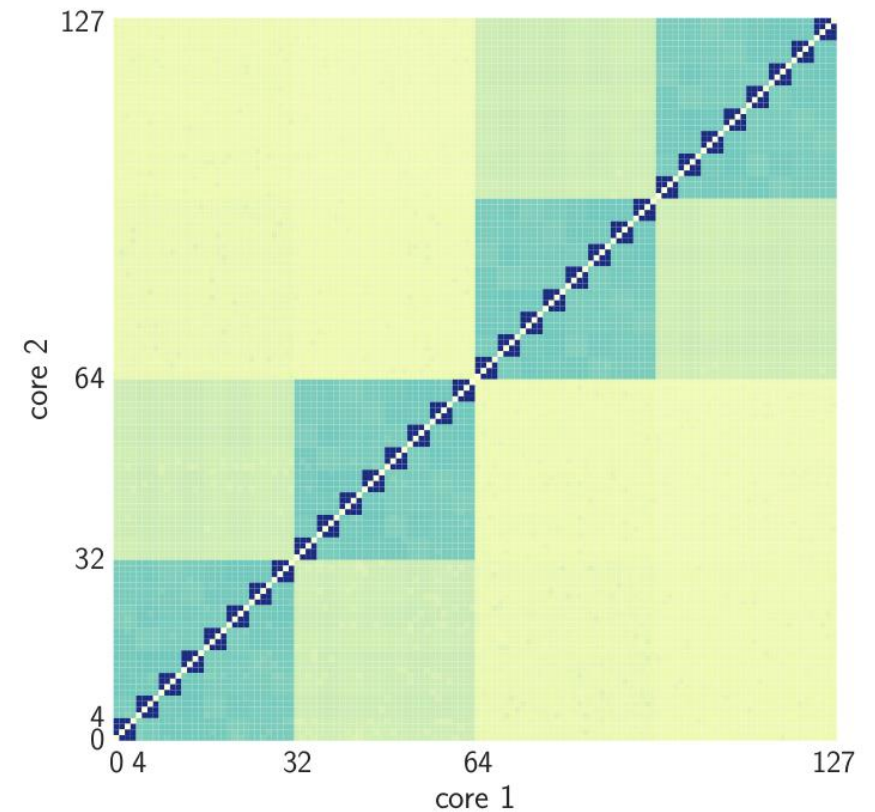
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

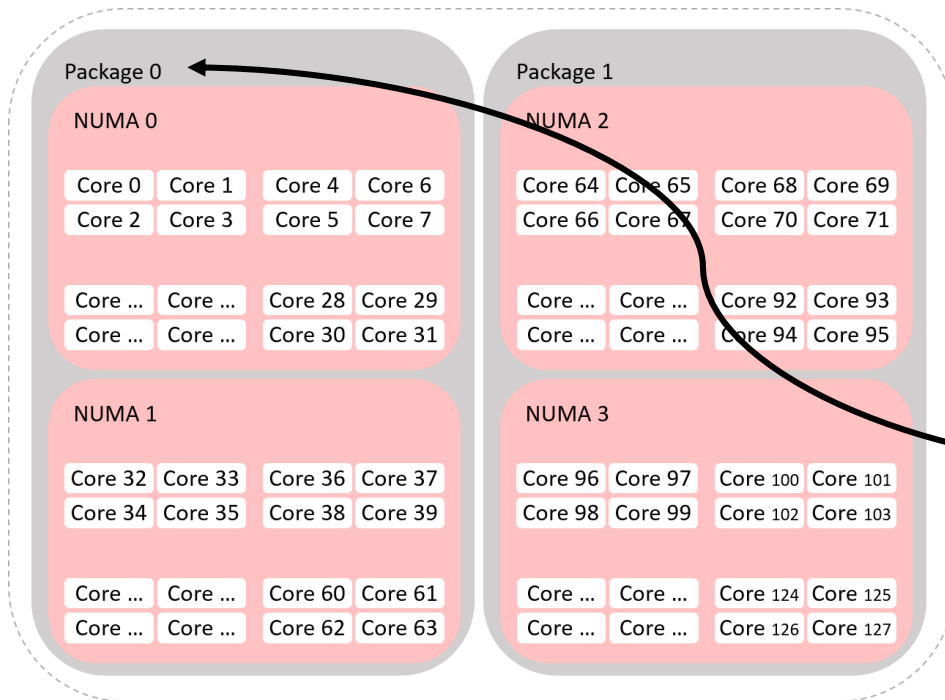


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

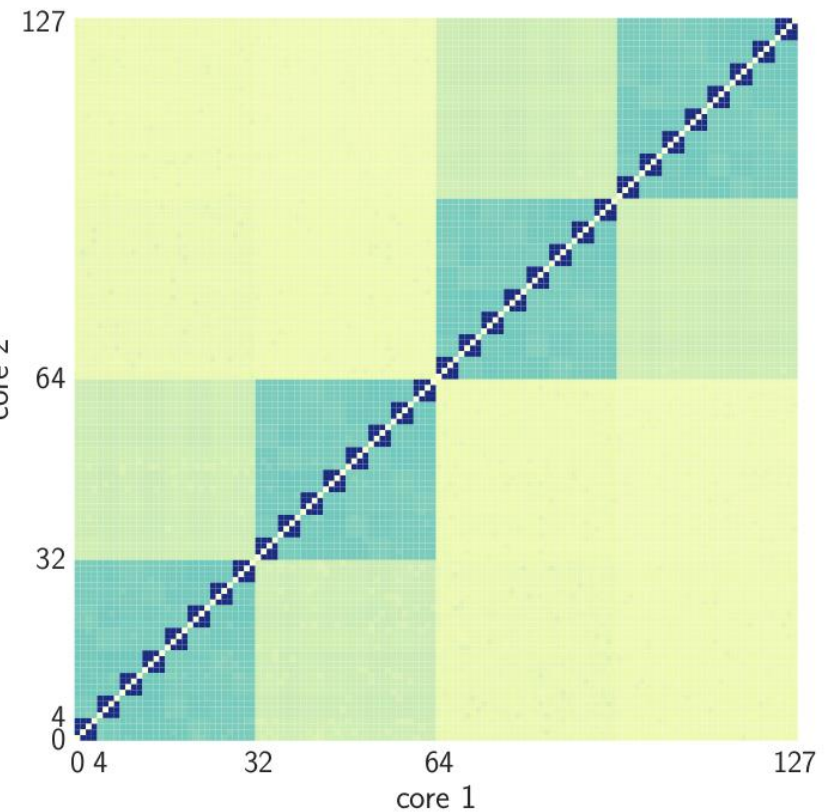
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput



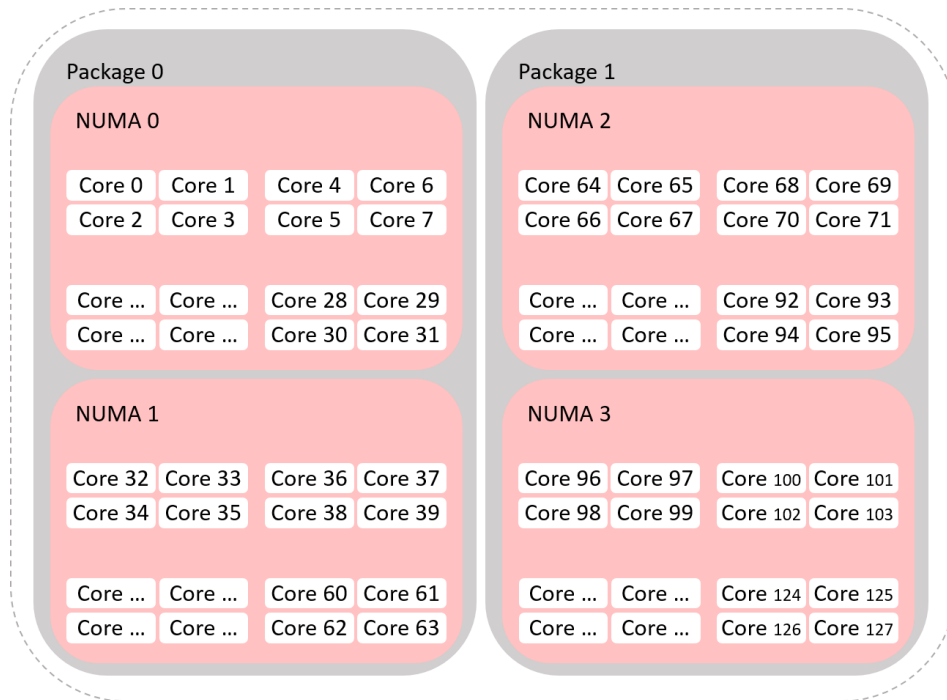


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

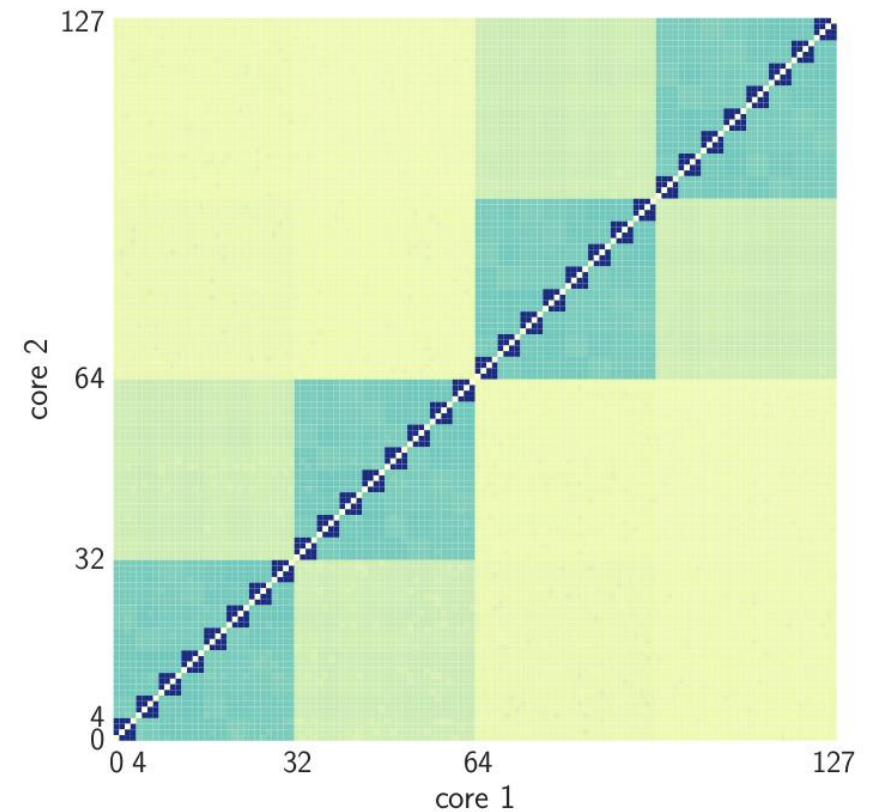
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput



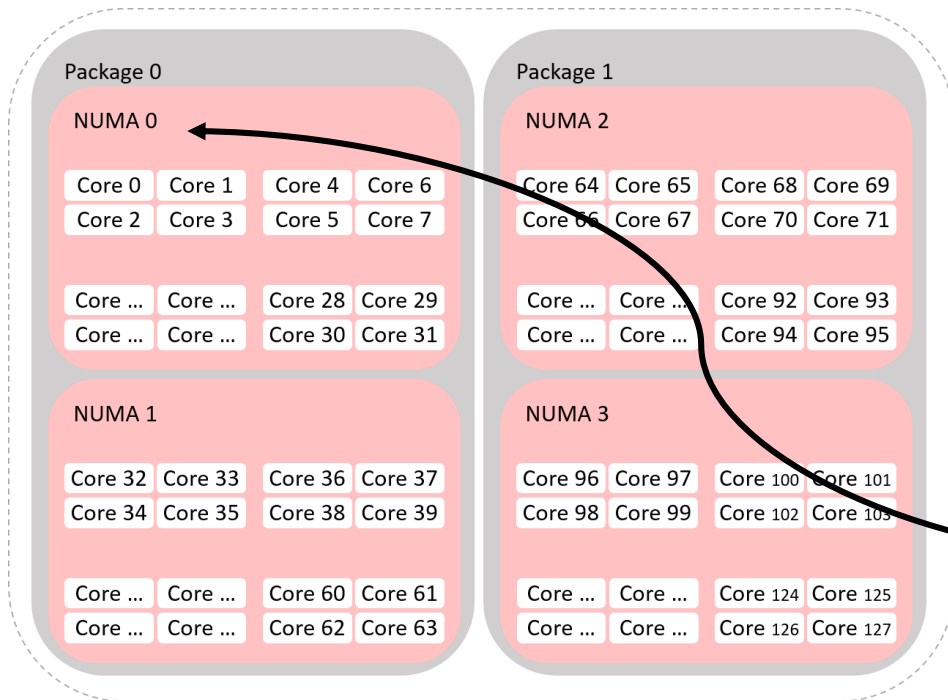


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

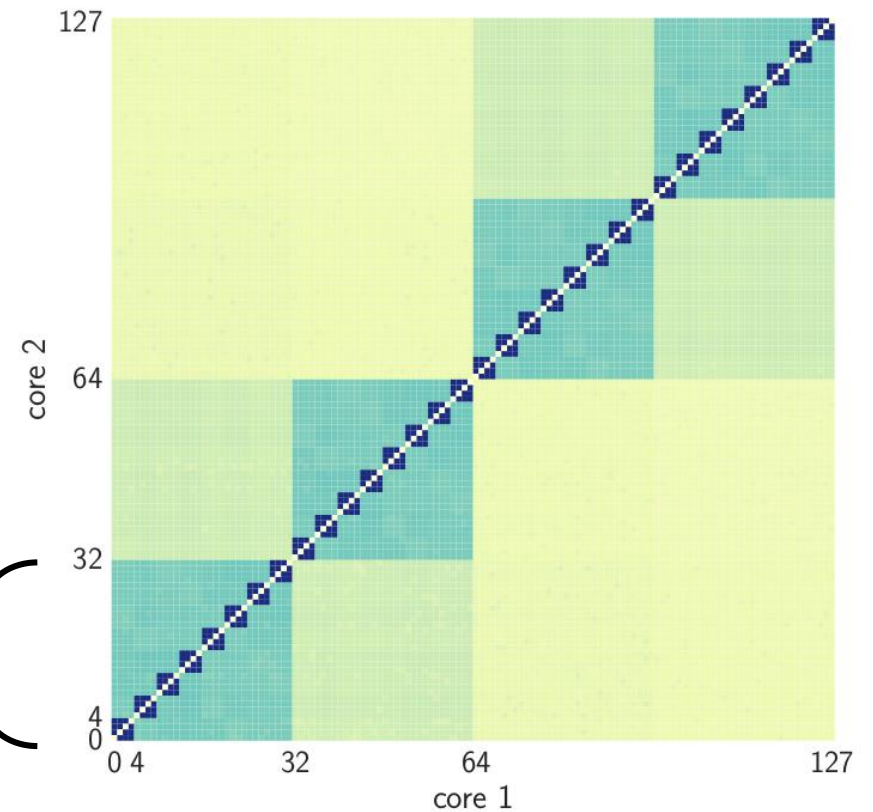
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

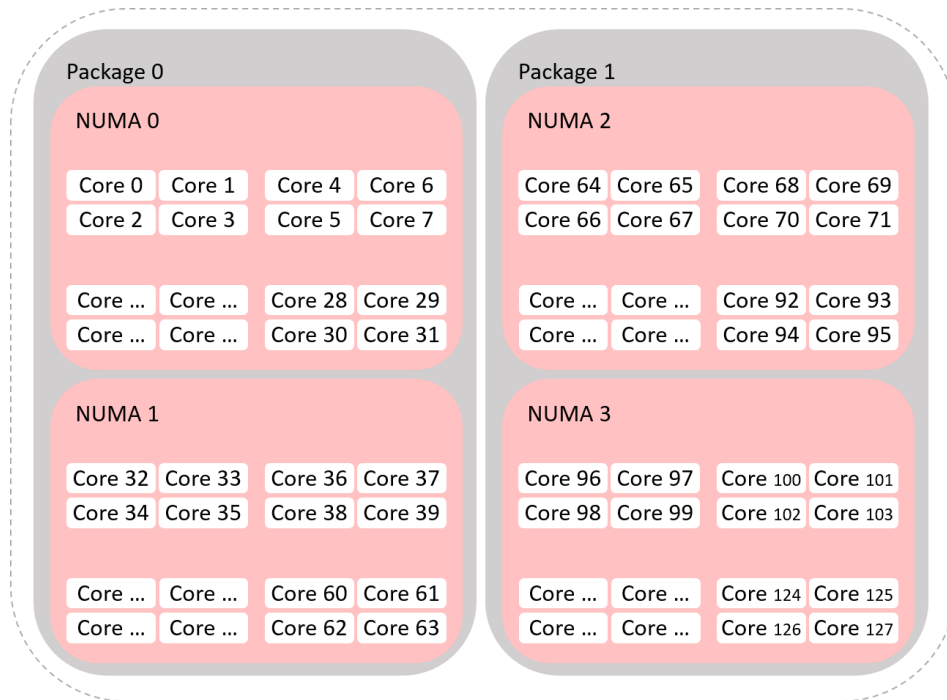


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

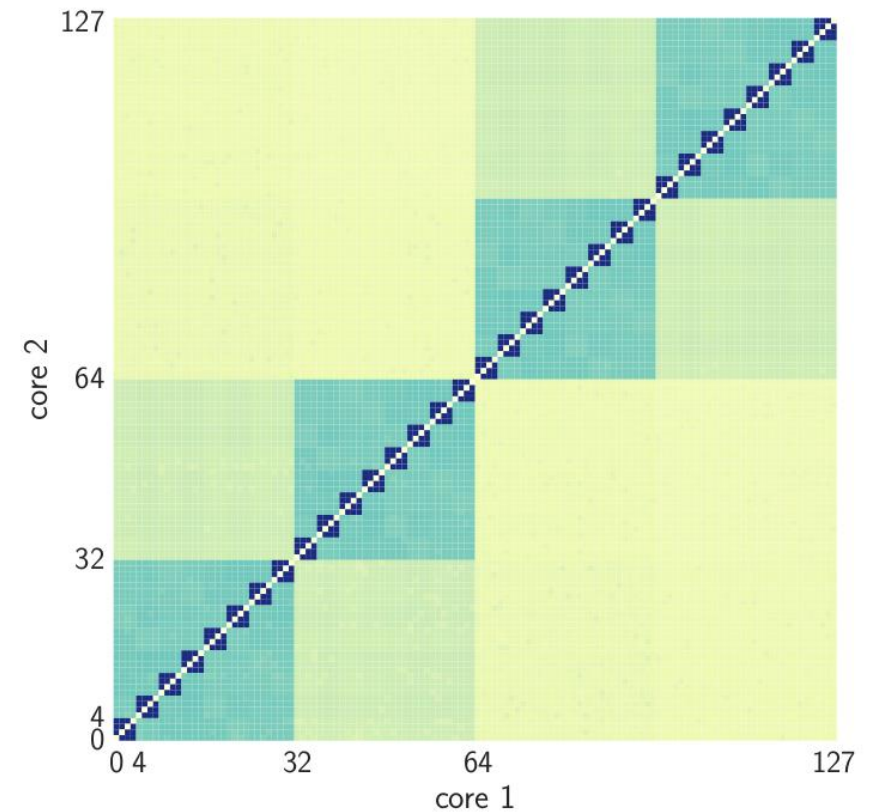
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

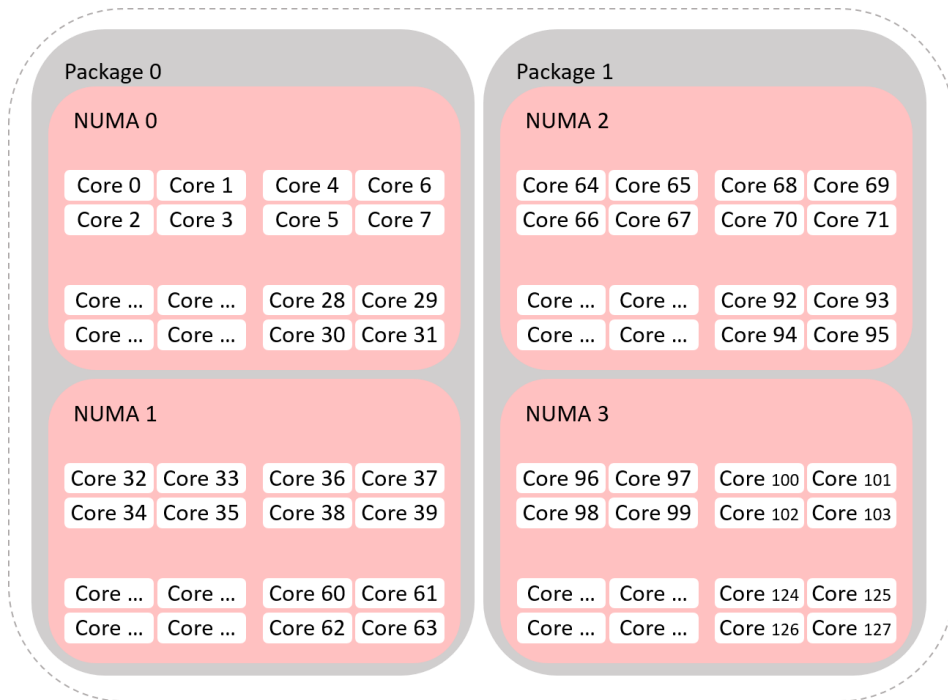


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

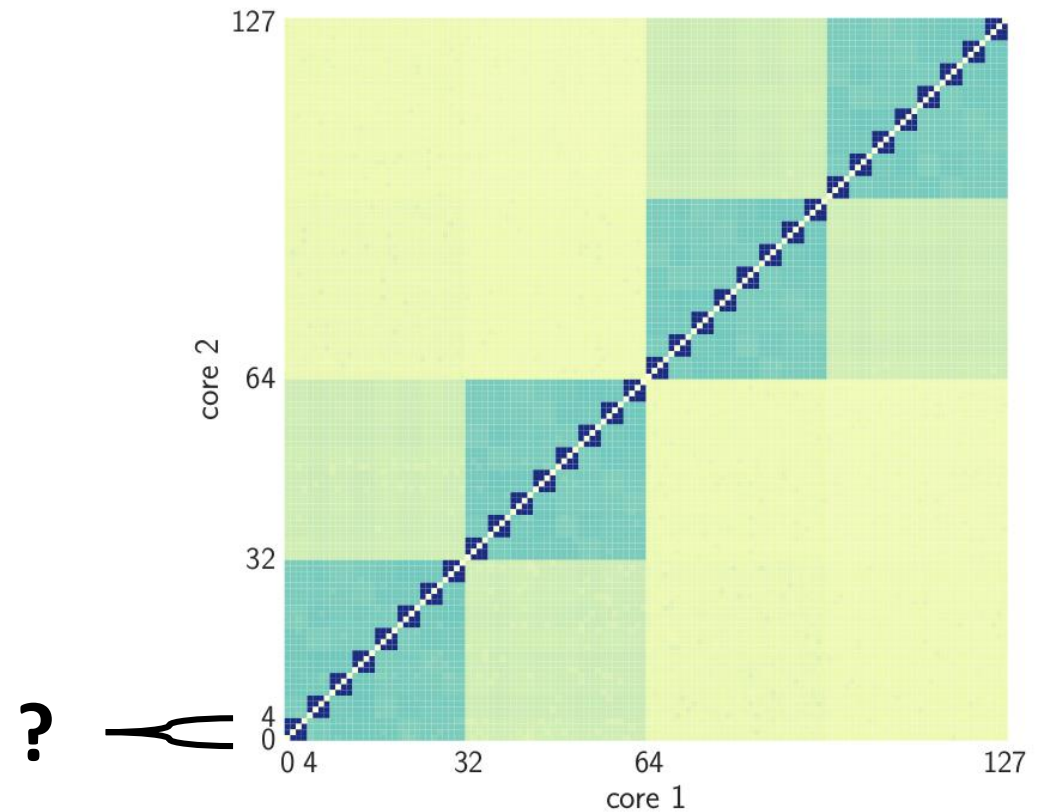
- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



## Experimental Discovery

- 2 threads alternately increment a shared counter
- Darker colors => Higher throughput

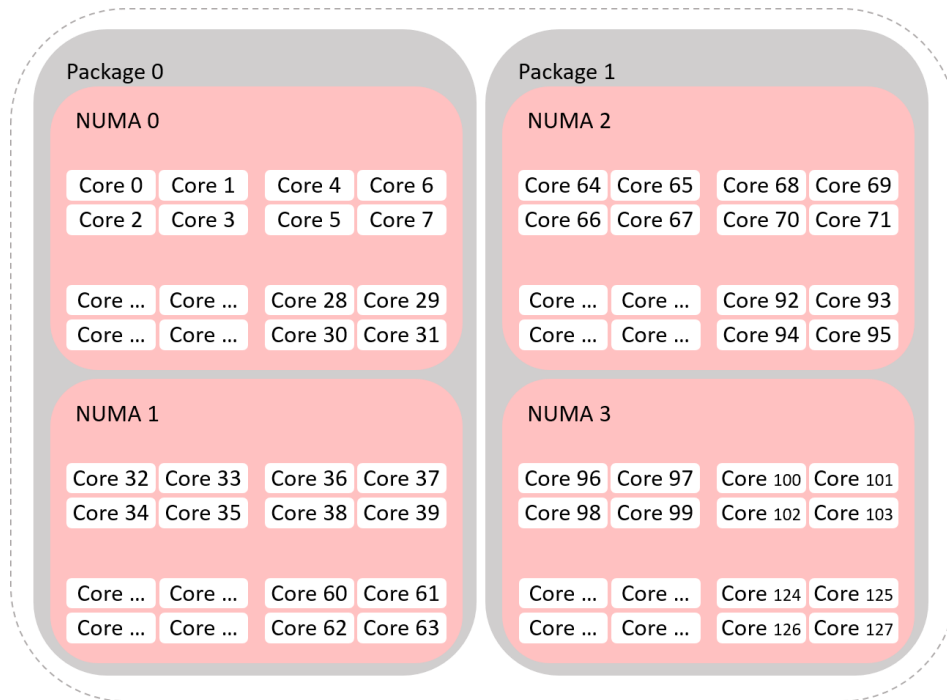


# Discovering the Memory Hierarchy

Operating systems know the hierarchy

- Visible with Linux's `lscpu` and `lstopo`

Taishan 200 server – Armv8



- Not shown by `lscpu`/`lstopo`
- Information is shown in processors' datasheet
  - Not efficient
  - Remained unstudied and unused
    - HMCS<4> includes hidden cache group level

# Agenda

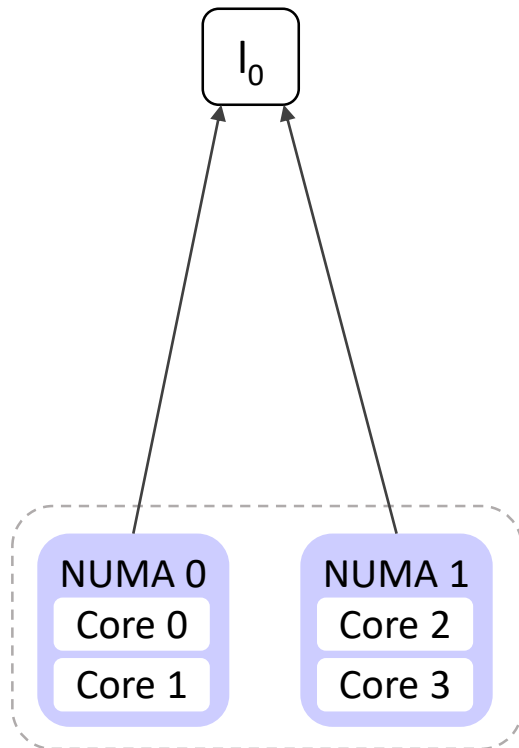
- How to figure out the hierarchy we need to use?
- **How does our CLoF Lock Generator works?**
- How do we know it is correct?
- How do we pick the best lock for the target platform?

# CLoF Lock Generator

Two NUMA node example

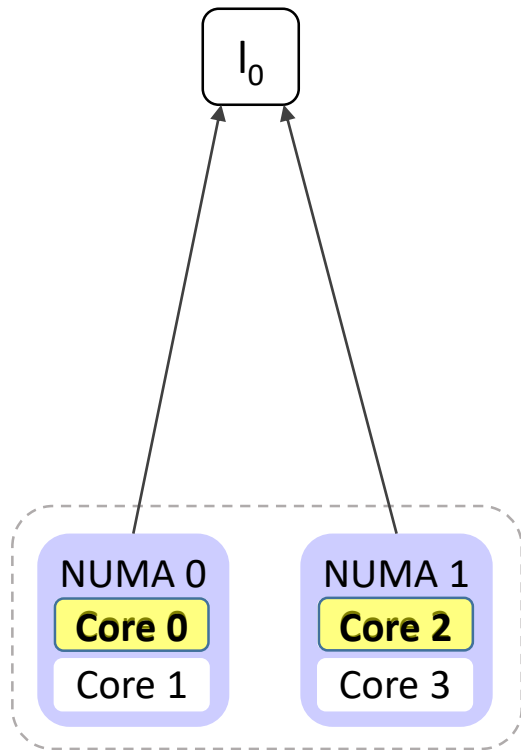
# CLoF Lock Generator

Two NUMA node example



# CLoF Lock Generator

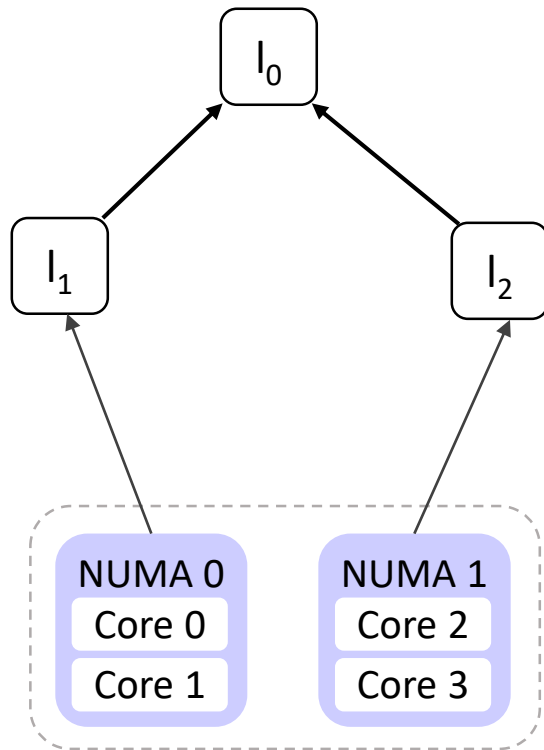
Two NUMA node example





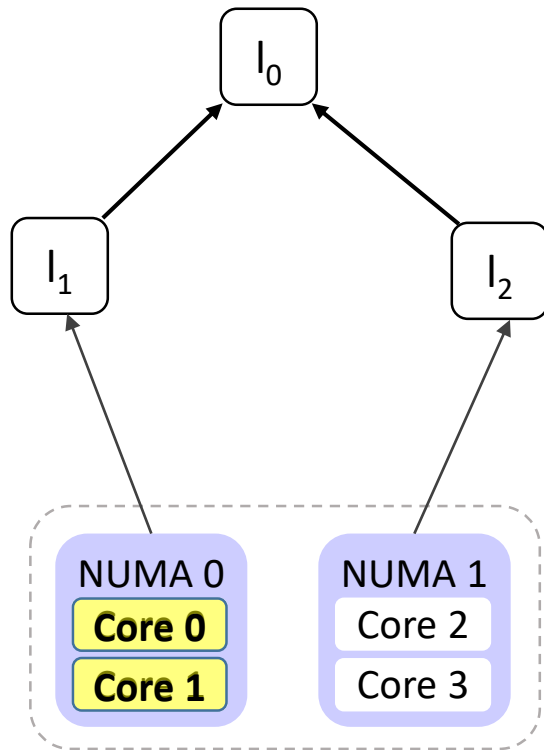
# CLoF Lock Generator

Two NUMA node example



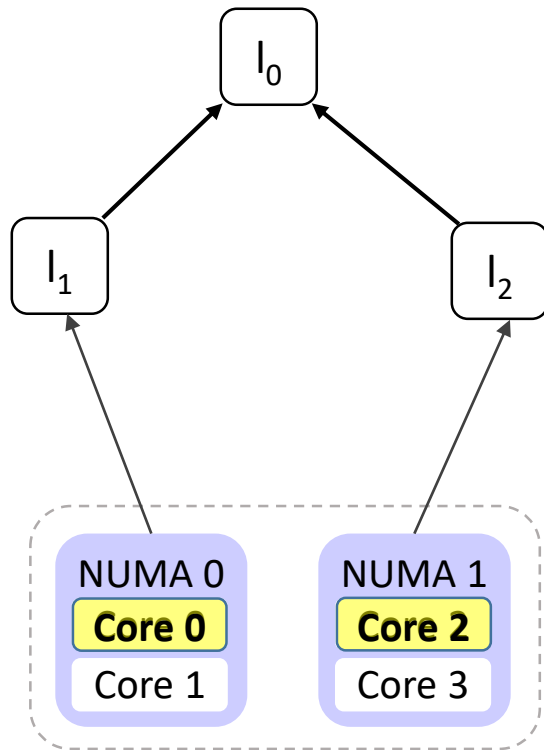
# CLoF Lock Generator

Two NUMA node example



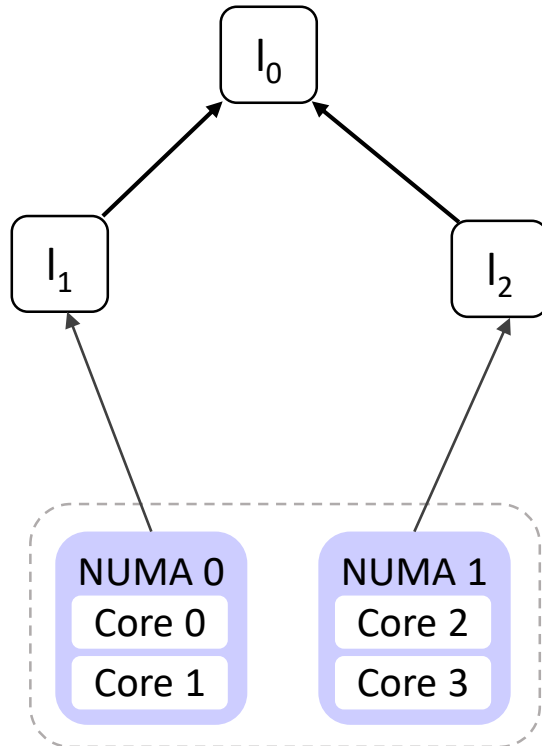
# CLoF Lock Generator

Two NUMA node example



# CLoF Lock Generator

Two NUMA node example



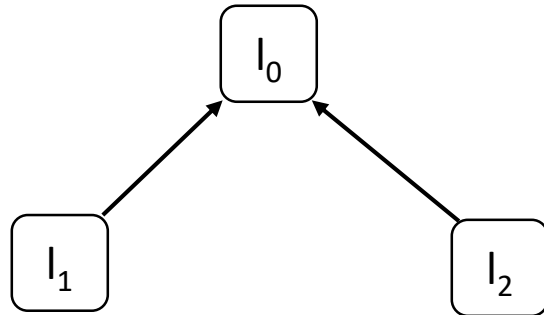
Simplified Code

```
CLoF(l, L)::acquire =  
    acquire l;  
    if(¬already has L)  
        acquire L;
```

```
CLoF(l, L)::release =  
    if(someone is waiting and  
       others won't starve)  
        release l;  
    else  
        release L;  
        release l;
```

# CLoF Lock Generator

## Compile-Time Syntactic Recursion



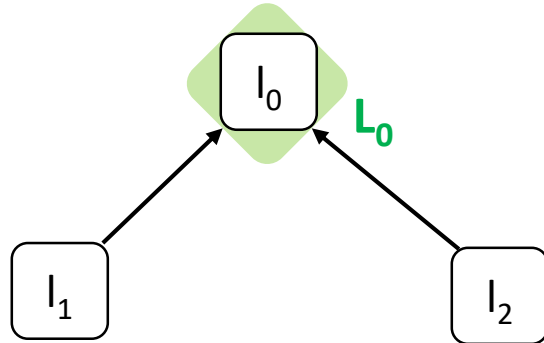
## Simplified Code

```
CLoF(l, L)::acquire =  
    acquire l;  
    if(¬already has L)  
        acquire L;
```

```
CLoF(l, L)::release =  
    if(someone is waiting and  
        others won't starve)  
        release l;  
    else  
        release L;  
        release l;
```

# CLoF Lock Generator

## Compile-Time Syntactic Recursion



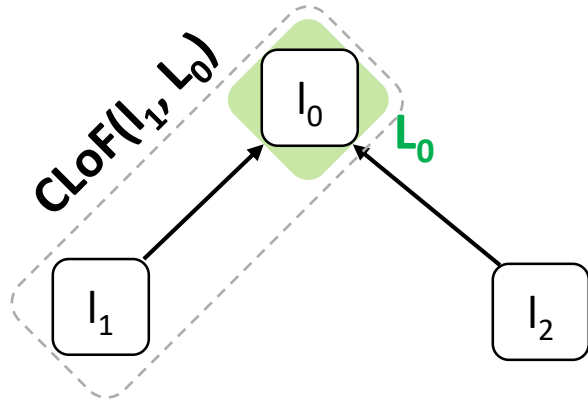
## Simplified Code

```
CLoF(l, L)::acquire =  
    acquire l;  
    if(¬already has L)  
        acquire L;
```

```
CLoF(l, L)::release =  
    if(someone is waiting and  
        others won't starve)  
        release l;  
    else  
        release L;  
        release l;
```

# CLoF Lock Generator

## Compile-Time Syntactic Recursion



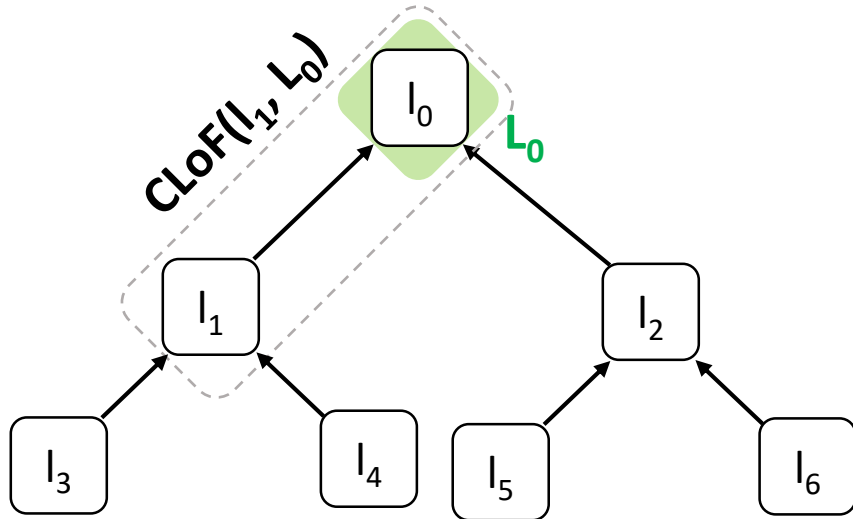
## Simplified Code

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

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

# CLoF Lock Generator

## Compile-Time Syntactic Recursion



## Simplified Code

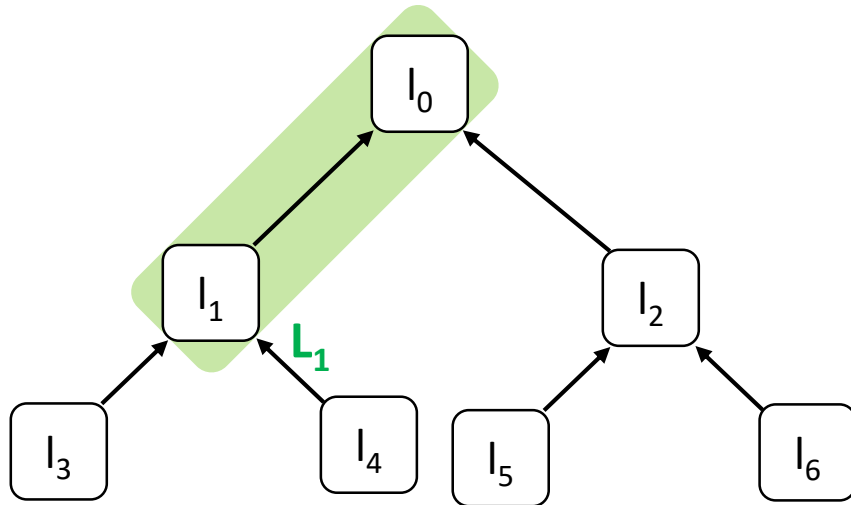
```
CLOF(l, L)::acquire =  
    acquire l;  
    if(¬already has L)  
        acquire L;
```

```
CLOF(l, L)::release =  
    if(someone is waiting and  
       others won't starve)  
        release l;  
    else  
        release L;  
        release l;
```



# CLoF Lock Generator

## Compile-Time Syntactic Recursion



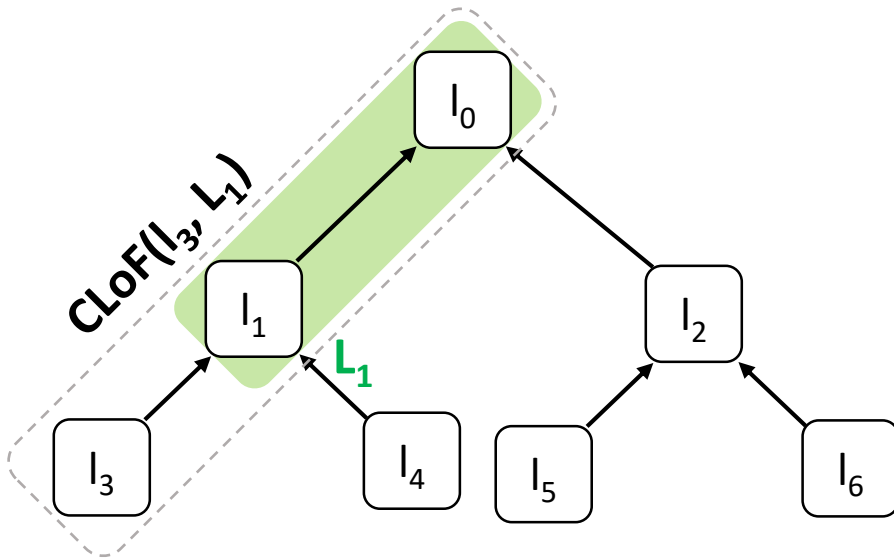
## Simplified Code

```
CLoF(l, L)::acquire =  
    acquire l;  
    if(¬already has L)  
        acquire L;
```

```
CLoF(l, L)::release =  
    if(someone is waiting and  
        others won't starve)  
        release l;  
    else  
        release L;  
        release l;
```

# CLoF Lock Generator

## Compile-Time Syntactic Recursion



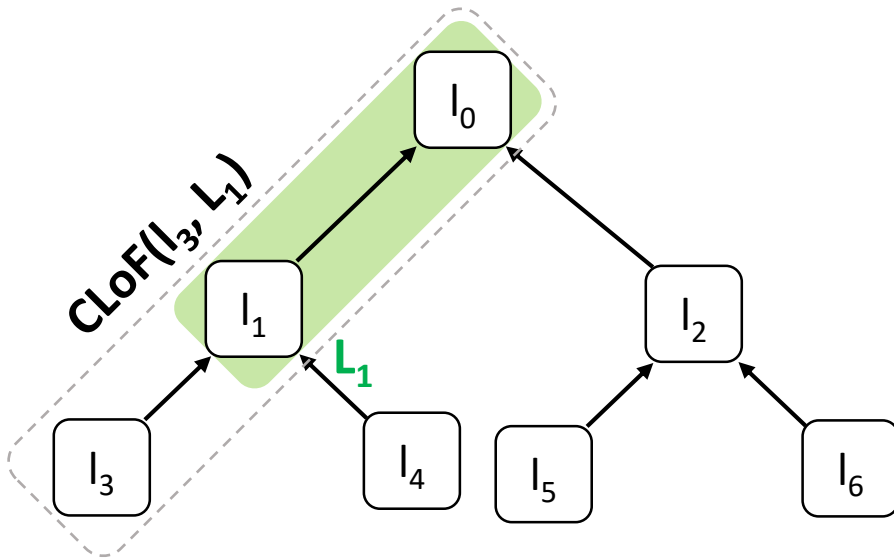
## Simplified Code

```
CloF(l, L)::acquire =  
    acquire l;  
    if(¬already has L)  
        acquire L;
```

```
CloF(l, L)::release =  
    if(someone is waiting and  
       others won't starve)  
        release l;  
    else  
        release L;  
        release l;
```

# CLoF Lock Generator

## Compile-Time Syntactic Recursion



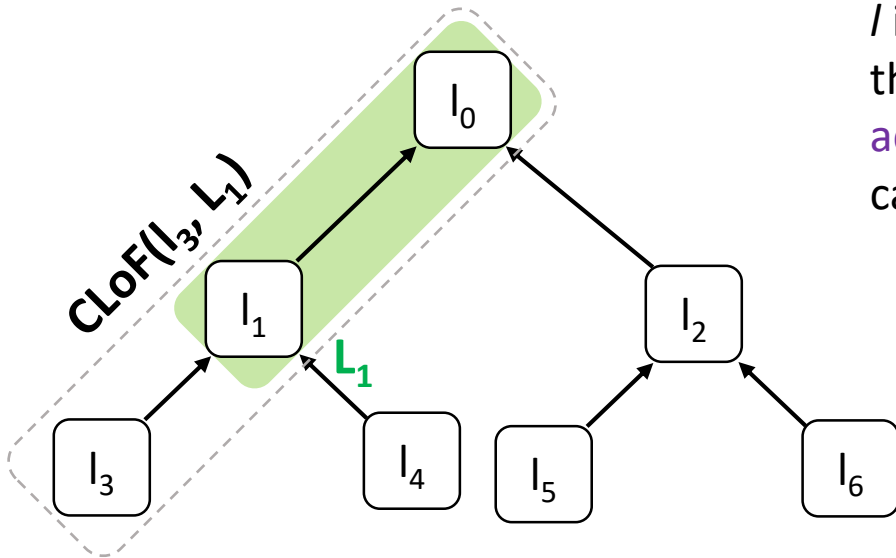
## Lock Abstraction

```
CLOF(l, L)::acquire =  
    acquire l;  
    if(¬already has L)  
        acquire L;
```

```
CLOF(l, L)::release =  
    if(someone is waiting and  
        others won't starve)  
        release l;  
    else  
        release L;  
        release l;
```

# CLoF Lock Generator

## Compile-Time Syntactic Recursion



$l$  is only accessed through `acquire/release` calls

## Lock Abstraction

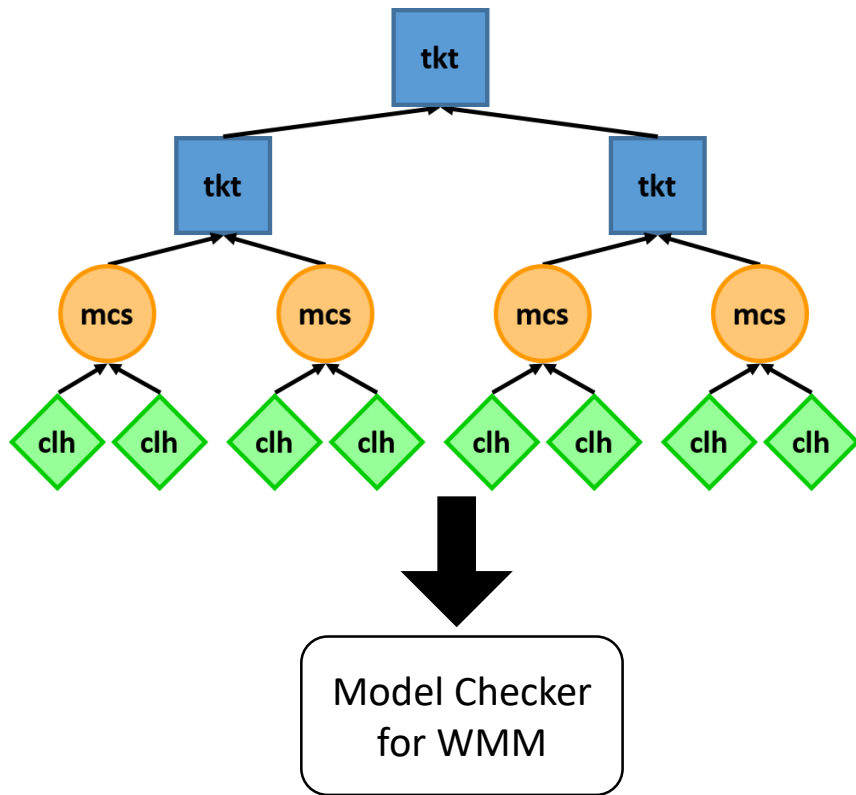
```
CLOF( $l, L$ )::acquire =  
→ acquire  $l$ ;  
if( $\neg$ already has  $L$ )  
→ acquire  $L$ ;
```

```
CLOF( $l, L$ )::release =  
→ if(someone is waiting and  
→ others won't starve)  
→ release  $l$ ;  
else  
→ release  $L$ ;  
→ release  $l$ ;
```

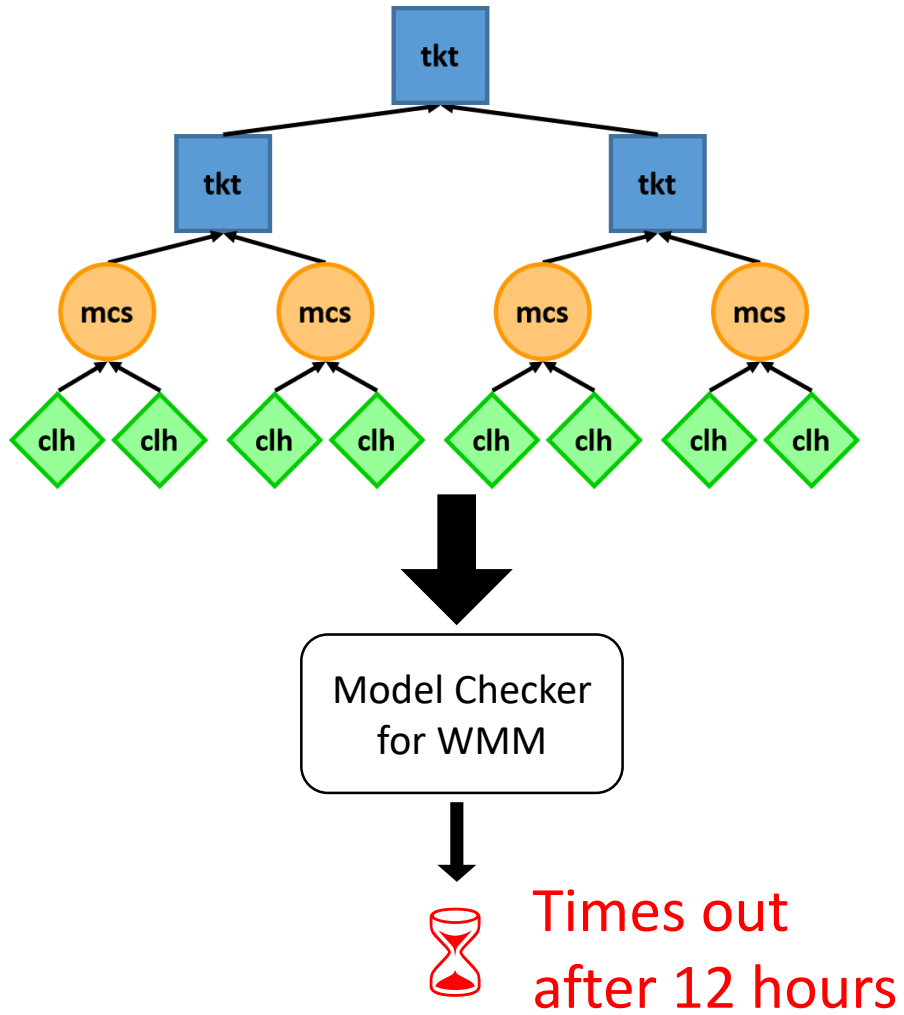
# Agenda

- How to figure out the hierarchy we need to use?
- How does our CLoF Lock Generator works?
- **How do we know it is correct?**
- How do we pick the best lock for the target platform?

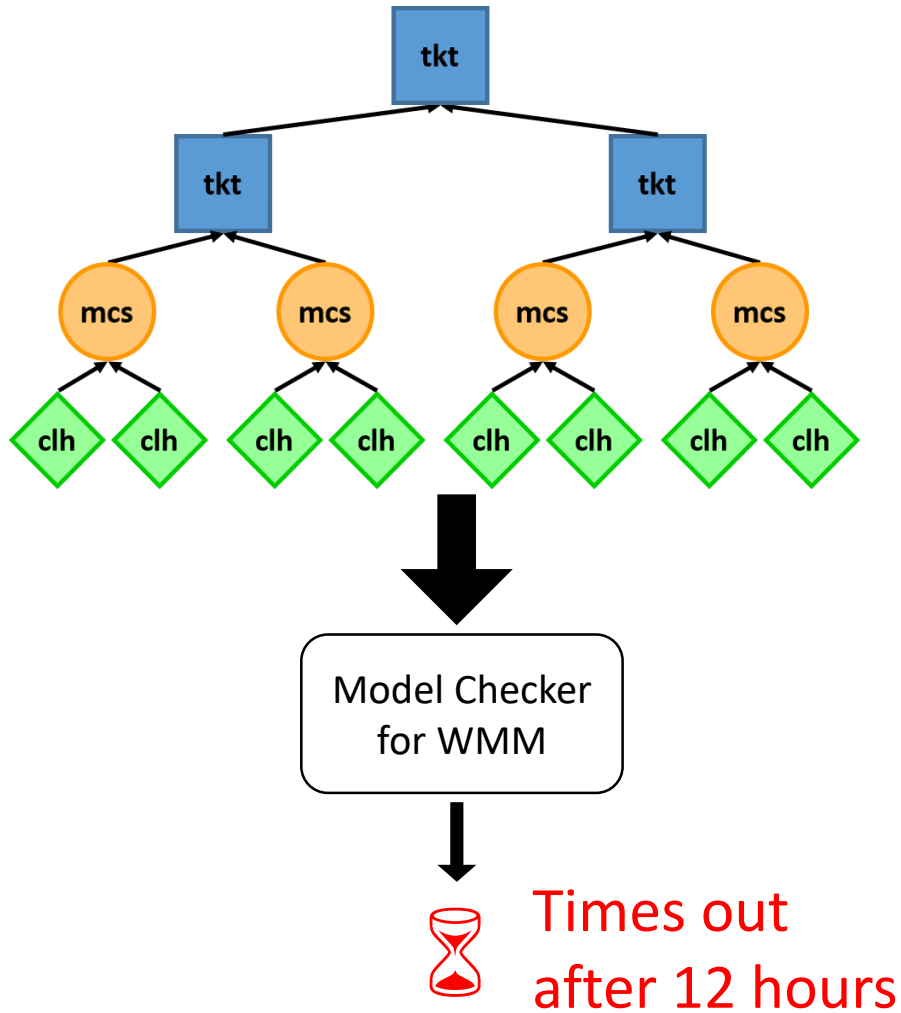
# CLoF Correctness



# CLoF Correctness



# CLoF Correctness



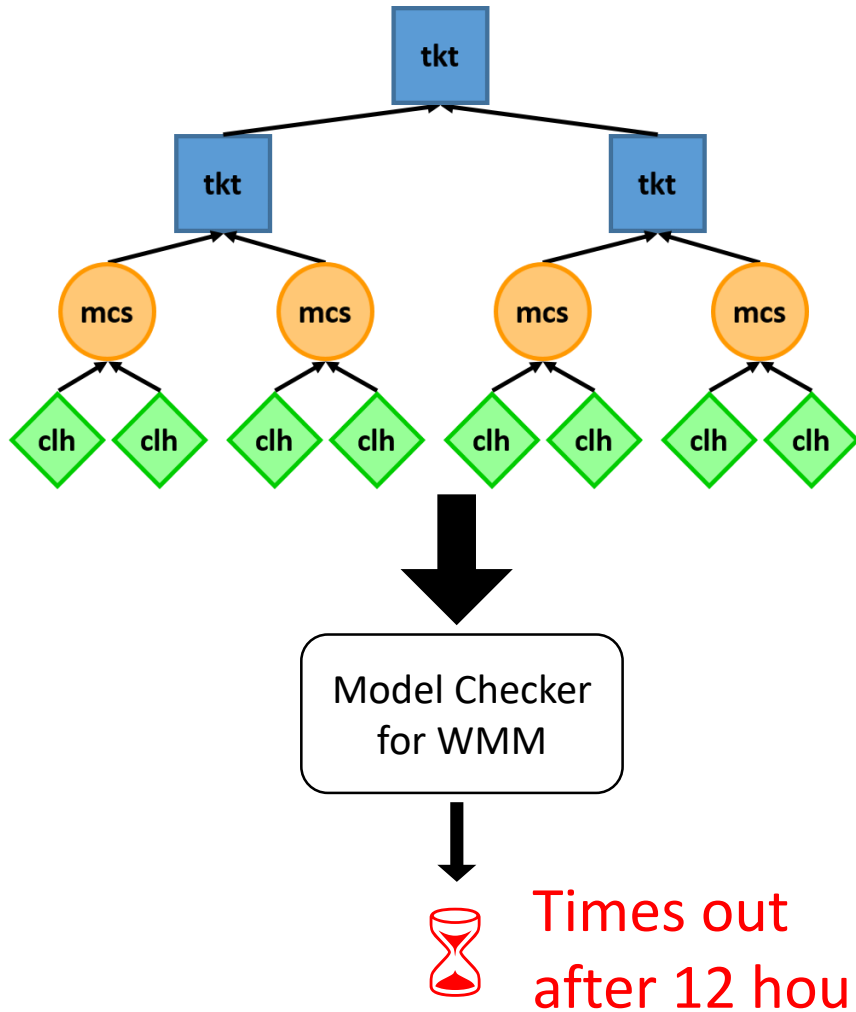
## Inductive Approach

Base Step

Inductive Step

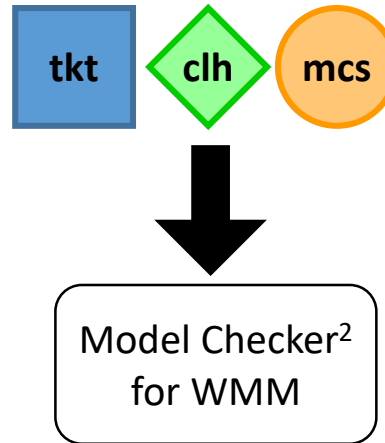


# CLoF Correctness



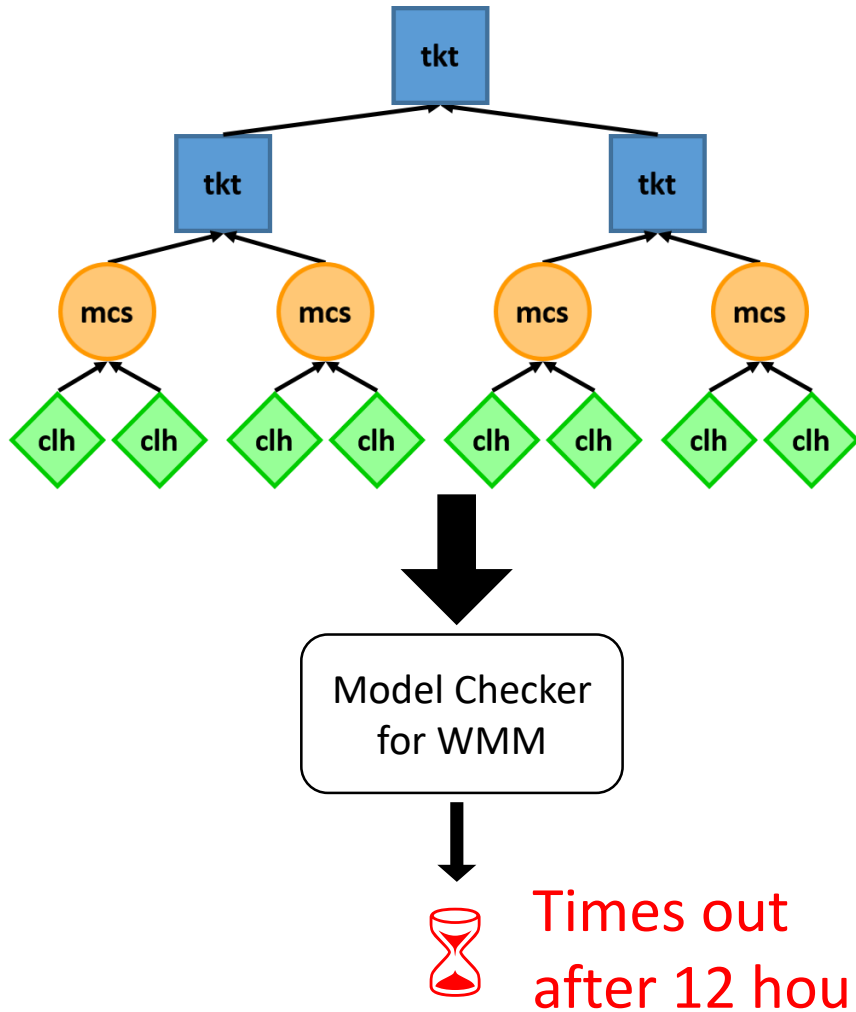
## Inductive Approach

### Base Step



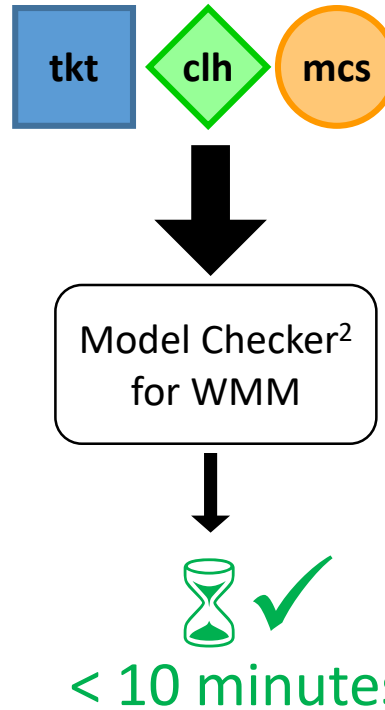
### Inductive Step

# CLoF Correctness



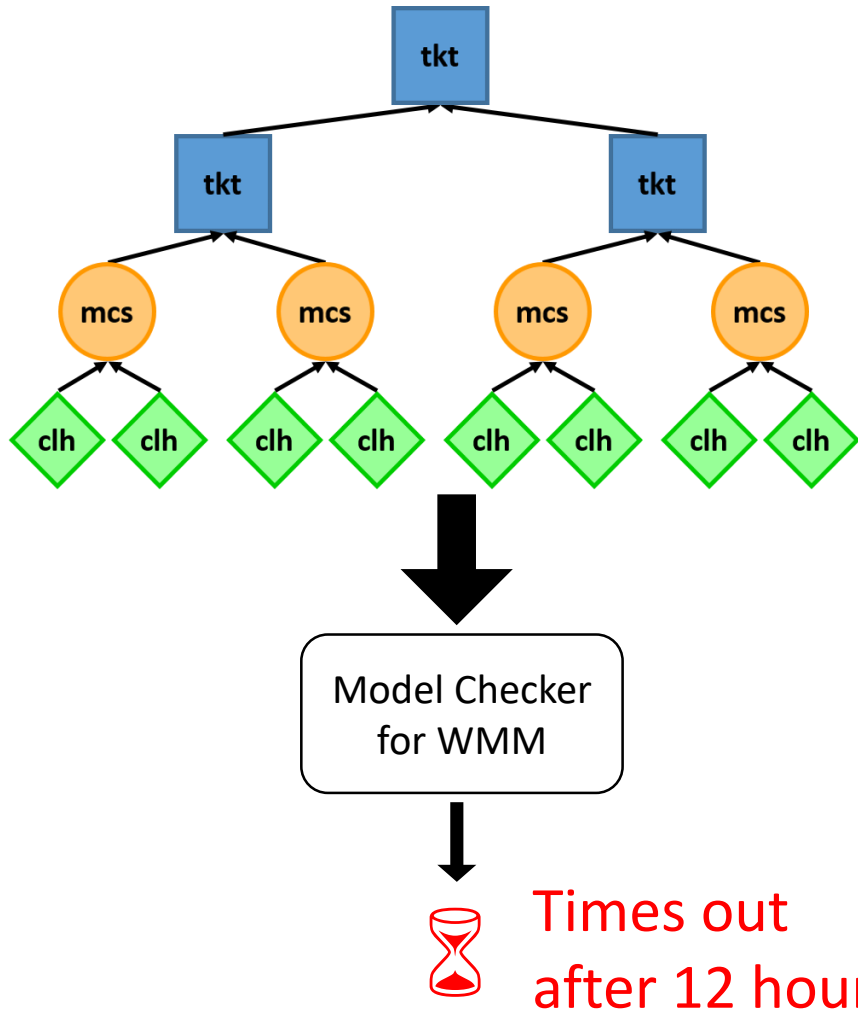
## Inductive Approach

### Base Step



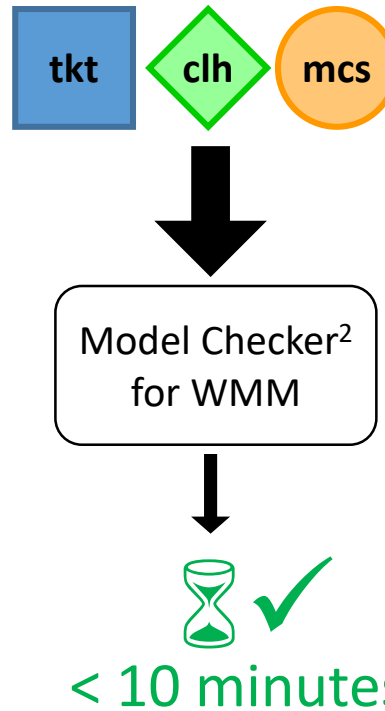
### Inductive Step

# CLoF Correctness

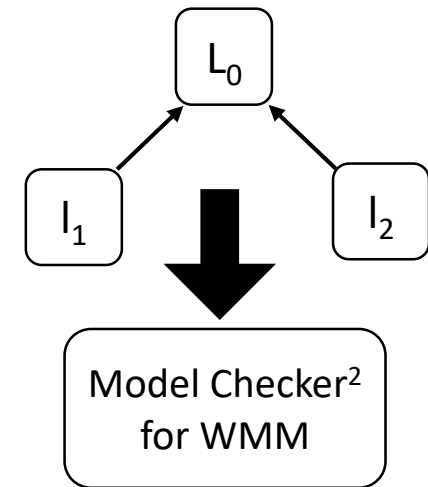


## Inductive Approach

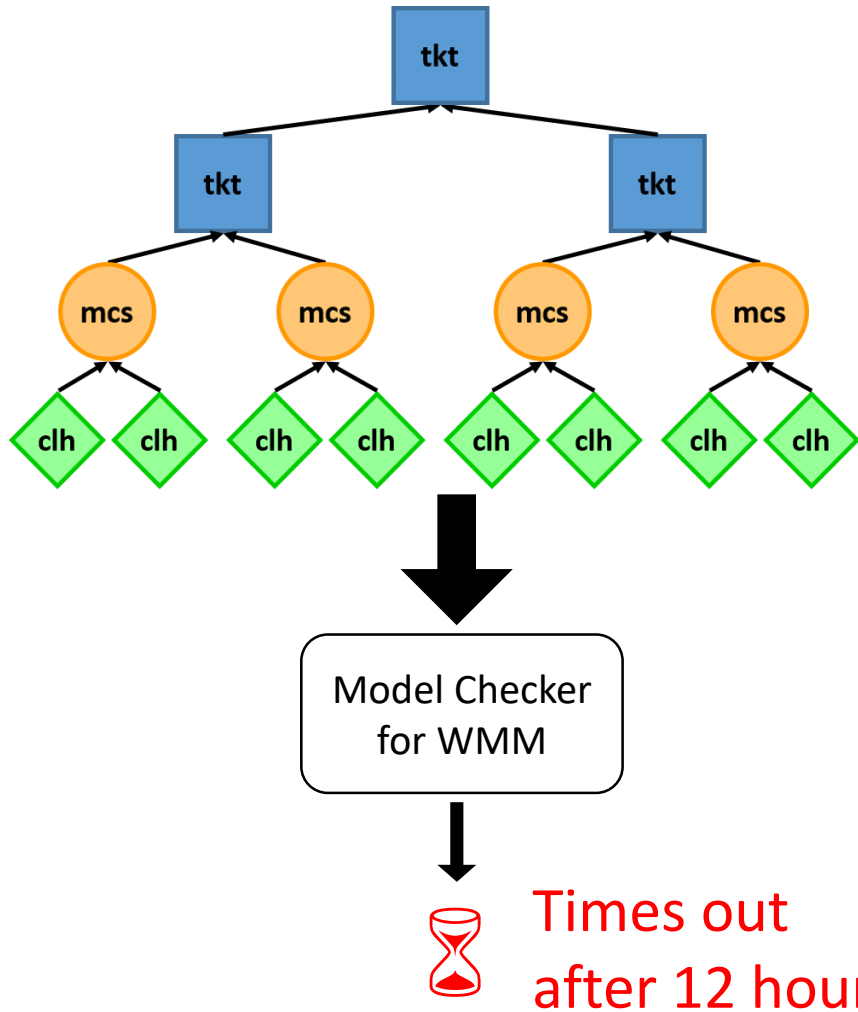
### Base Step



### Inductive Step

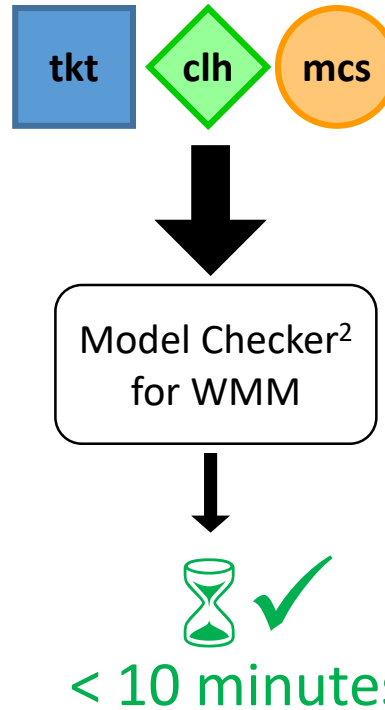


# CLoF Correctness

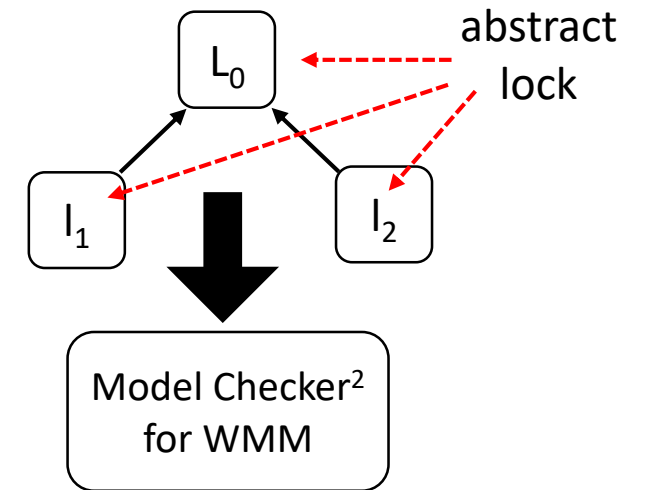


## Inductive Approach

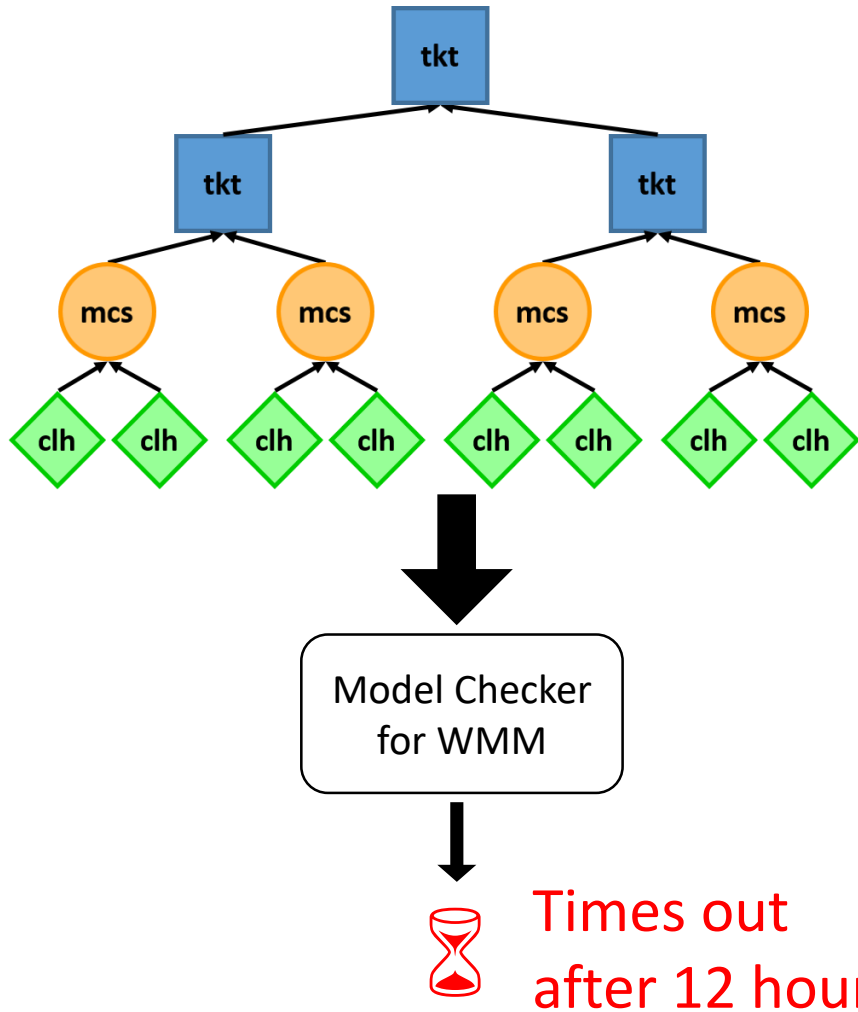
### Base Step



### Inductive Step

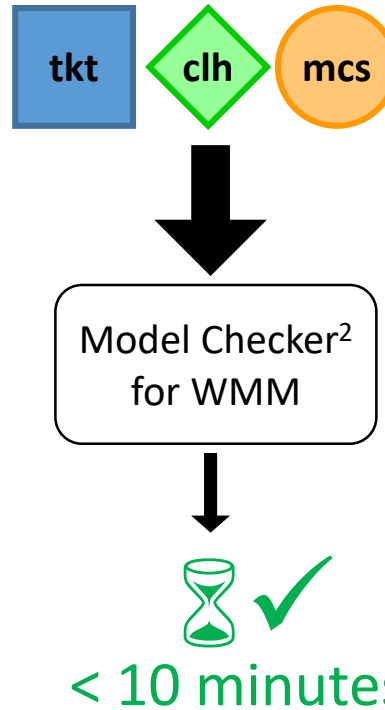


# CLoF Correctness

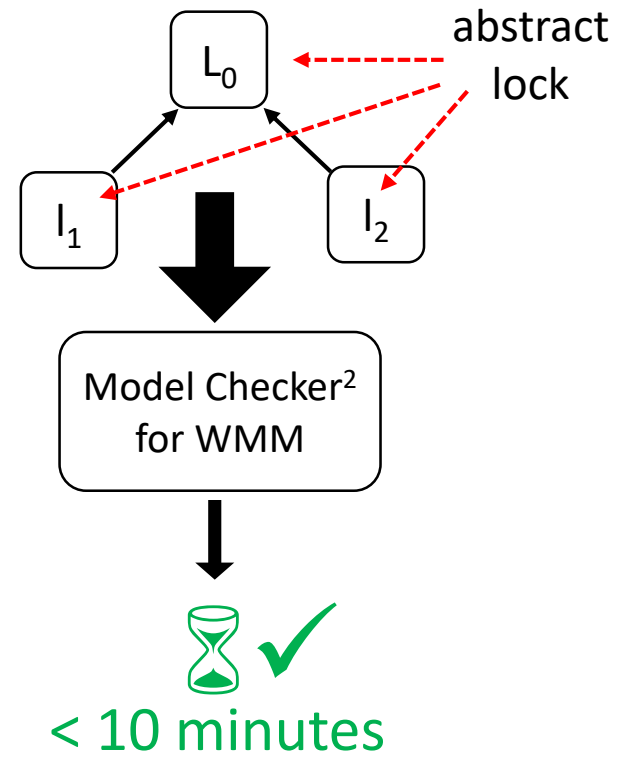


## Inductive Approach

### Base Step



### Inductive Step



# Agenda

- How to figure out the hierarchy we need to use?
- How does our CLoF Lock Generator works?
- How do we know it is correct?
- How do we pick the best lock for the target platform?

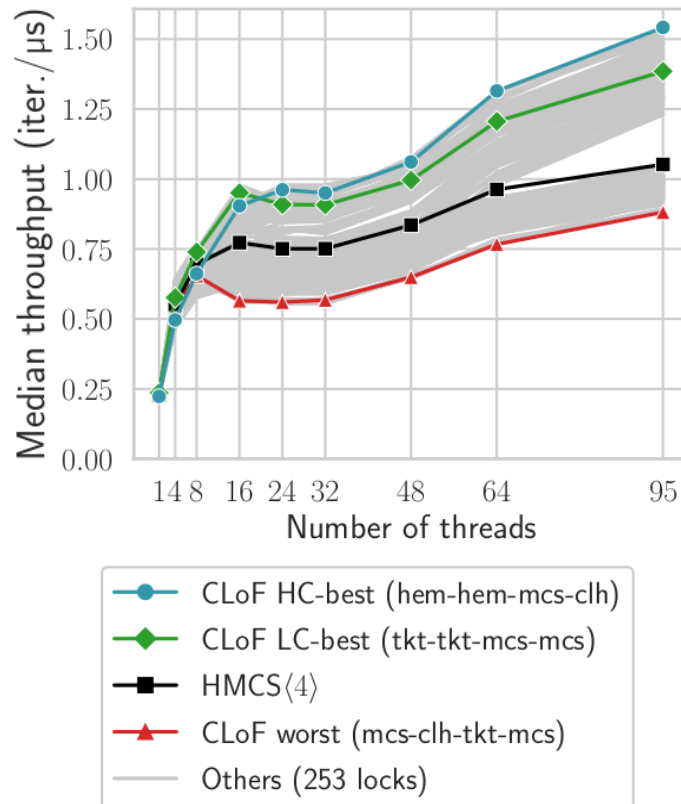
# Run Scripted Benchmark

- Run all lock combinations composed by the CLoF lock generator
- For  $B = 4$  locks and  $L = 4$  levels, we have  $B^L = 4^4 = 256$  combinations
  - Up to 1 hour in a platform with 128 cores

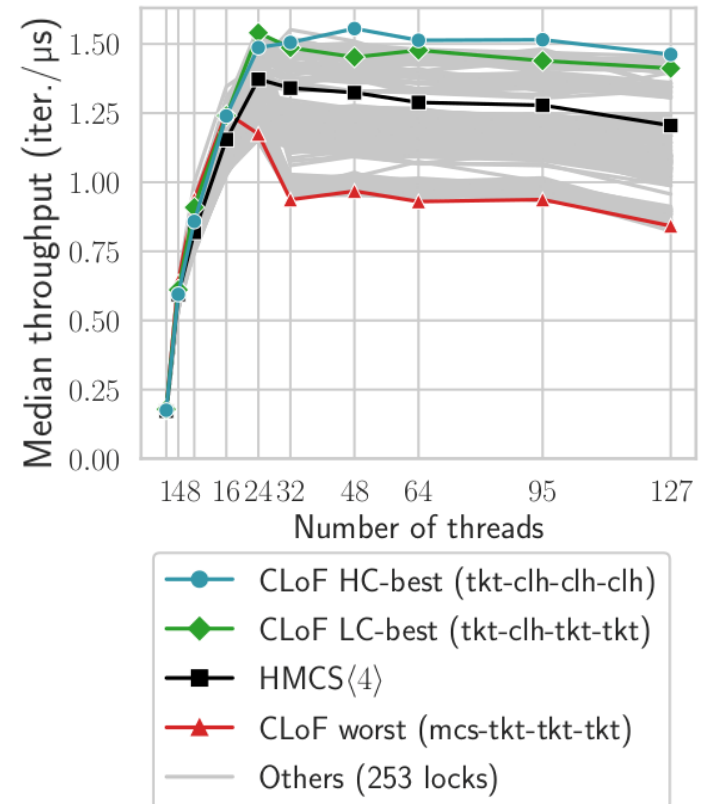
# Run Scripted Benchmark

- Run all lock combinations composed by the CLoF lock generator
- For  $B = 4$  locks and  $L = 4$  levels, we have  $B^L = 4^4 = 256$  combinations
  - Up to 1 hour in a platform with 128 cores

x86 server – levelDB readrandom benchmark



arm server – levelDB readrandom benchmark

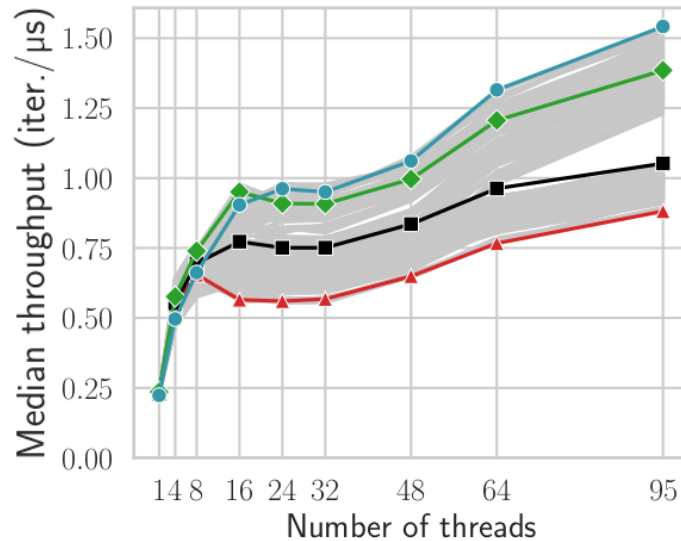




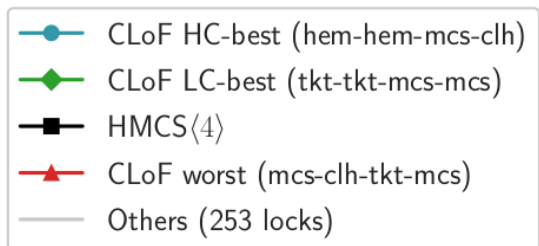
# Run Scripted Benchmark

- Run all lock combinations composed by the CLoF lock generator
- For  $B = 4$  locks and  $L = 4$  levels, we have  $B^L = 4^4 = 256$  combinations
  - Up to 1 hour in a platform with 128 cores

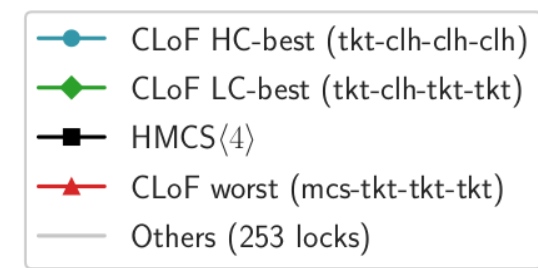
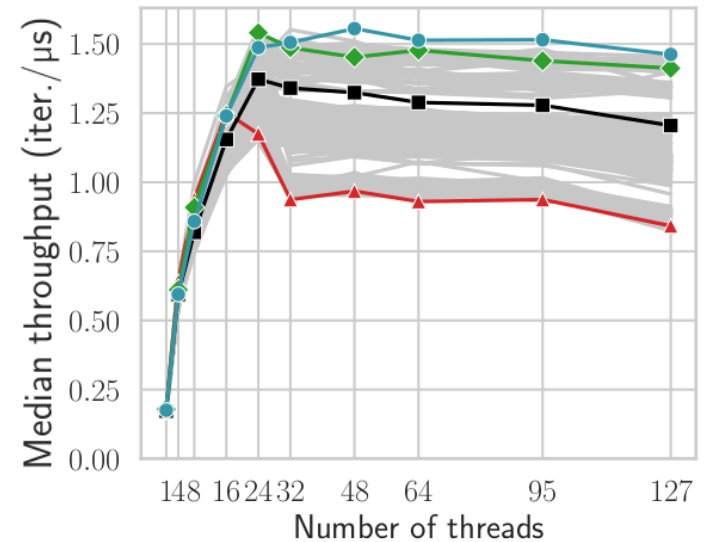
x86 server – levelDB readrandom benchmark



No combination is better for all contention levels



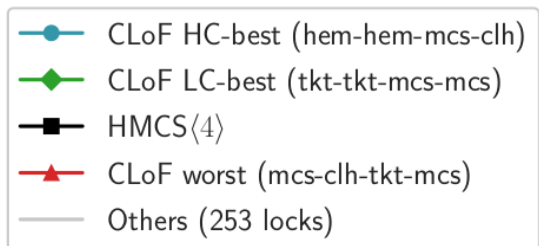
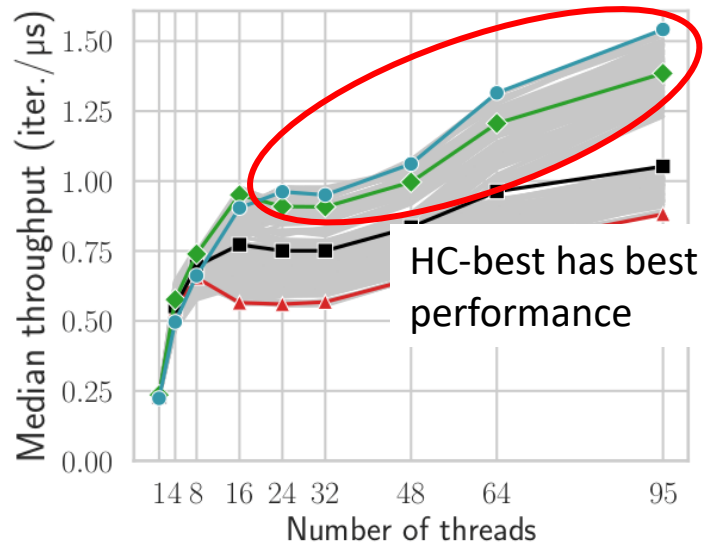
arm server – levelDB readrandom benchmark



# Run Scripted Benchmark

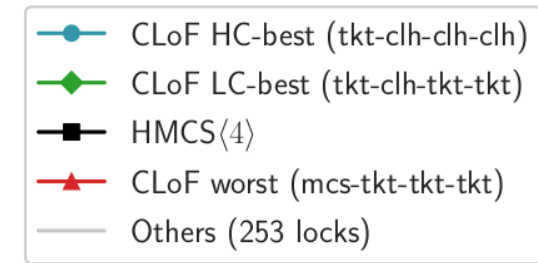
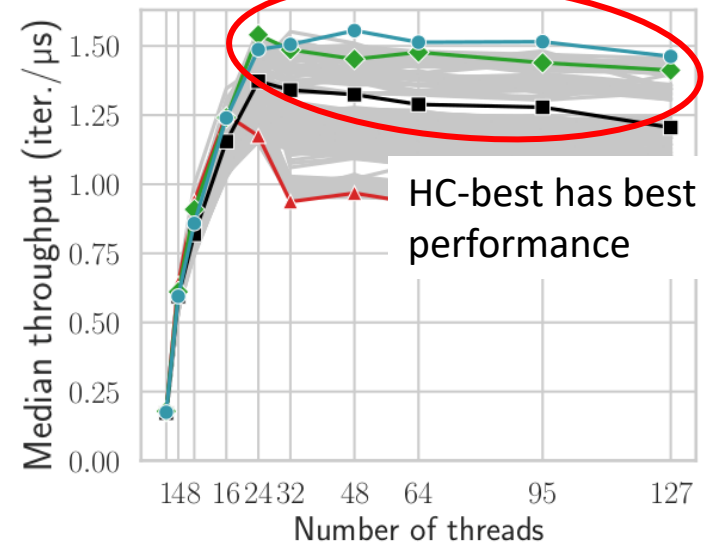
- Run all lock combinations composed by the CLoF lock generator
- For  $B = 4$  locks and  $L = 4$  levels, we have  $B^L = 4^4 = 256$  combinations
  - Up to 1 hour in a platform with 128 cores

x86 server – levelDB readrandom benchmark



No combination is better for all contention levels

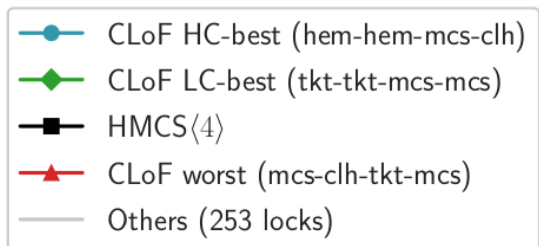
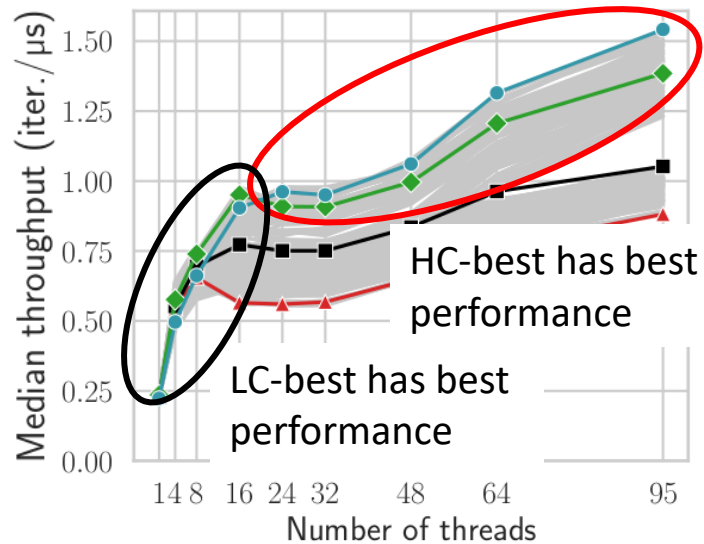
arm server – levelDB readrandom benchmark



# Run Scripted Benchmark

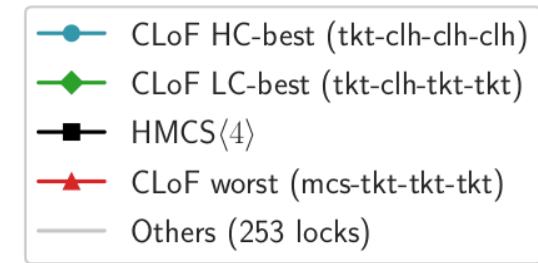
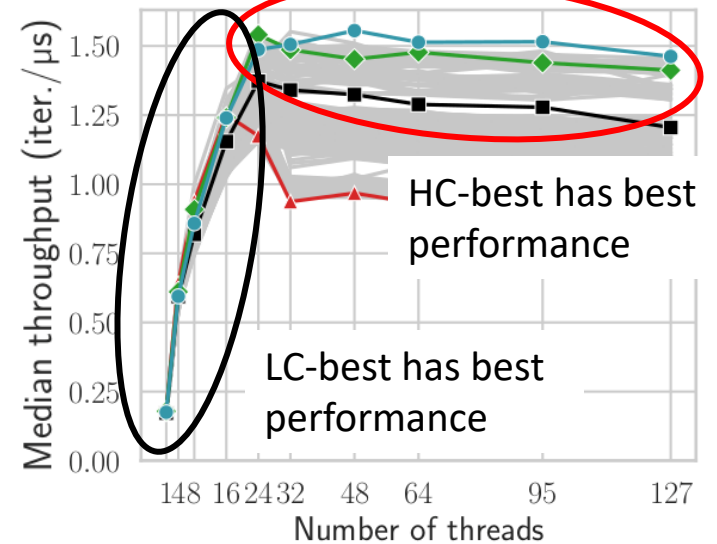
- Run all lock combinations composed by the CLoF lock generator
- For  $B = 4$  locks and  $L = 4$  levels, we have  $B^L = 4^4 = 256$  combinations
  - Up to 1 hour in a platform with 128 cores

x86 server – levelDB readrandom benchmark



No combination is better for all contention levels

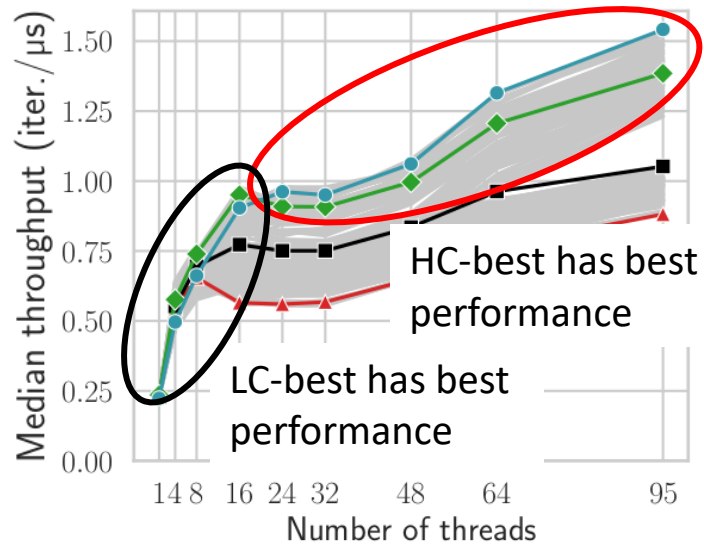
arm server – levelDB readrandom benchmark



# Run Scripted Benchmark

- Run all lock combinations composed by the CLoF lock generator
- For  $B = 4$  locks and  $L = 4$  levels, we have  $B^L = 4^4 = 256$  combinations
  - Up to 1 hour in a platform with 128 cores

x86 server – levelDB readrandom benchmark

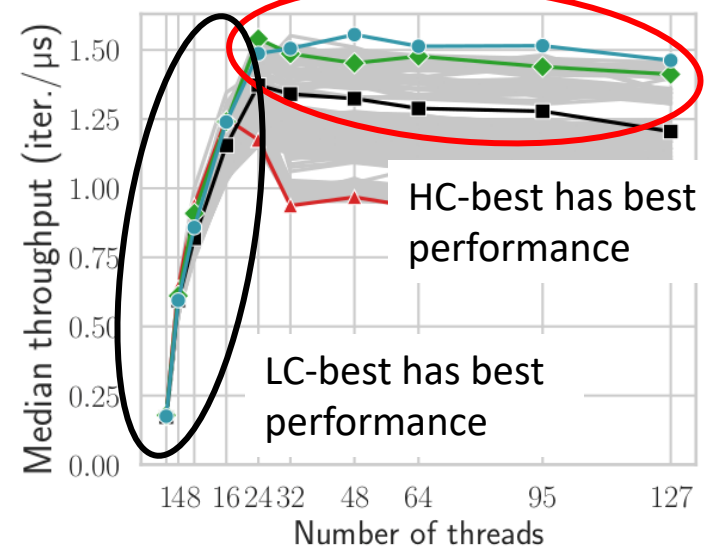


- CLoF HC-best (hem-hem-mcs-clh)
- ◆ CLoF LC-best (tk-tkt-mcs-mcs)
- HMCS(4)
- ▲ CLoF worst (mcs-clh-tkt-mcs)
- Others (253 locks)

No combination is better for all contention levels

Which one should be chosen?

arm server – levelDB readrandom benchmark

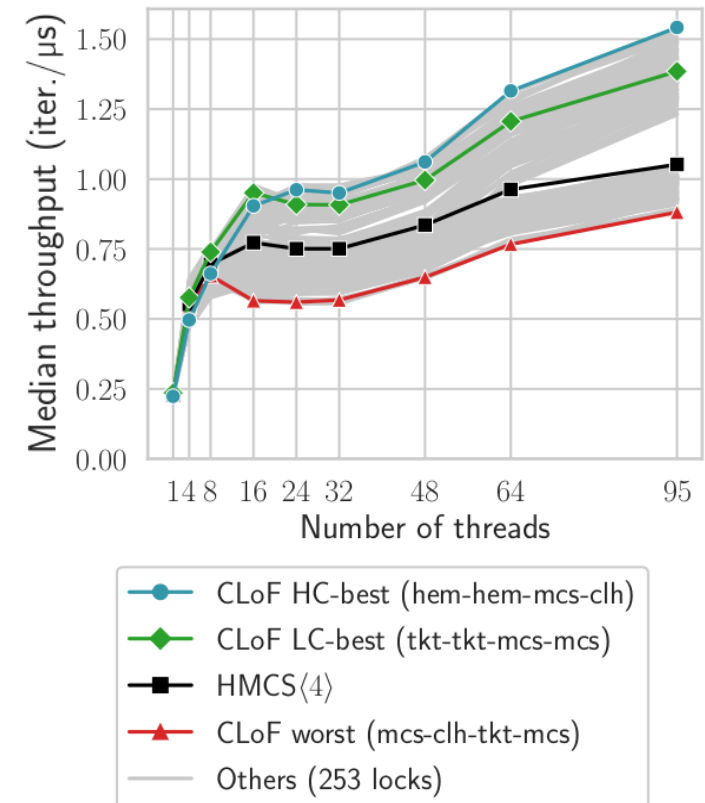


- CLoF HC-best (tk-clh-clh-clh)
- ◆ CLoF LC-best (tk-clh-tkt-tkt)
- HMCS(4)
- ▲ CLoF worst (mcs-tkt-tkt-tkt)
- Others (253 locks)

# Tuning Point: Choose Selection Policy

- CLoF establishes 2 selection policies for the best lock:
  - HC-best prioritizes performance at high contention
  - LC-best prioritizes performance at low contention

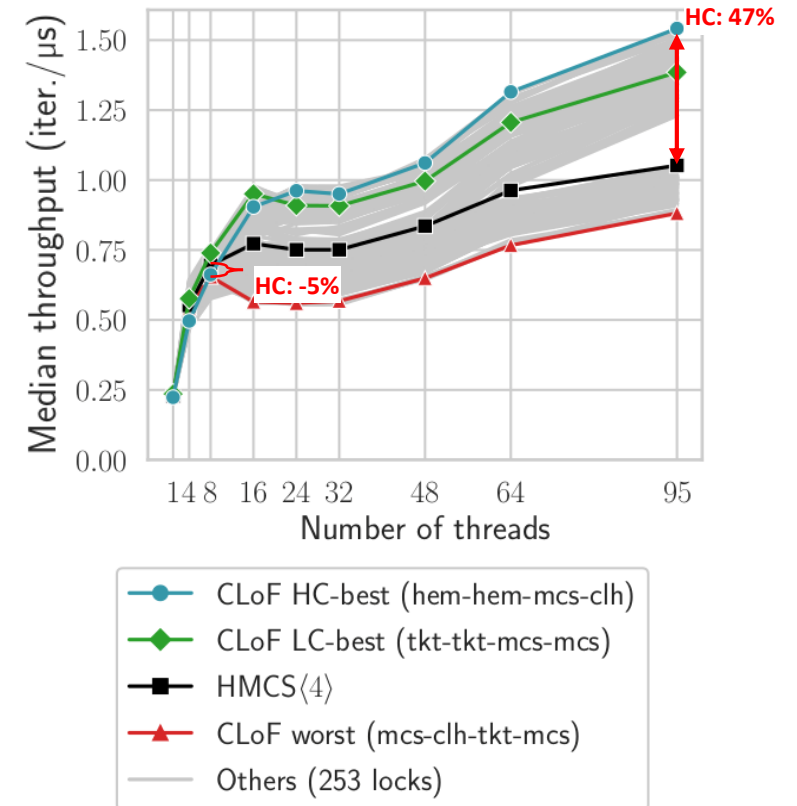
x86 server – levelDB readrandom benchmark



# Tuning Point: Choose Selection Policy

- CLoF establishes 2 selection policies for the best lock:
  - HC-best prioritizes performance at high contention
  - LC-best prioritizes performance at low contention

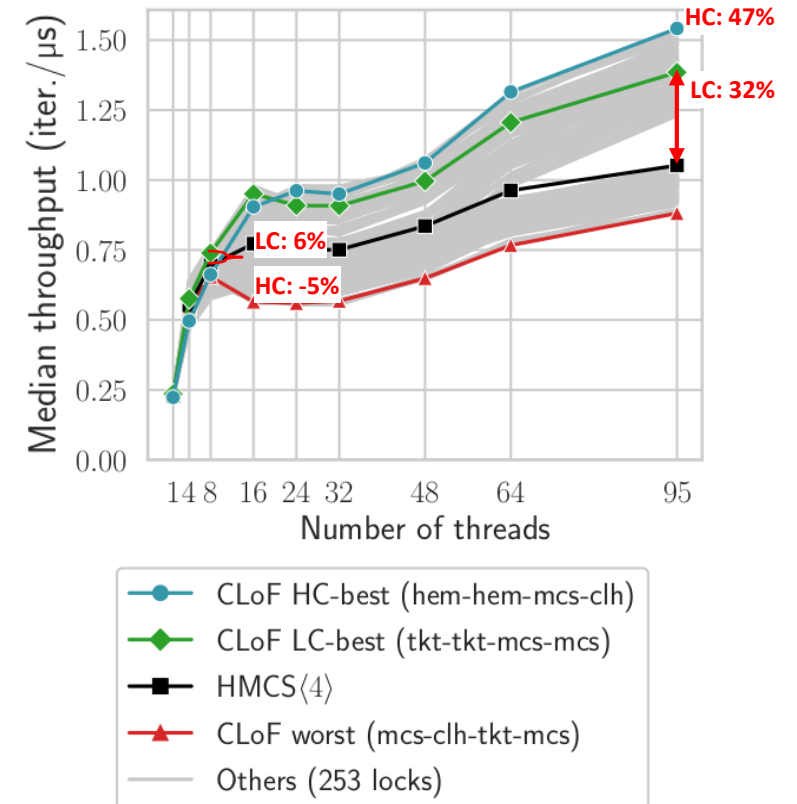
x86 server – levelDB readrandom benchmark



# Tuning Point: Choose Selection Policy

- CLoF establishes 2 selection policies for the best lock:
  - HC-best prioritizes performance at high contention
  - LC-best prioritizes performance at low contention

x86 server – levelDB readrandom benchmark



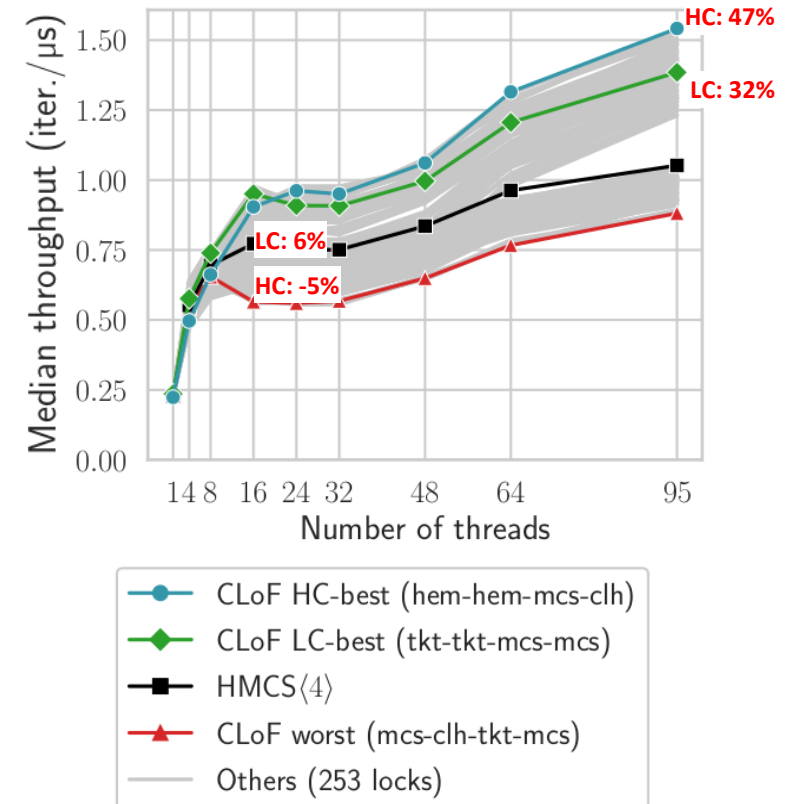
# Tuning Point: Choose Selection Policy

- CLoF establishes 2 selection policies for the best lock:
  - HC-best prioritizes performance at high contention
  - LC-best prioritizes performance at low contention

User can tune which selection policy is desired

- Without the need to re-run the benchmark

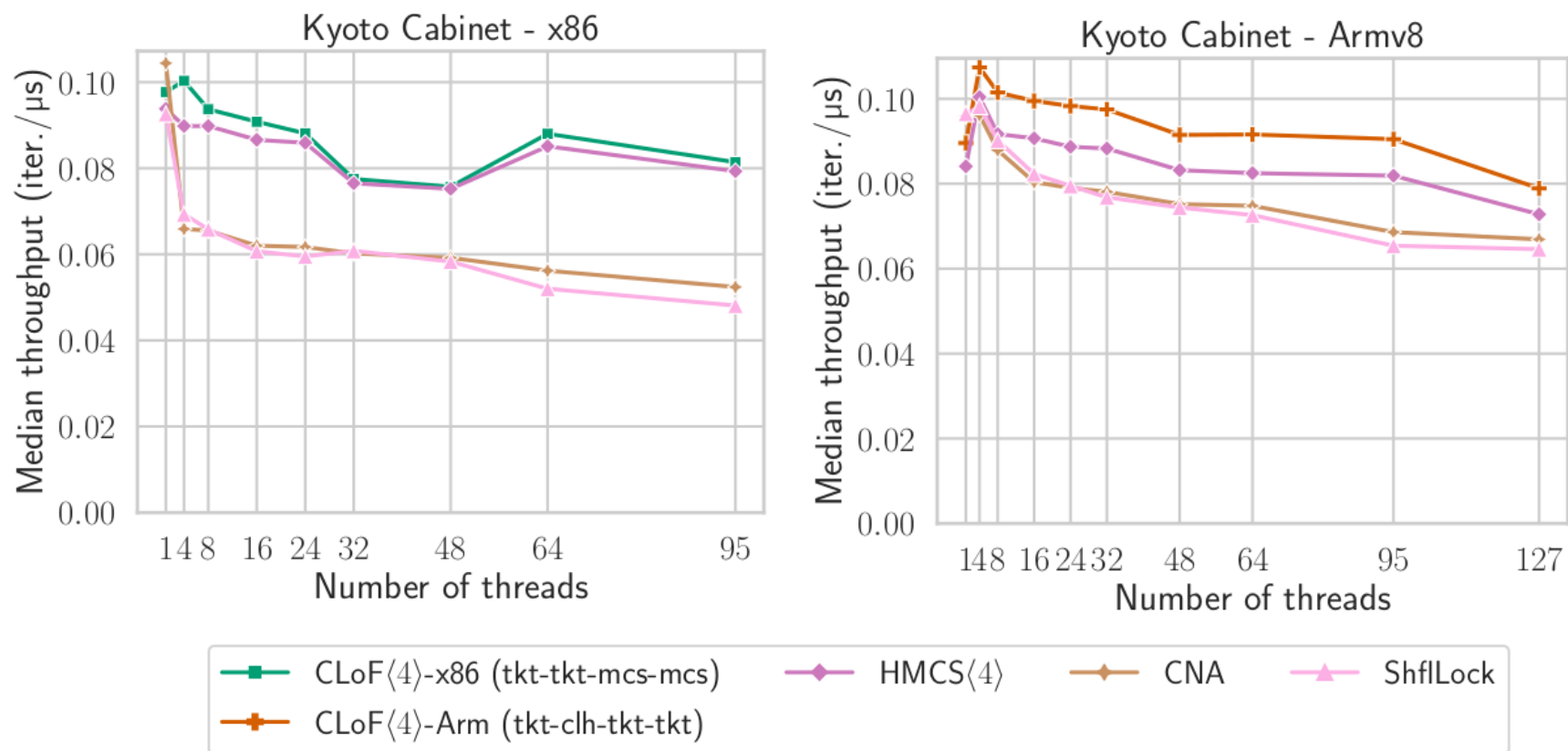
x86 server – levelDB readrandom benchmark





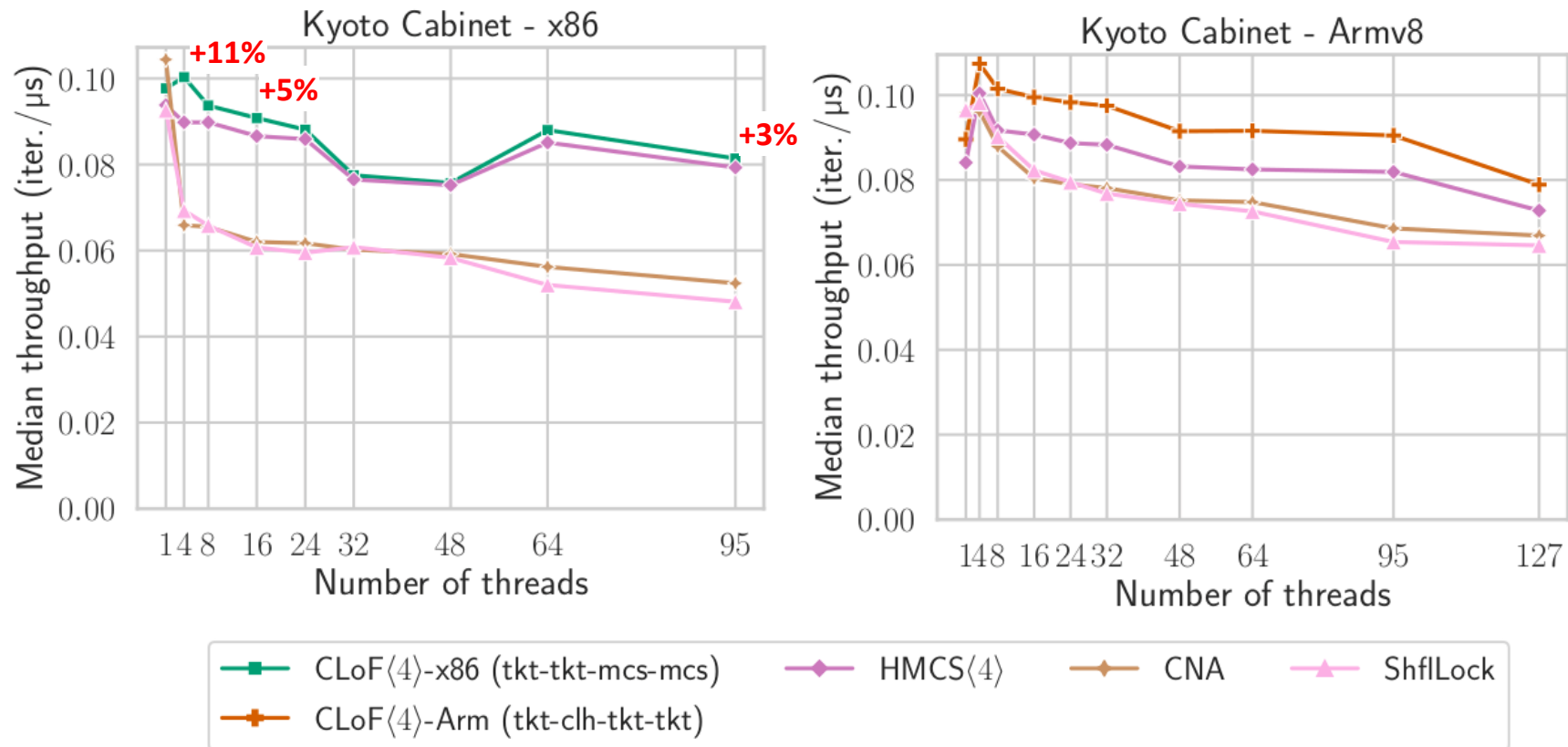
# Evaluation

- Another benchmark, called Kyoto Cabinet, is used to cross-validate results
  - Display LC-best



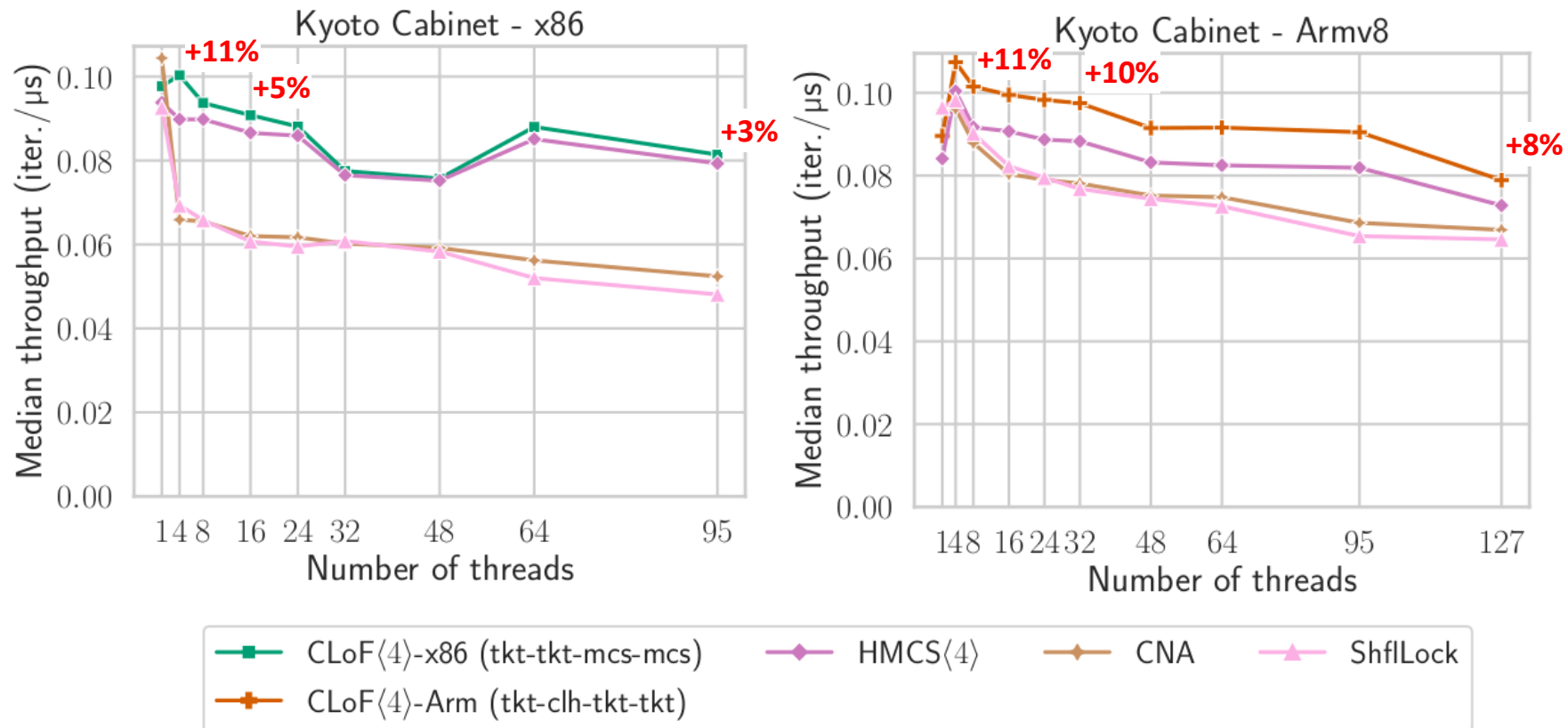
# Evaluation

- Another benchmark, called Kyoto Cabinet, is used to cross-validate results
  - Display LC-best



# Evaluation

- Another benchmark, called Kyoto Cabinet, is used to cross-validate results
  - Display LC-best



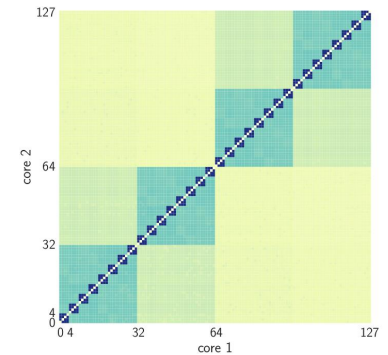
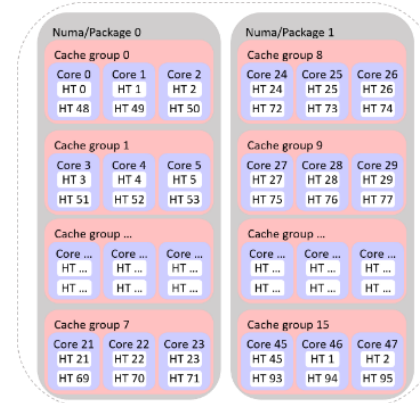
# Conclusion and Future Work

- CLoF locks
  - fully leverage deep hierarchy
  - are level heterogeneity
  - can be optimized for target platform
  - are correct-by-construction on Weak Memory Models
- Don't miss the details!
  - platform-specific optimizations
  - analysis of lock combinations
  - ...
- Future work
  - CLoF in the Linux kernel
  - big.LITTLE platforms

Thanks!

# Tuning Point: Choosing the Hierarchy Levels

- Not all levels will always be used
  - Application can disable hyperthreads – x86 server
- Some levels may have small improvement – package level on Kunpeng 920



# Tuning Point: Choosing the Hierarchy Levels

- Not all levels will always be used
  - Application can disable hyperthreads – x86 server
  - Some levels may have small improvement – package level on Kunpeng 920
- User can tune levels that are wanted
  - Include/Remove levels found at discovery

