# Rethinking Secure NVMs in the Age of CXL

**Samuel Thomas**

# CXL is here!

Intel, Google and others join forces for CXL interconnect

CXL is designed to create a high-speed, low latency interconnect between the CPU and workload accelerators

March 13, 2019   By: Will Calvert   💬 Have your say

FINALLY, A COHERENT INTERCONNECT STRATEGY: CXL ABSORBS GEN-Z

November 23, 2021   Timothy Prickett Morgan

Compute Express Link (CXL) 3.0 Debuts, Wins CPU Interconnect Wars

News   By Paul Alcorn last updated August 2, 2022

CXL emerges as the clear winner of the CPU interconnect wars.

# CXL is here!

Intel, Google and others join for... for CXL interconnect

CXL is designed to create a high-speed, low latency interconnect between the CPU and workload accelerators

March 13, 2019   By: Will Calvert   💬 Have your say

**FINALLY, A COHERENT INTERCONNECT STRATEGY: CXL ABSORBS GEN-Z**

November 23, 2021   Timothy Prickett Morgan

Compute Express Link (CXL) 3.0 Debuts, Wins CPU Interconnect Wars

News   By Paul Alcorn last updated August 2, 2022

CXL emerges as the clear winner of the CPU interconnect wars.

# CXL is here!

Samsung Electronics Introduces Industry's First 512GB CXL Memory Module

Korea on May 10, 2022

Intel Sapphire Rapids CXL with Emmitsburg PCH Shown at SC21

By Patrick Kennedy - December 7, 2021

# CXL is here!

## Intel, Google and others join forces for CXL interconnect

CXL is designed to create a high-speed, low latency interconnect between the CPU and workload accelerators

March 13, 2019   By: Will Calvert   💬 Have your say

## FINALLY, A COHERENT INTERCONNECT STRATEGY: CXL ABSORBS GEN-Z

November 23, 2021   Timothy Prickett Morgan

## Compute Express Link (CXL) 3.0 Debuts, Wins CPU Interconnect Wars

News   By Paul Alcorn last updated August 2, 2022

CXL emerges as the clear winner of the CPU interconnect wars.

**How should we think about the CXL architecture?**

# CXL is here!

**How should we design systems for CXL?**

**Samsung Electronics Introduces Industry's First 512GB CXL Memory Module**

Korea on May 10, 2022

Intel Sapphire Rapids Emmitsburg PCH Show

By **Patrick Kennedy**  -  December 7, 2021

**Pond: CXL-Based Memory**

**Demyst Genuine CXL**

**TPP: Transparent Page Placement for CXL-Enabled Tiered-Memory**

| Hasan Al Maruf<br>University of Michigan<br>USA | Hao Wang<br>NVIDIA<br>USA | Abhishek Dhanotia<br>Meta Inc.<br>USA |
| Johannes Weiner<br>Meta Inc.<br>USA | Niket Agarwal<br>NVIDIA<br>USA | Pallab Bhattacharya<br>NVIDIA<br>USA |
| Chris Petersen<br>Meta Inc.<br>USA | Mosharaf Chowdhury<br>University of Michigan<br>USA | Shobhit Kanaujia<br>Meta Inc.<br>USA |
| | Prakash Chauhan<br>Meta Inc.<br>USA | |

# Emerging Technologies

# Outline



Applications influence architecture!

1. Secure Memory Architecture → 2. Secure NVMs → 3. Application-Awareness → 4. Changes due to CXL

# Outline

Applications influence architecture!

1. *Secure Memory Architecture*

2. Secure NVMs

3. Application-Awareness

4. Changes due to CXL

arr[i] = n;

main() {
    int x = arr[i];
}

arr[i] = n;

main() {
    int x = arr[i];
}

main() {
    while (1)
        clflush(arr[i]);
}

CPU

Memory

n

arr[i] = n;

main() {
    int x = arr[i];
}

main() {
    while (1)
        clflush(arr[i]);
}

CPU

Memory

n

5

arr[i] = n;

main() {
    int x = arr[i];
}

main() {
    while (1)
        clflush(arr[i]);
}

CPU

Memory

n'

5

On-Chip, Trusted

Off-Chip, Untrusted

**CPU**

**Memory**

On-Chip, Trusted

Off-Chip, Untrusted

CPU

H(n)

n

Memory

On-Chip, Trusted

Off-Chip, Untrusted

n

CPU

H(n)

Memory

# On-Chip, Trusted

# Off-Chip, Untrusted

**CPU**

n

...

...

...

**Memory**

7

On-Chip, Trusted

Off-Chip, Untrusted

CPU

n'

...

...

...

H(n')

H(...)

H(...)

H(...)

Memory

On-Chip, Trusted

**CPU**

Off-Chip, Untrusted

H(a)    H(b)         H(c)    H(d)

On-Chip, Trusted

CPU

root

Off-Chip, Untrusted

H(a,b)

H(c,d)

H(a)

H(b)

H(c)

H(d)

On-Chip, Trusted

CPU

Fetch d!

root

Off-Chip, Untrusted

H(a,b)

H(c,d)

H(a)

H(b)

H(c)

H(d)

On-Chip, Trusted

CPU

Fetch d!

root

Off-Chip, Untrusted

H(a,b)

H(c,d)

H(c,d)

H(a)

H(b)

H(c)

H(d)

H(d)

On-Chip, Trusted

CPU

root

Metadata Cache

Off-Chip, Untrusted

H(a,b)

H(c,d)
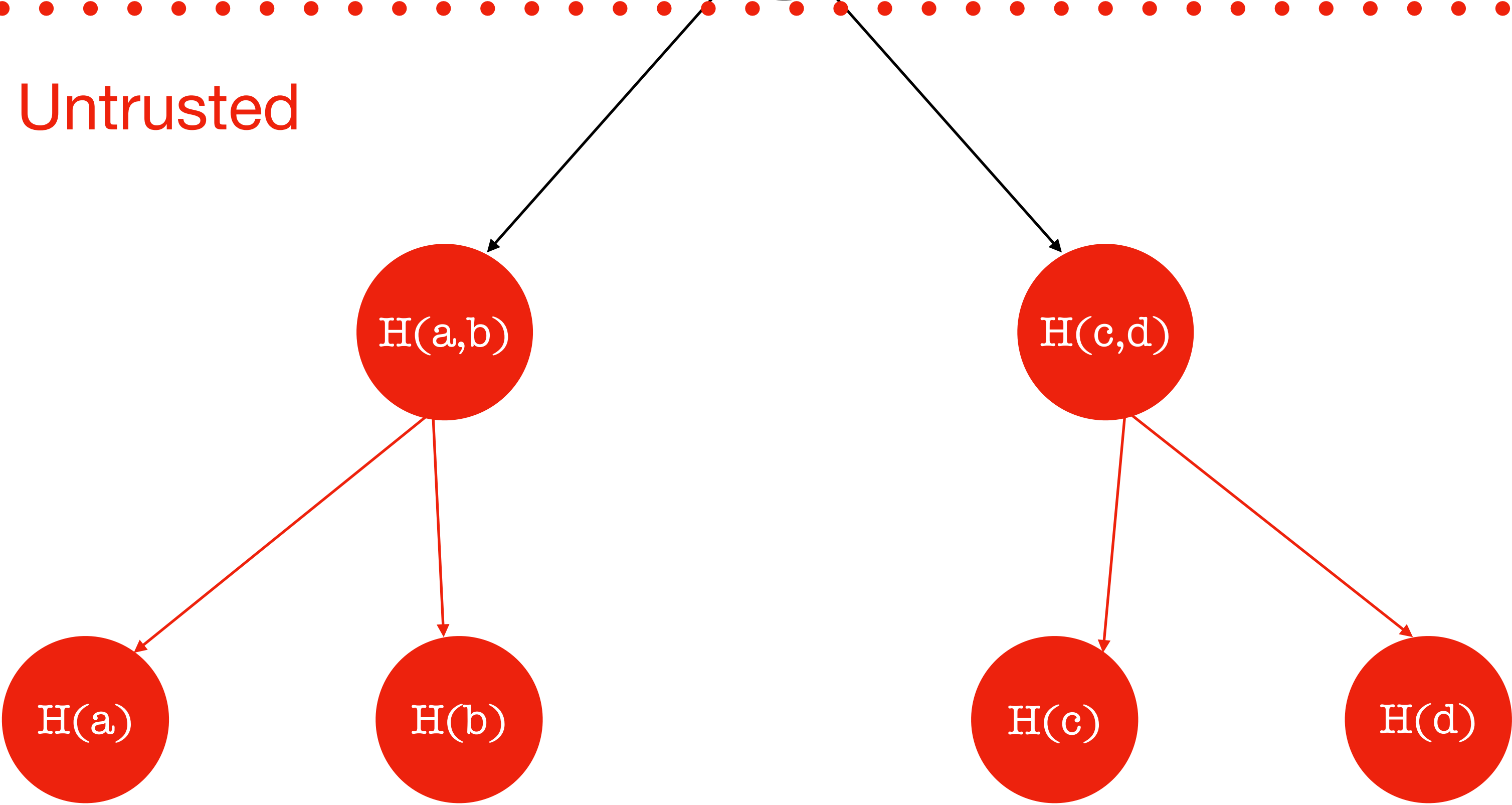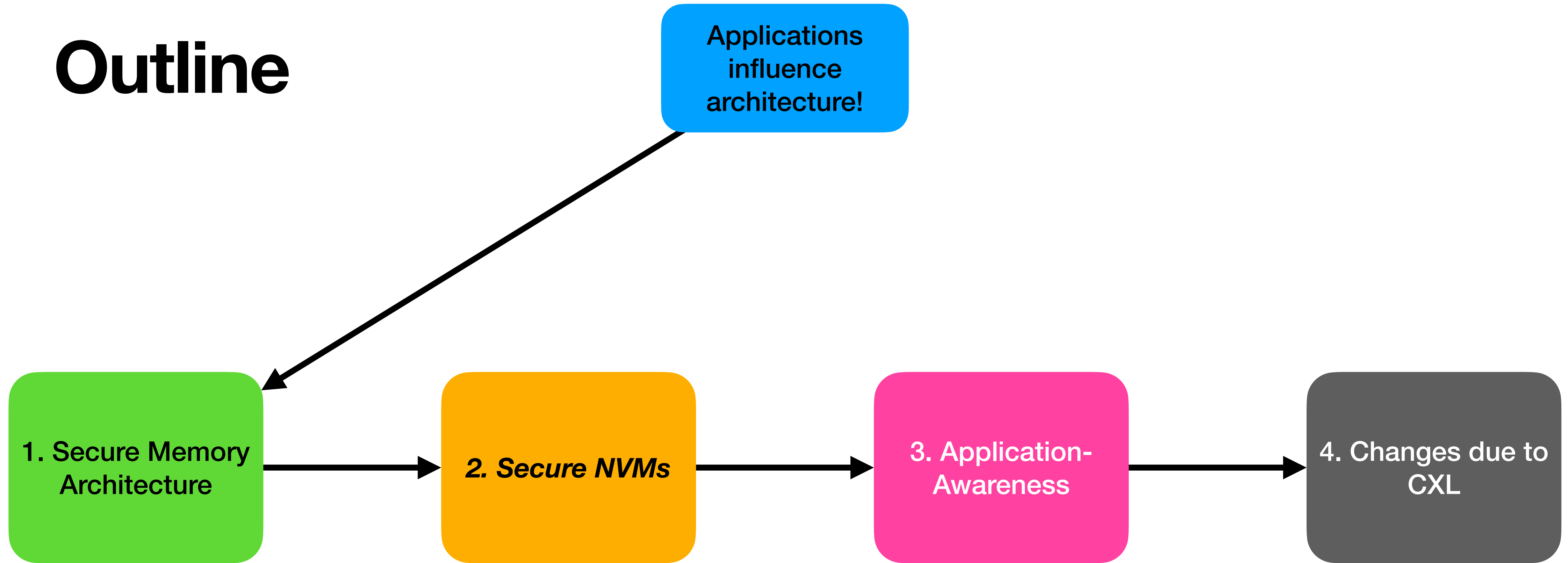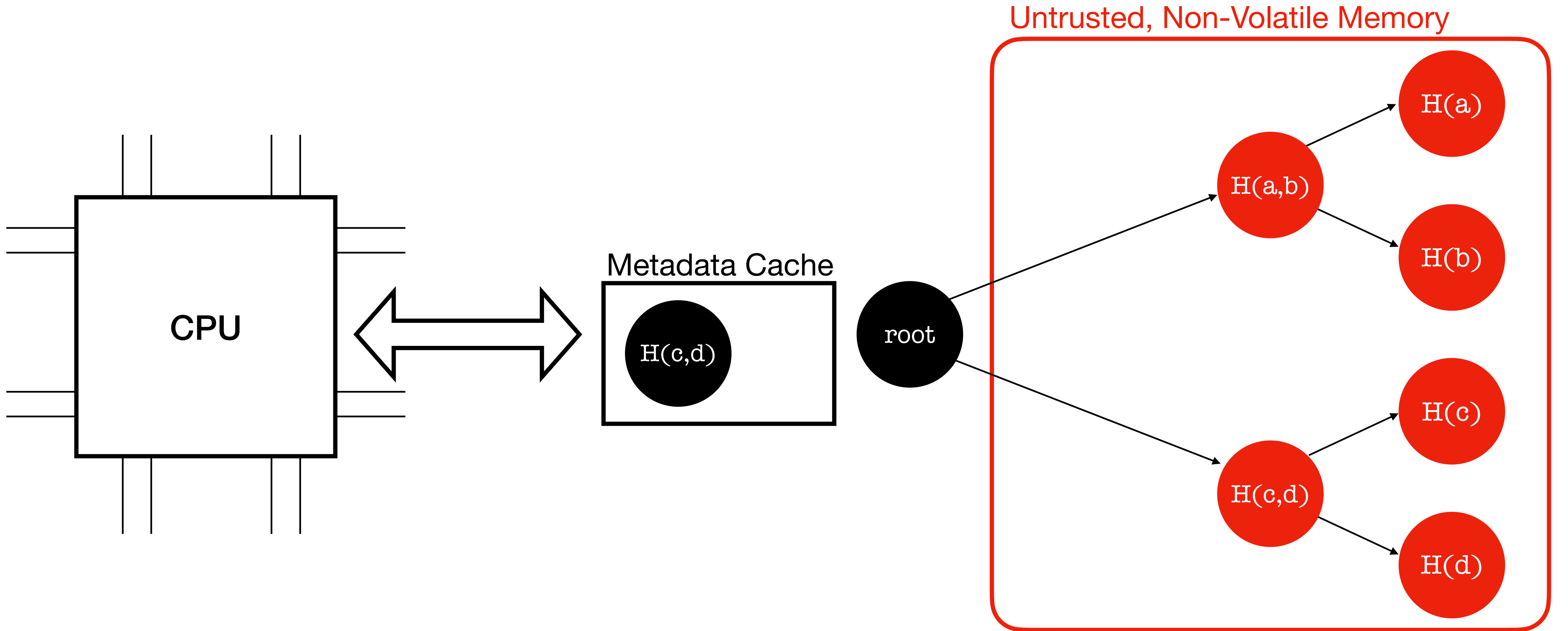
H(a)

H(b)

H(c)

H(d)

On-Chip, Trusted

CPU

root

Metadata Cache

H(c,d)
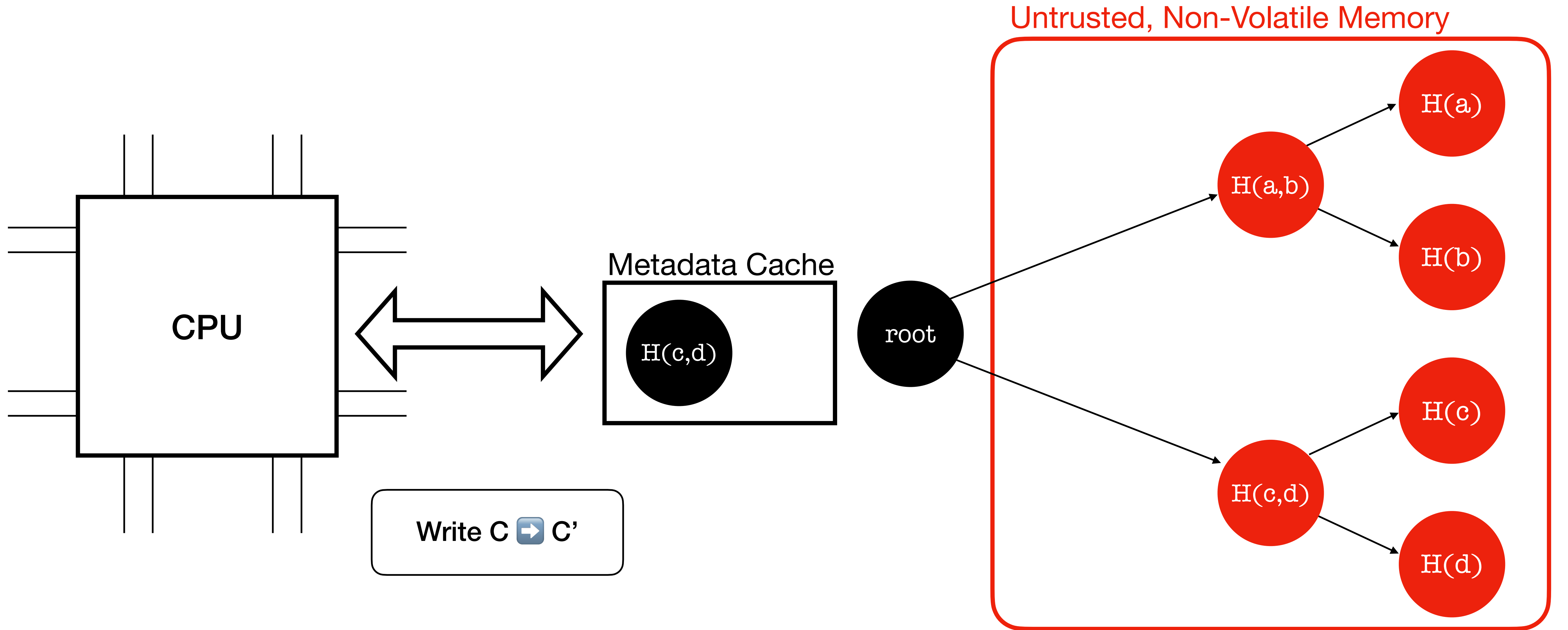
Off-Chip, Untrusted

H(a,b)

H(c,d)

H(a)

H(b)

H(c)

H(d)

On-Chip, Trusted

CPU

root

Metadata Cache

H(c,d)

Off-Chip, Untrusted

H(a,b)

H(c,d)

H(a)

H(b)

H(c)

H(d)

# Outline

Applications influence architecture!

1. Secure Memory Architecture

2. *Secure NVMs*

3. Application-Awareness

4. Changes due to CXL

Untrusted, Non-Volatile Memory

CPU

Metadata Cache

H(c,d)

root

Write C ➡ C'

H(a,b)

H(a)

H(b)

H(c,d)

H(c)

H(d)

CPU

Metadata Cache

H(c',d)

root

H(a,b)

H(a)

H(b)

H(c,d)

H(c')

H(d)

Write C ➡️ C'

Untrusted, Non-Volatile Memory

Metadata Cache

Write C ➡️ C'

CPU

11

Untrusted, Non-Volatile Memory

CPU

Strict Persistence Metadata Cache

H(c',d)

root

Write C ➡ C'

H(a,b)

H(a)

H(b)

H(c',d)

H(c')

H(d)

12

✓ Crash Consistency

Untrusted, Non-Volatile Memory

CPU

Strict Persistence
Metadata Cache

H(c',d)
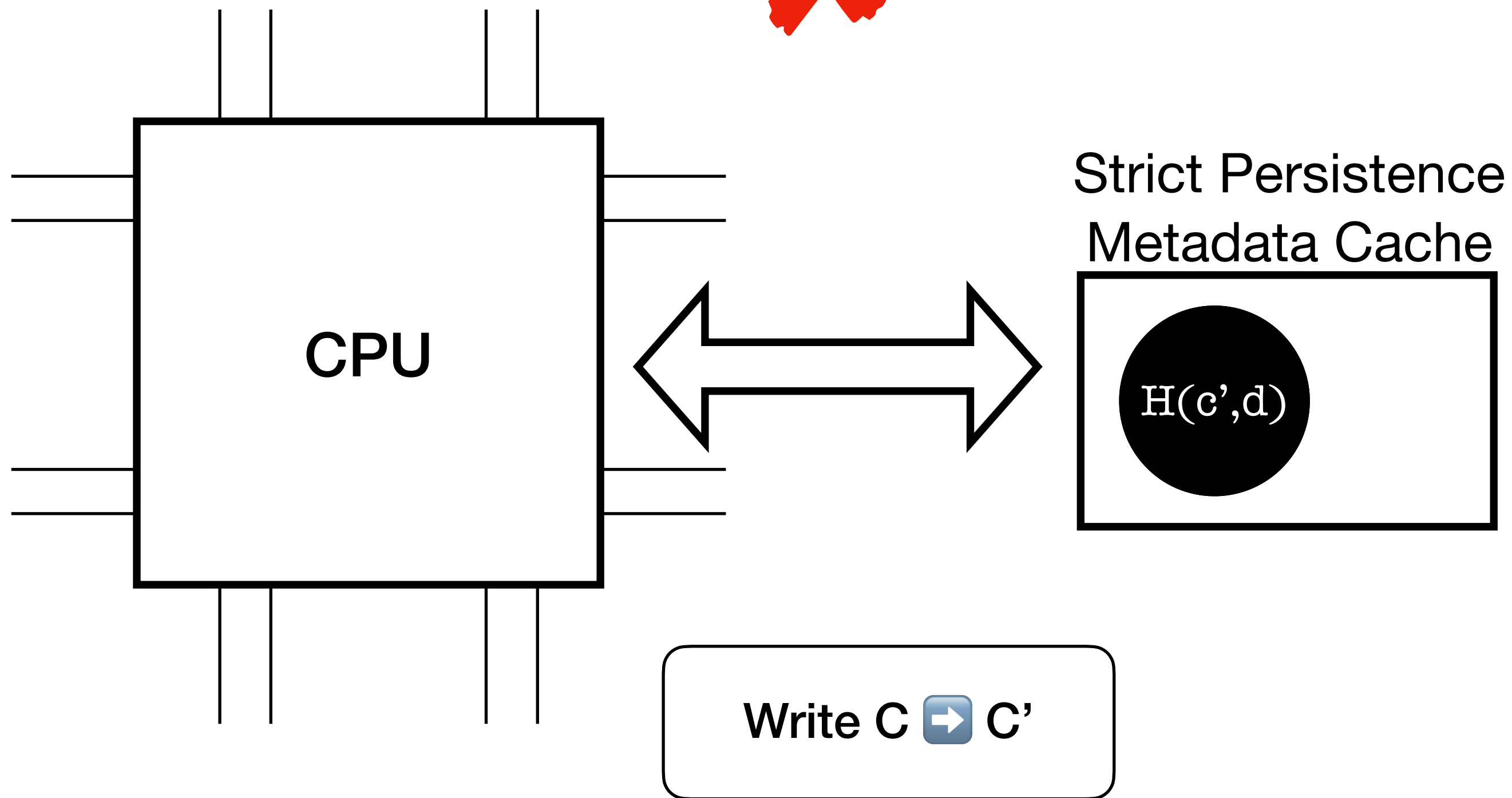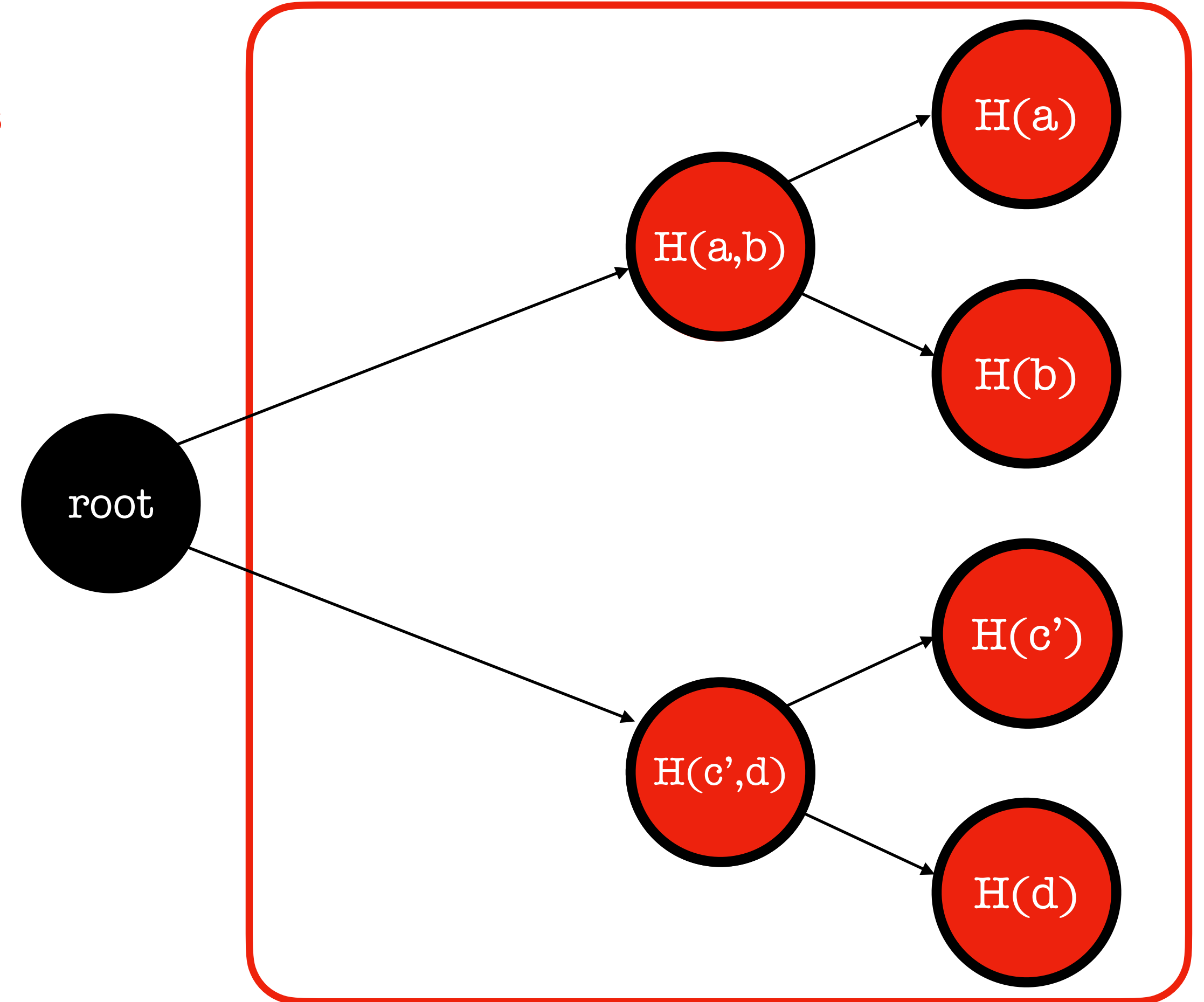
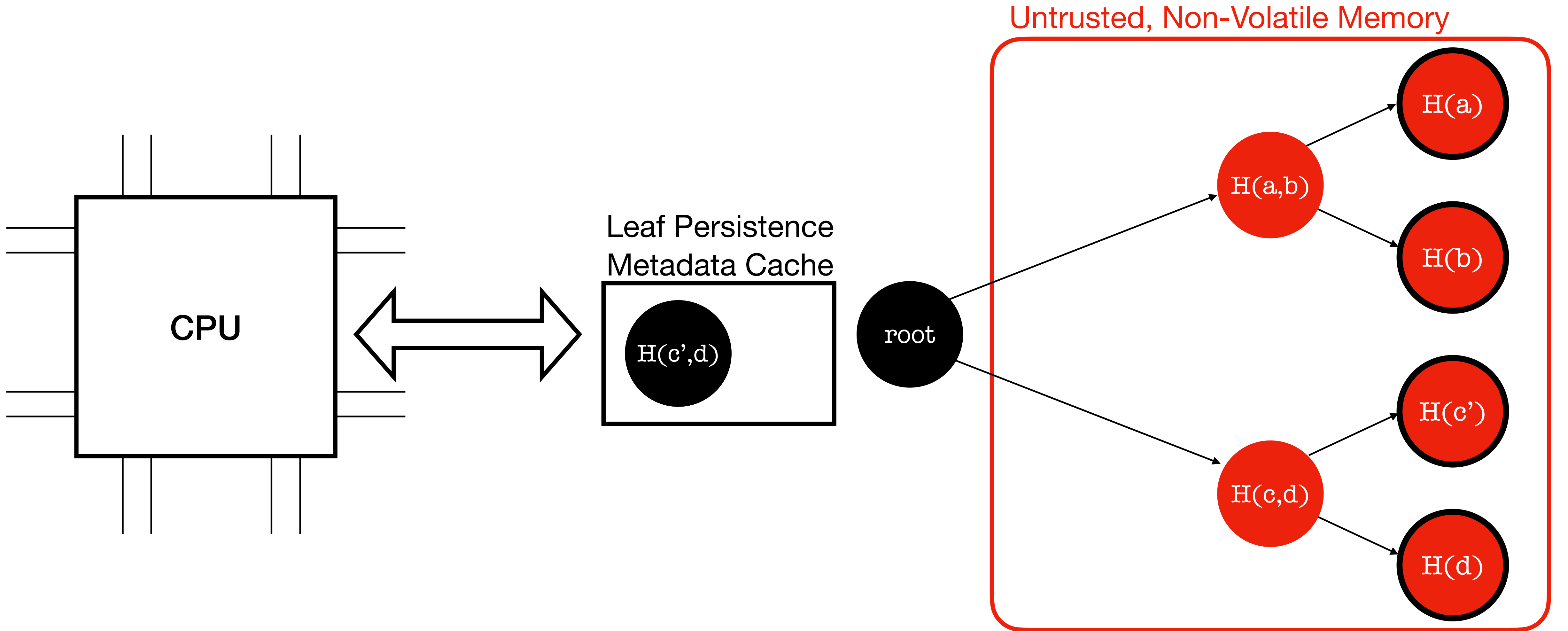Write C ➡ C'

root

H(a,b)

H(a)

H(b)

H(c',d)

H(c')

H(d)

✅ Crash Consistency
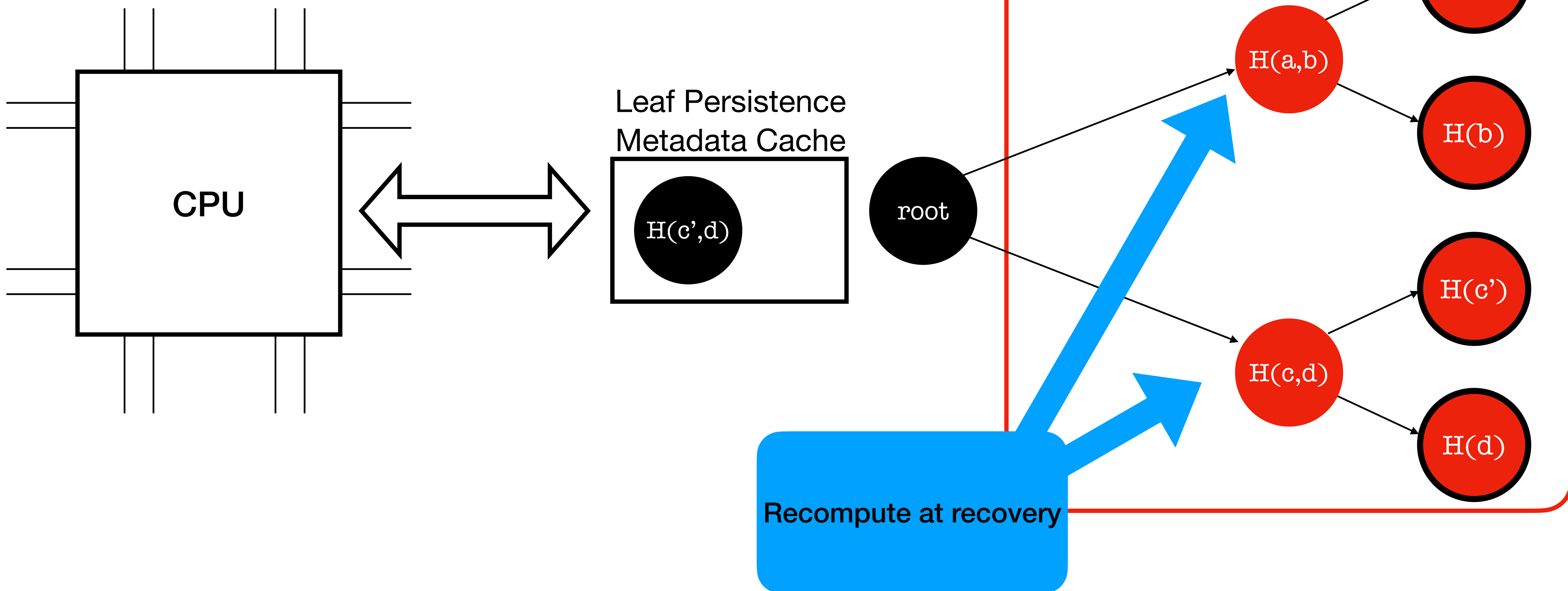
❌ Performance Limitations

Untrusted, Non-Volatile Memory

CPU

Strict Persistence
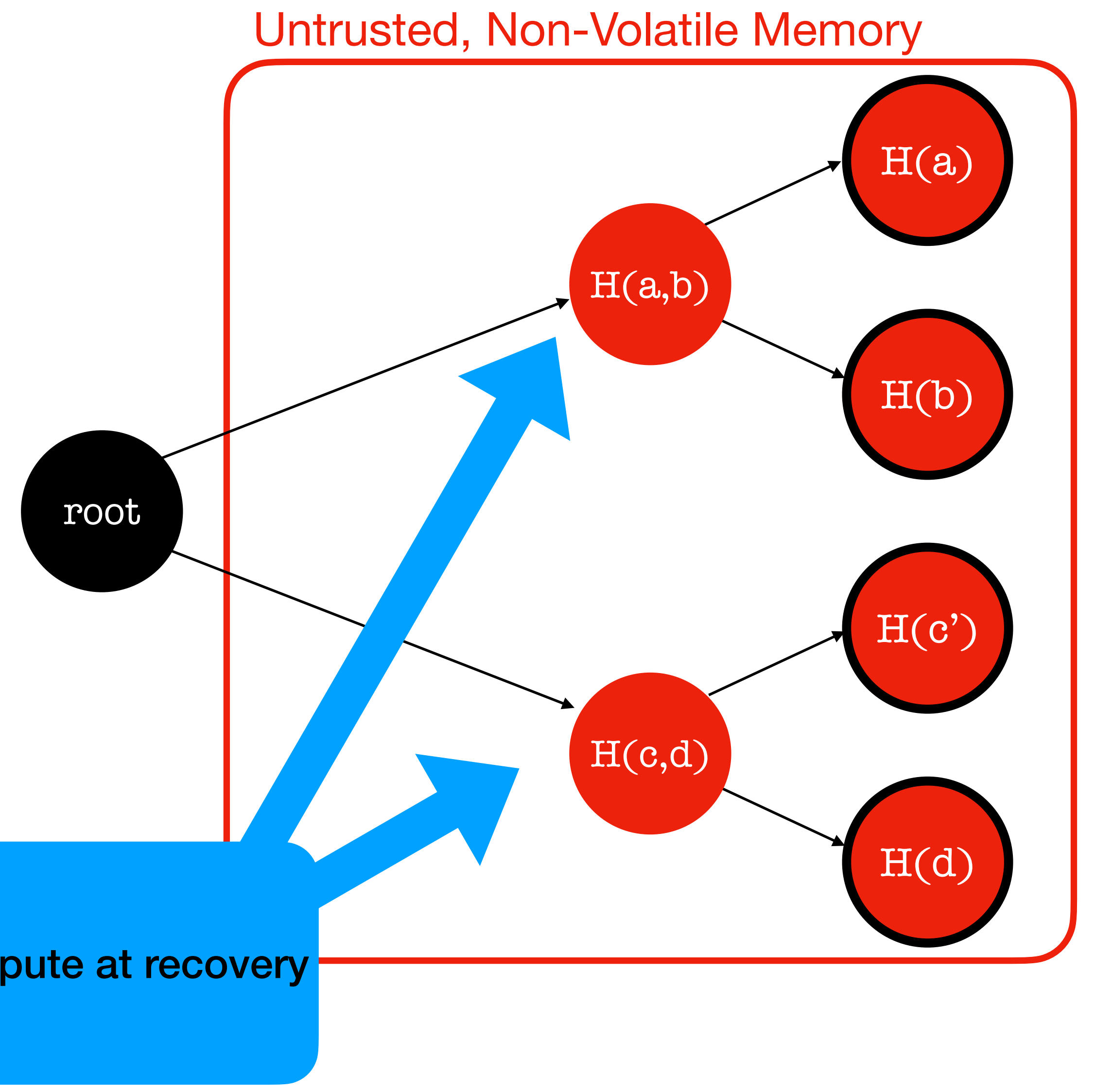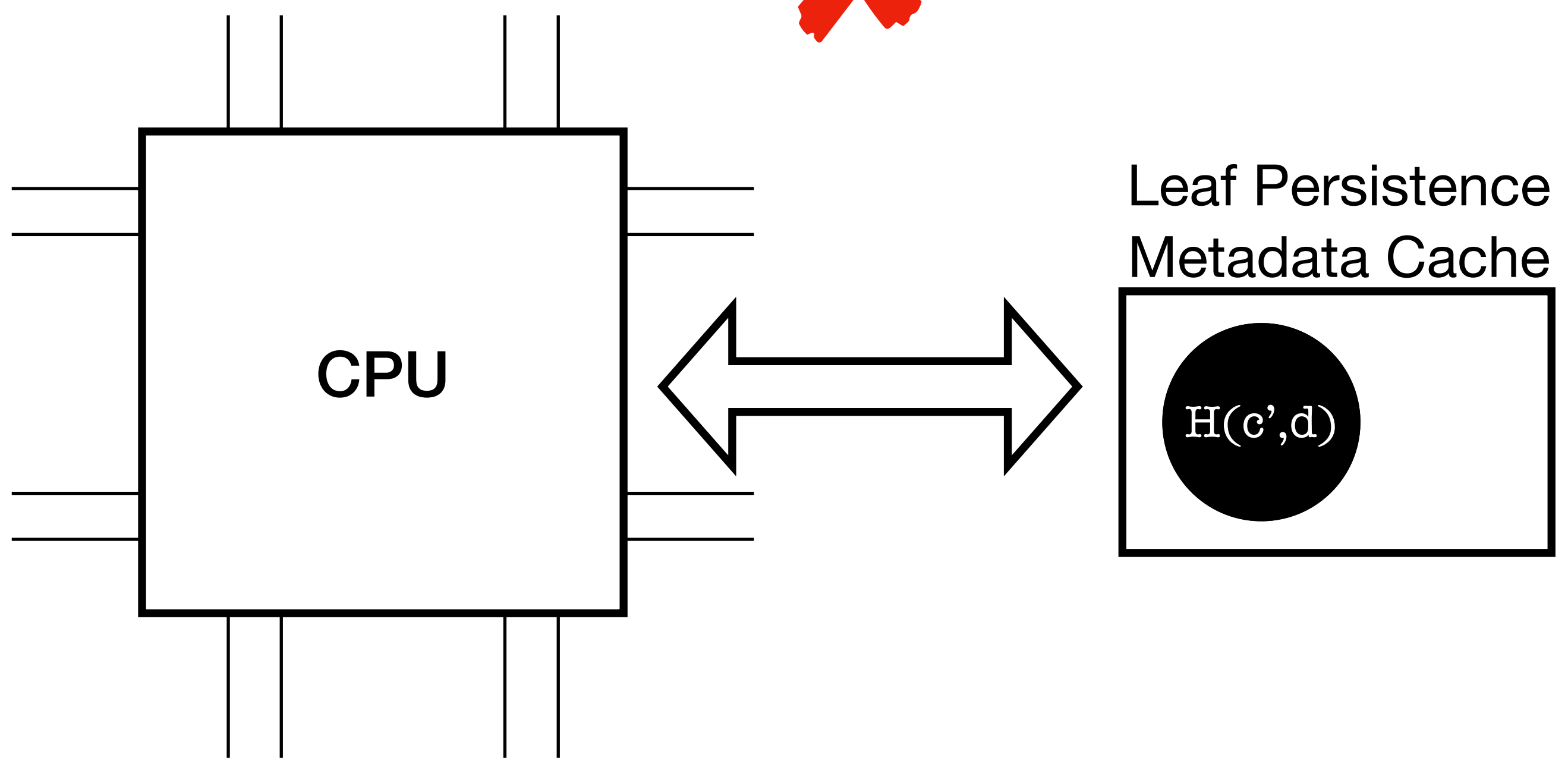Metadata Cache

H(c',d)

root

Write C ➡️ C'

H(a,b)

H(a)

H(b)

H(c',d)

H(c')

H(d)

13

# Outline

Applications influence architecture!

1. Secure Memory Architecture → 2. Secure NVMs → *3. Application-Awareness* → 4. Changes due to CXL

Samuel Thomas
*BROWN*

Kidus Workneh, Joseph Izraelevitz, Tamara Lehman
*University of Colorado Boulder*
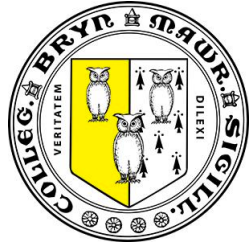
Jac McCarty

R. Iris Bahar

QR code to the paper!

Performance ⟷ Recovery

Samuel Thomas

Kidus Workneh, Joseph Izraelevitz, Tamara Lehman
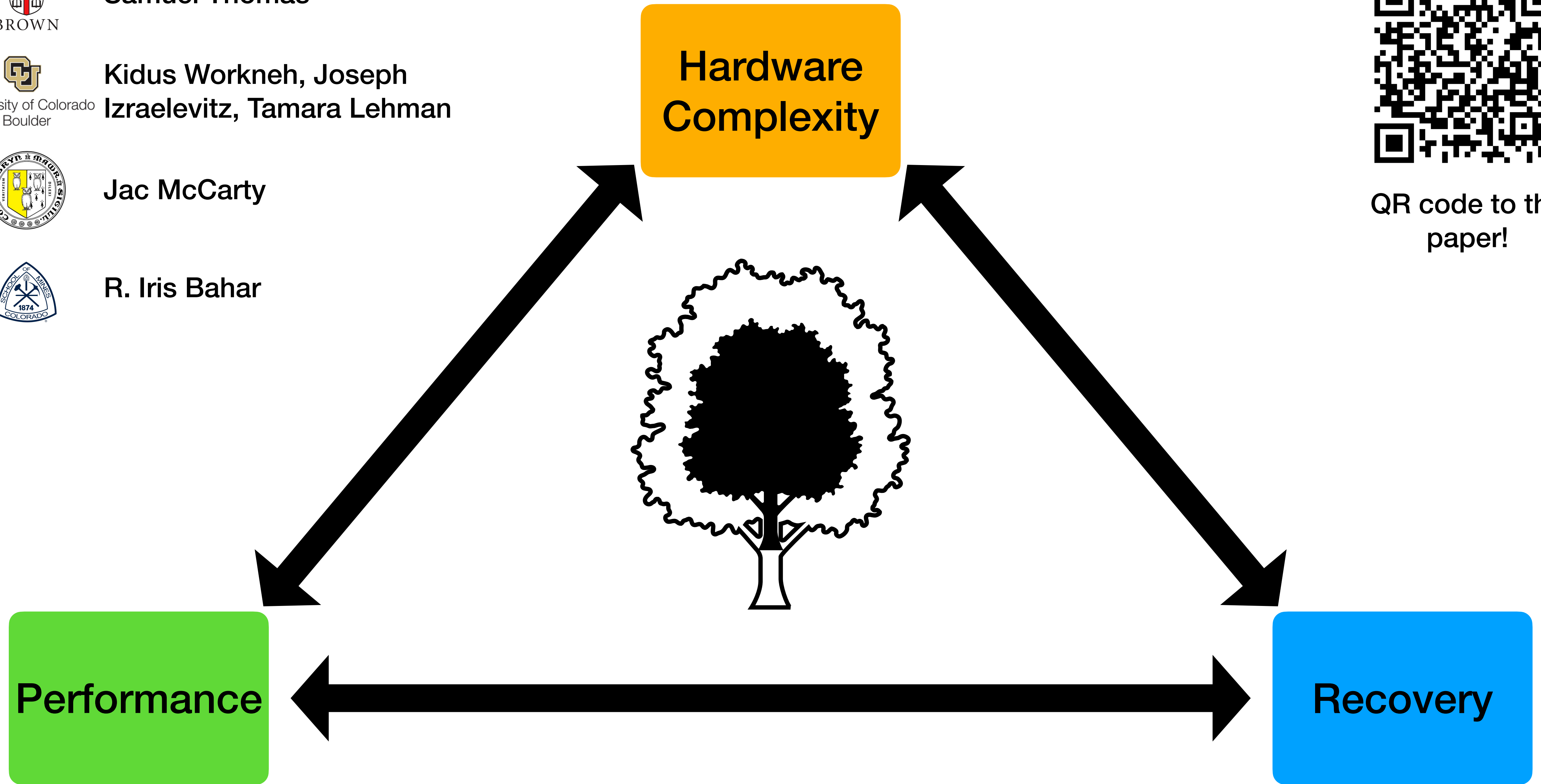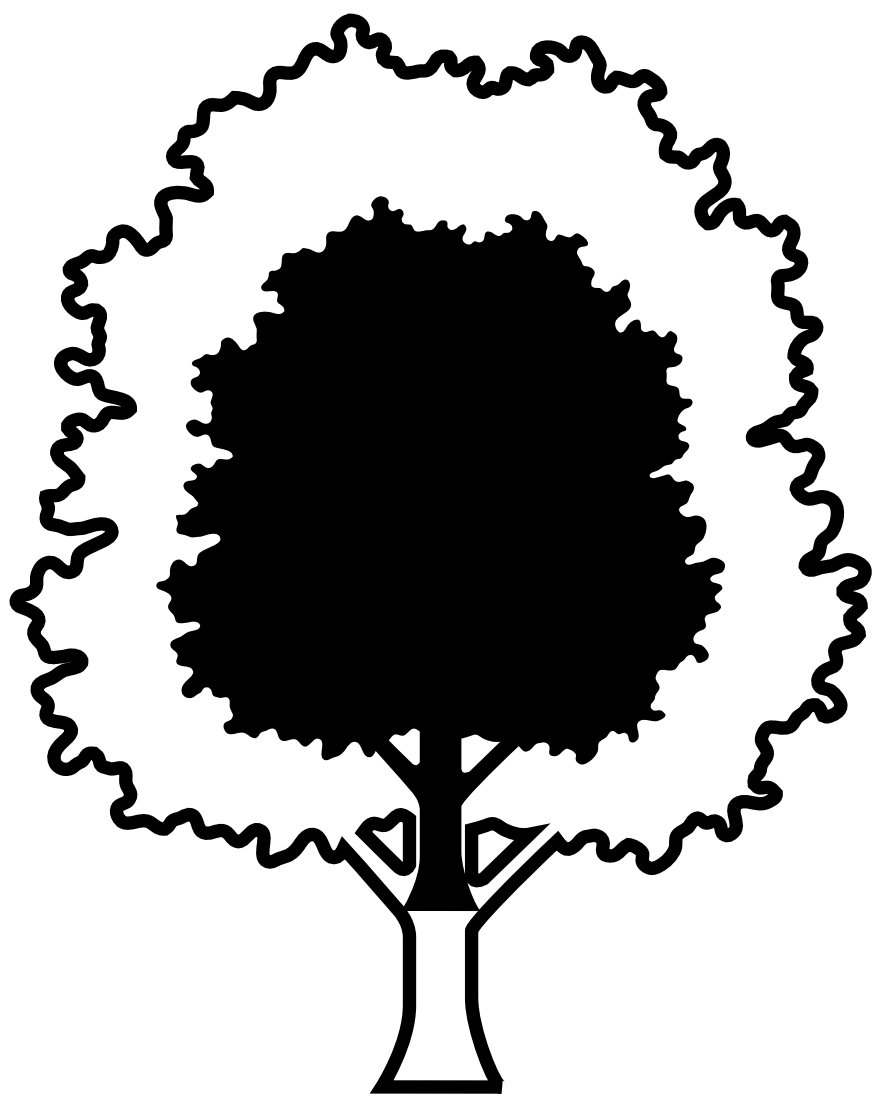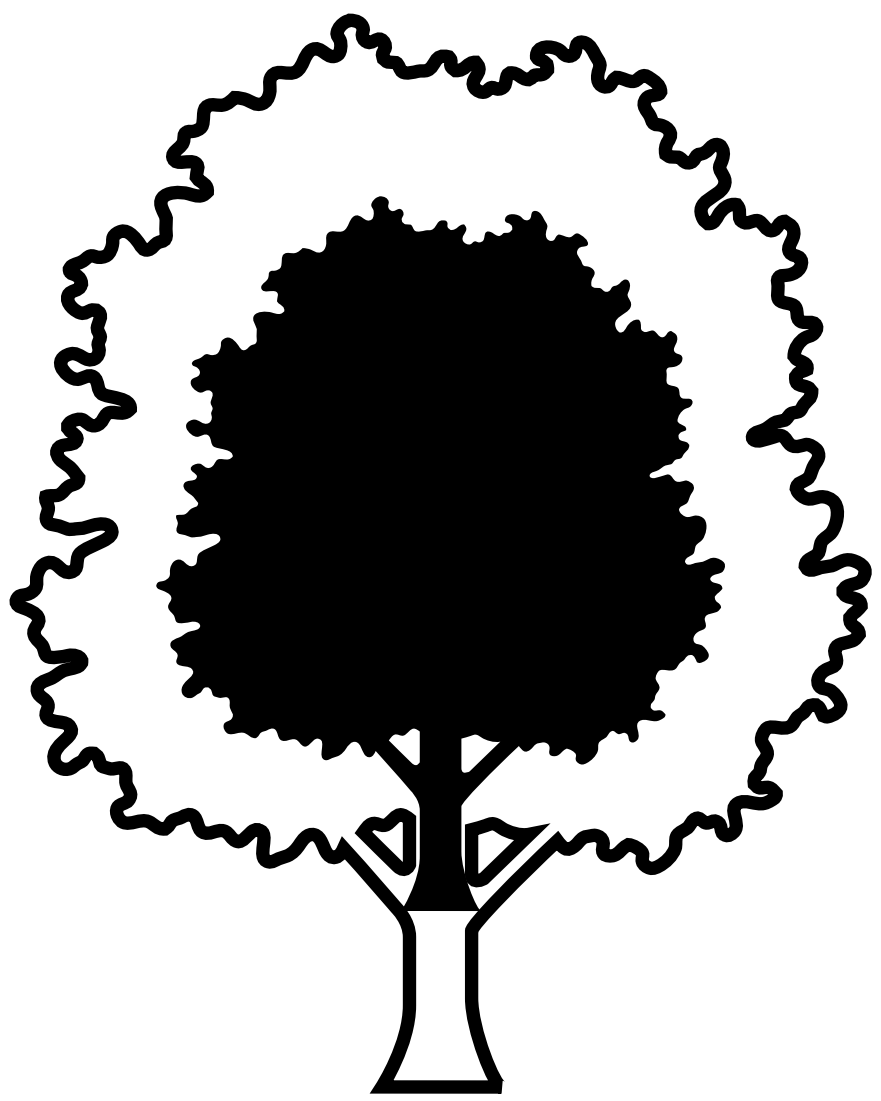
Jac McCarty

R. Iris Bahar

QR code to the paper!

Hardware Complexity

Performance

Recovery

Samuel Thomas

Kidus Workneh, Joseph
Izraelevitz, Tamara Lehman

Jac McCarty

R. Iris Bahar

QR code to the
paper!

**Hardware
Complexity**

**Performance**

**Recovery**

# A Midsummer Night's Tree



Metadata Cache

H(c,d)

root

fast subtree

H(a,b)

H(c,d)

H(a)

H(b)

H(c)

H(d)

# A Midsummer Night's Tree

Rethinking Secure NVMs in the Age of CXL

# A Midsummer Night's Tree



Metadata Cache

H(c,d)

root

fast subtree

Strict Persistence

H(a,b)

H(a) 🔥 H(b)

Leaf Persistence

H(c,d)

H(c) 🥶 H(d)

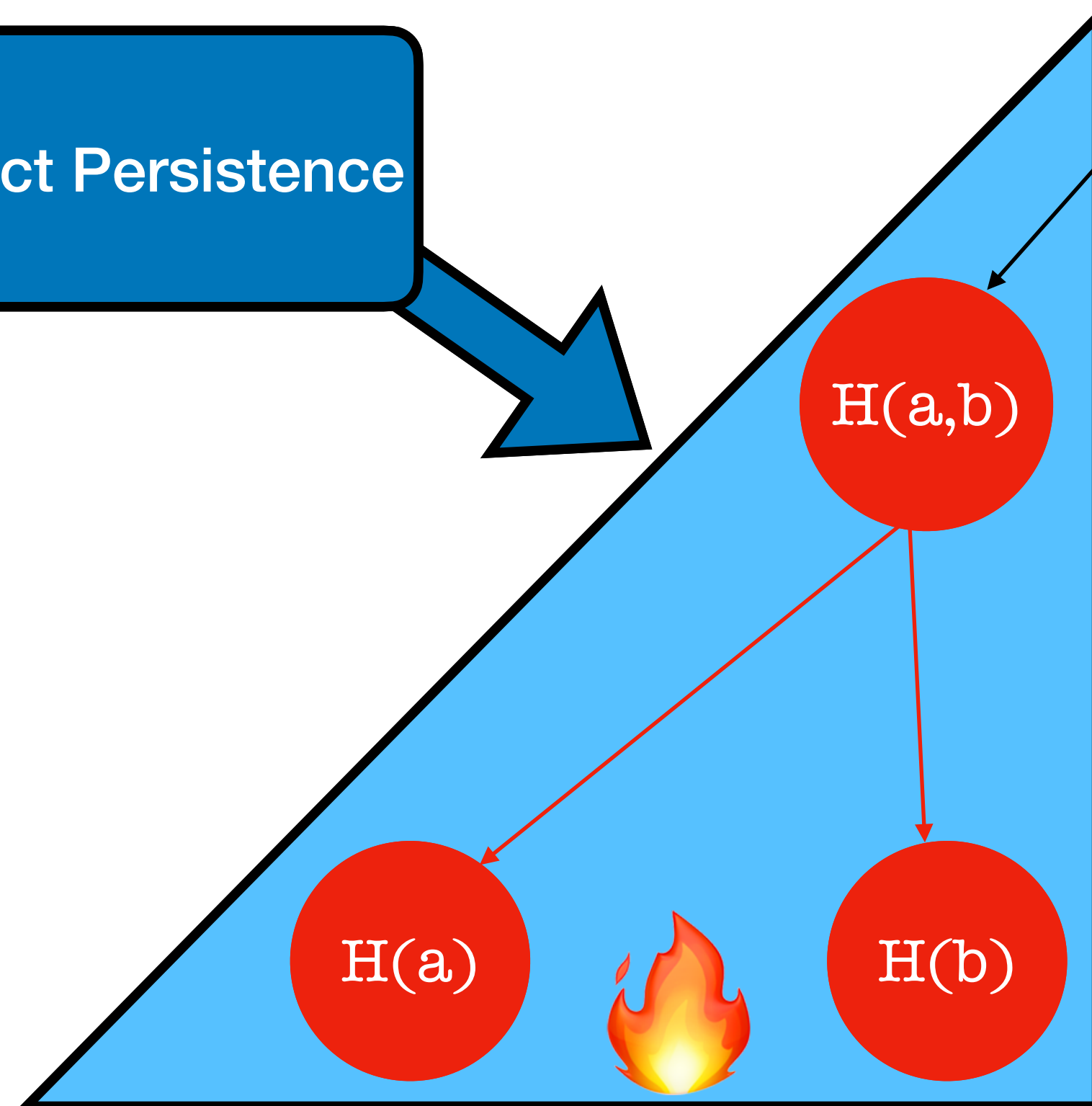# A Midsummer Night's Tree



root

fast subtree

Metadata Cache

H(c,d)

Strict Persistence

H(a,b)

H(a) 🔥 H(b)

Strict Persistence

H(c,d)

H(c) 🥶 H(d)

# A Midsummer Night's Tree

perlbench

# Biased Physical Page Allocation

struct **free_area** {

  ...

};

# Biased Physical Page Allocation



```
struct free_area {

    ...

};
```

18

# Biased Physical Page Allocation

struct **free_area** {

  ...

};

unbiased

Samuel Thomas, et al.

A Midsummer Night's Tree: Efficient and High Performance Secure SCM

biased

# Hardware Overhead

| | Volatile Overhead | Non-Volatile Overhead |
|---|---|---|
| **AMNT** | 96 bytes | 64 bytes |
| Anubis, ISCA19 | 37 kB | 64 bytes |
| BMF, MICRO21 | 768 bytes | 4 kB |

# Performance Overhead

Samuel Thomas, et al.

A Midsummer Night's Tree: Efficient and High Performance Secure SCM

# Performance Overhead

# Performance Overhead

# Performance Overhead

# Recovery

Intel® Optane™ Persistent Memory 200 Series Enables Fast Tiered Memory, Delivering 32 Percent More Bandwidth on Average[1] with up to 6 TB Total Memory per Socket[2].

| SKU[*] | 128 GB | 256 GB | 512 GB |
|---|---|---|---|
| USER CAPACITY[*] | 126.7 GB | 253.7 GB | 507.7 GB |
| BANDWIDTH 67% READ; 33% WRITE 15W 64B | 1.06 GB/s | 1.41 GB/s | 1.15 GB/s |



|  | 128GB | 256GB | 512GB | ... | 128TB |
|---|---|---|---|---|---|
| 64B words to fetch | 38.3M | 76.7M | 153.4M | ... | 39.3B |
| time to recover (leaf) | 1.89 sec | 2.84 sec | 6.9 sec | ... | 30 min |
| time to recover (AMNT) | 0.03 sec | 0.04 sec | 0.11 sec | ... | 32 sec |

# Outline



Applications influence architecture!

1. Secure Memory Architecture → 2. Secure NVMs → 3. Application-Awareness → 4. Changes due to CXL

CPU

memory bus

24

**CPU**

**Limited by the number of DIMM slots!**

**CPU**

**Limited by the number of DIMM slots!**

Access memory devices via I/O

**CXL Memory**

**CXL Memory**

**CXL Memory**

CPU

Limited by the number of DIMM slots!

Access memory devices via I/O

PCIe (+/- 90ns)

```
→ ~ sudo numactl --hardware
node 0 cpus: 0 1 2 3 4 5 6 7 8 9 10 11
node 0 size: 24124 MB
node 0 free: 23546 MB
node 1 cpus:
node 1 size: 8038 MB      zNUMA
node 1 free: 7999 MB
```
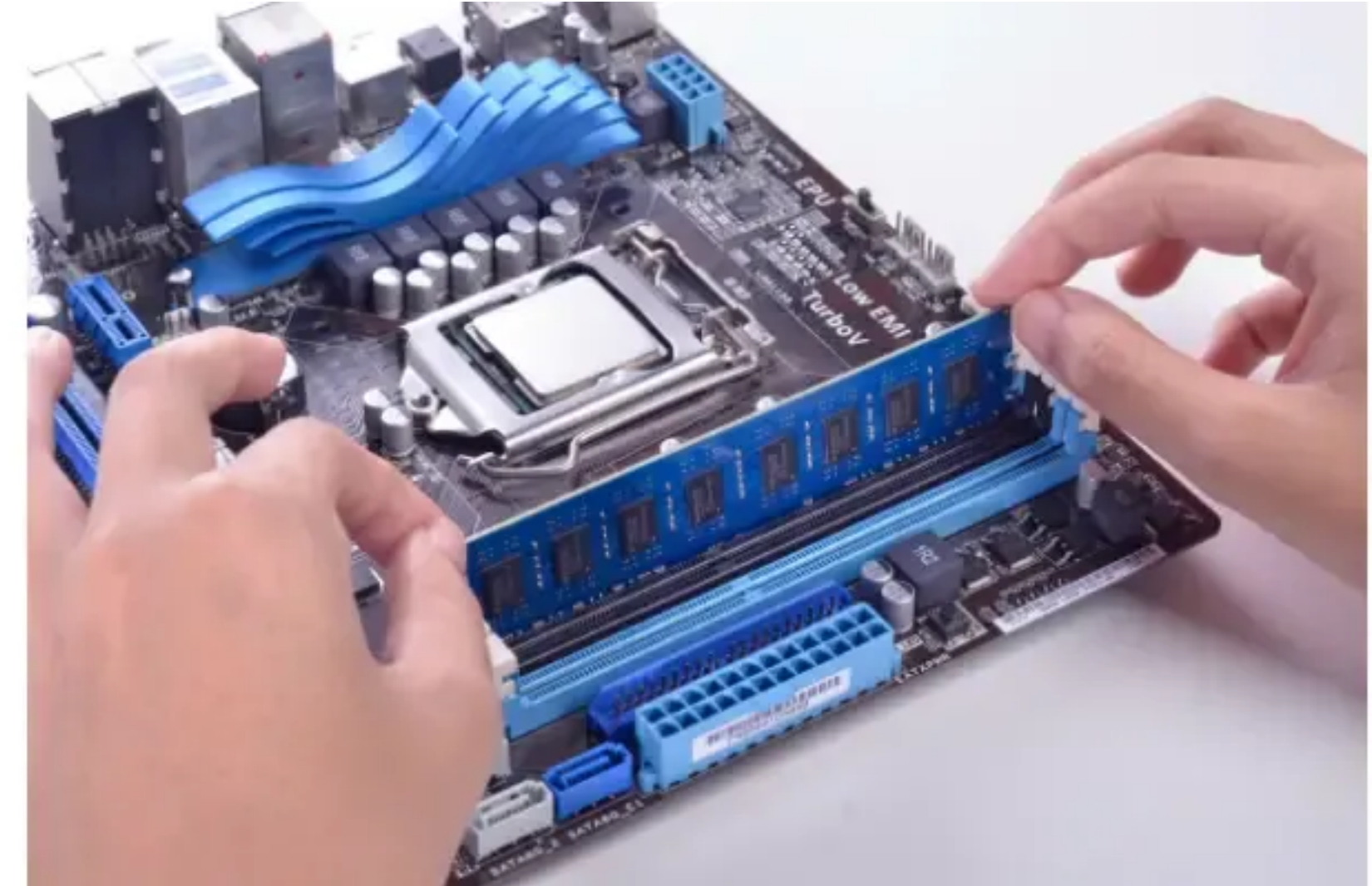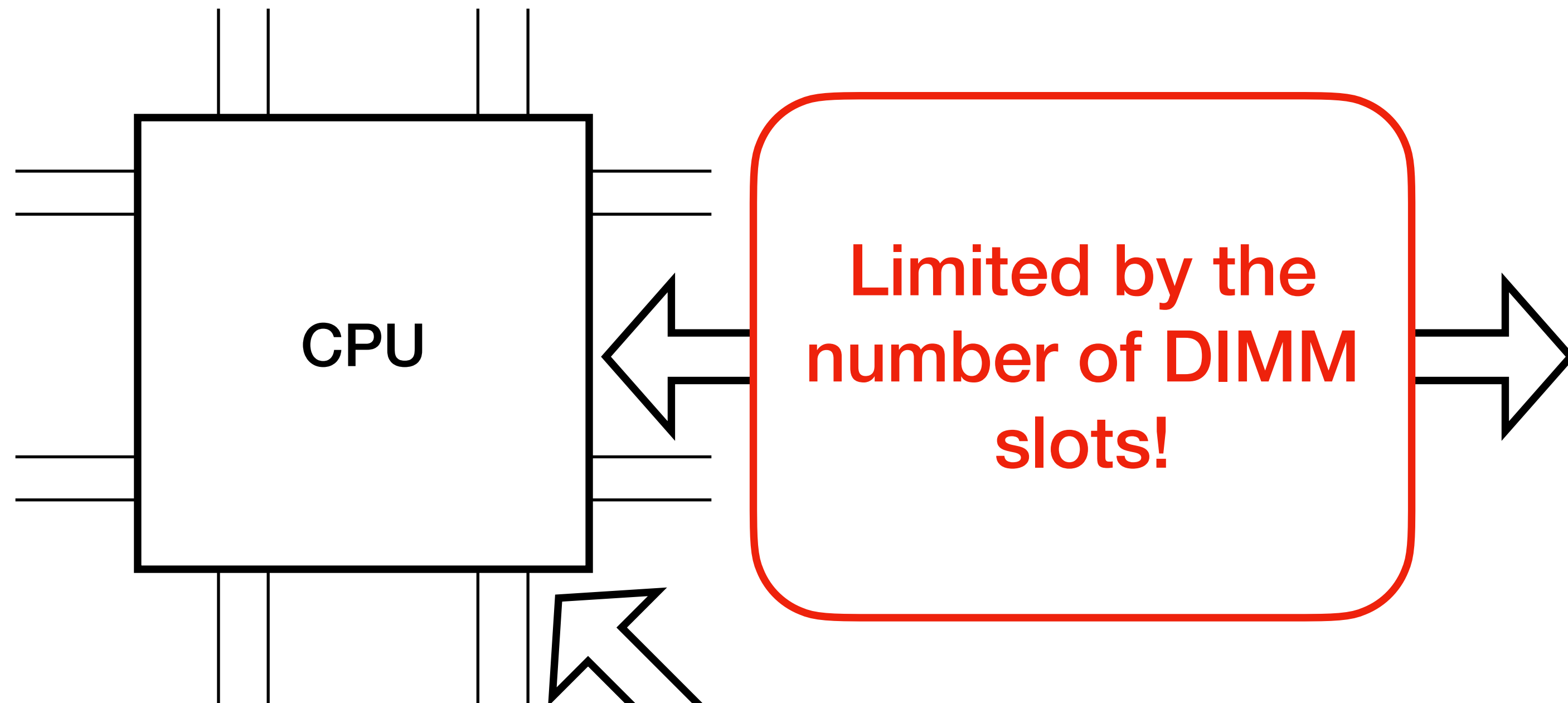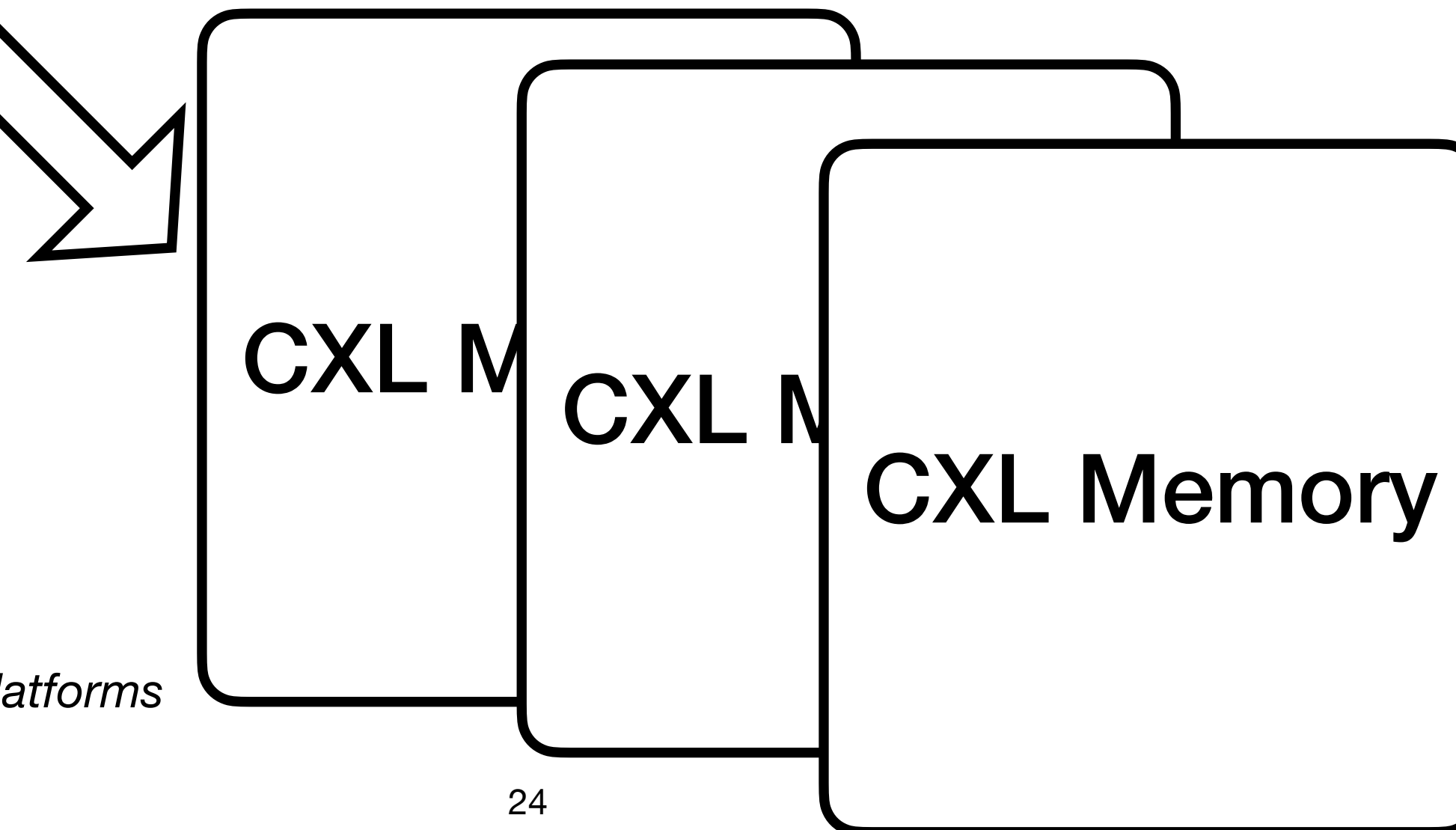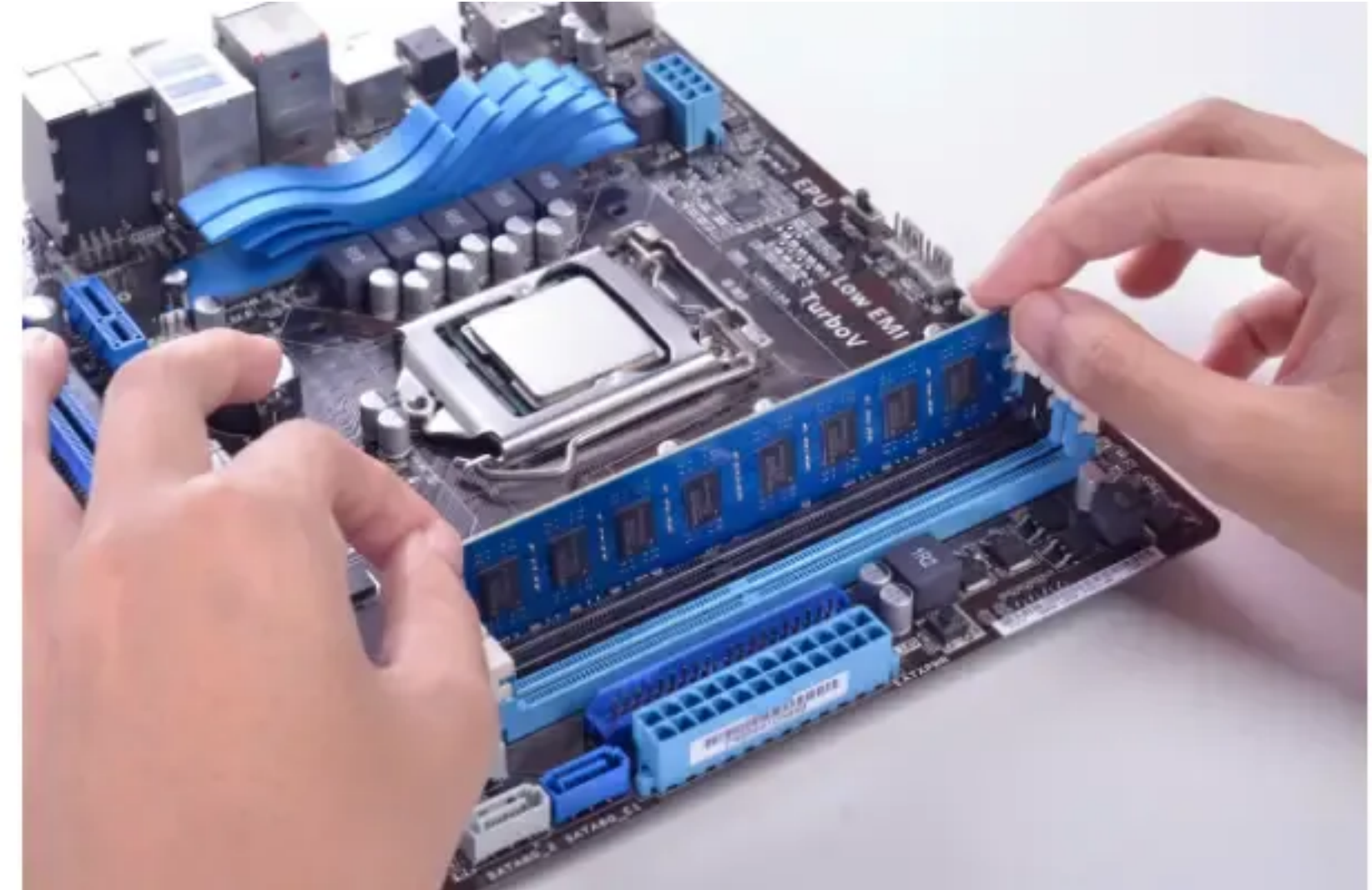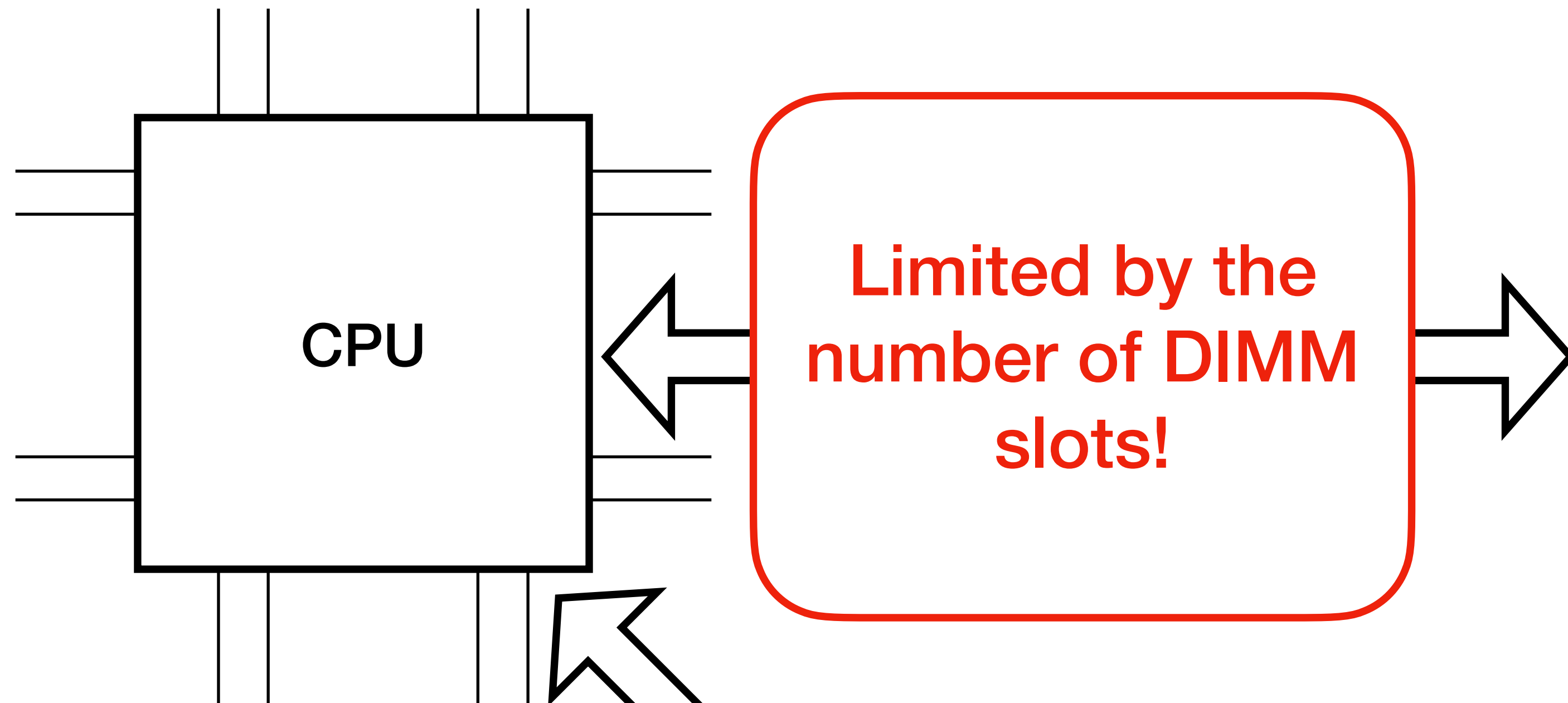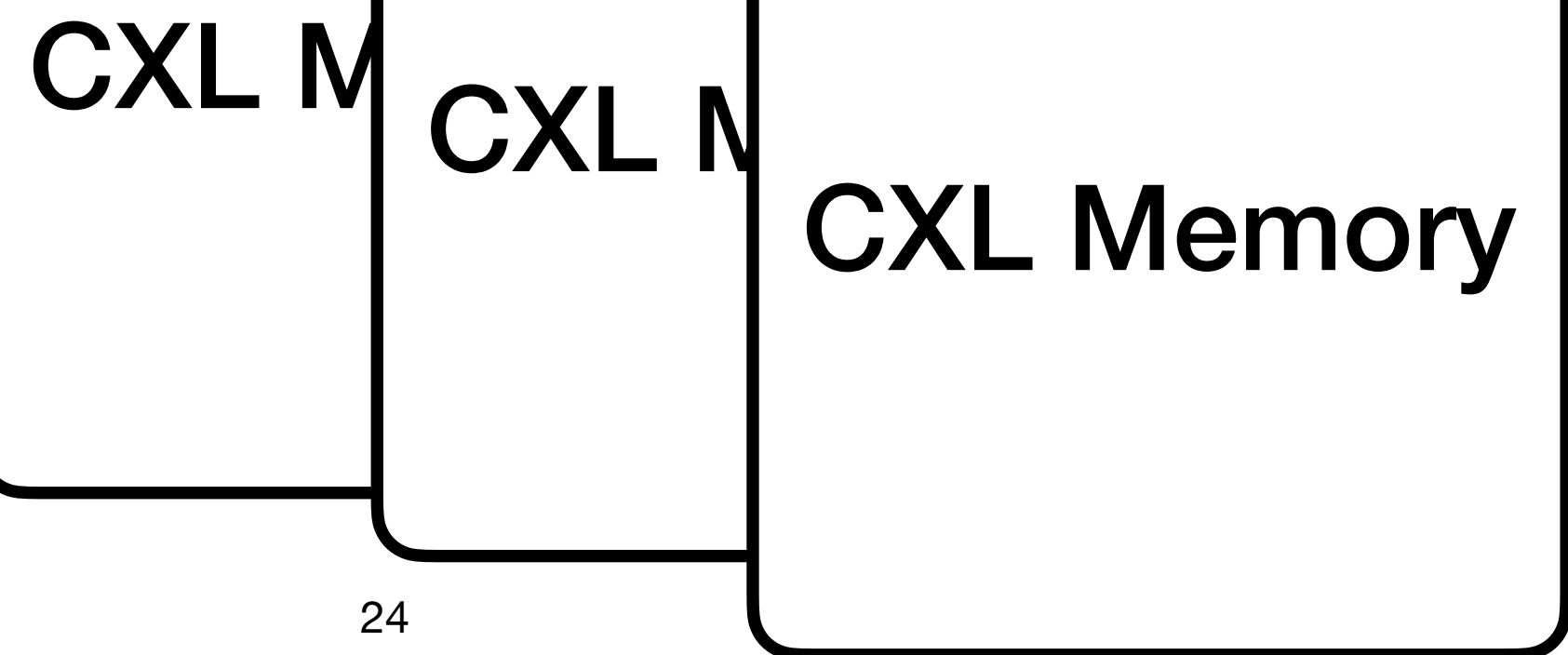
```
node distances:
node   0   1
  0:  10  20
  1:  20  10
```

Li, et al. *Pond: CXL-based memory pooling systems for cloud platforms*

CXL Memory

CXL Memory

CXL Memory

CPU

**Limited by the number of DIMM slots!**

Access memory devices via I/O

Scalable capacity ✓

Increased potential bandwidth ✓

No good programming model (yet)! ✗

PCIe (+/- 90ns)

```
→ ~ sudo numactl --hardware        node distances:
node 0 cpus: 0 1 2 3 4 5 6 7 8 9 10 11   node    0    1
node 0 size: 24124 MB                      0:   10   20
node 0 free: 23546 MB                      1:   20   10
node 1 cpus:
node 1 size: 8038 MB    zNUMA
node 1 free: 7999 MB
```
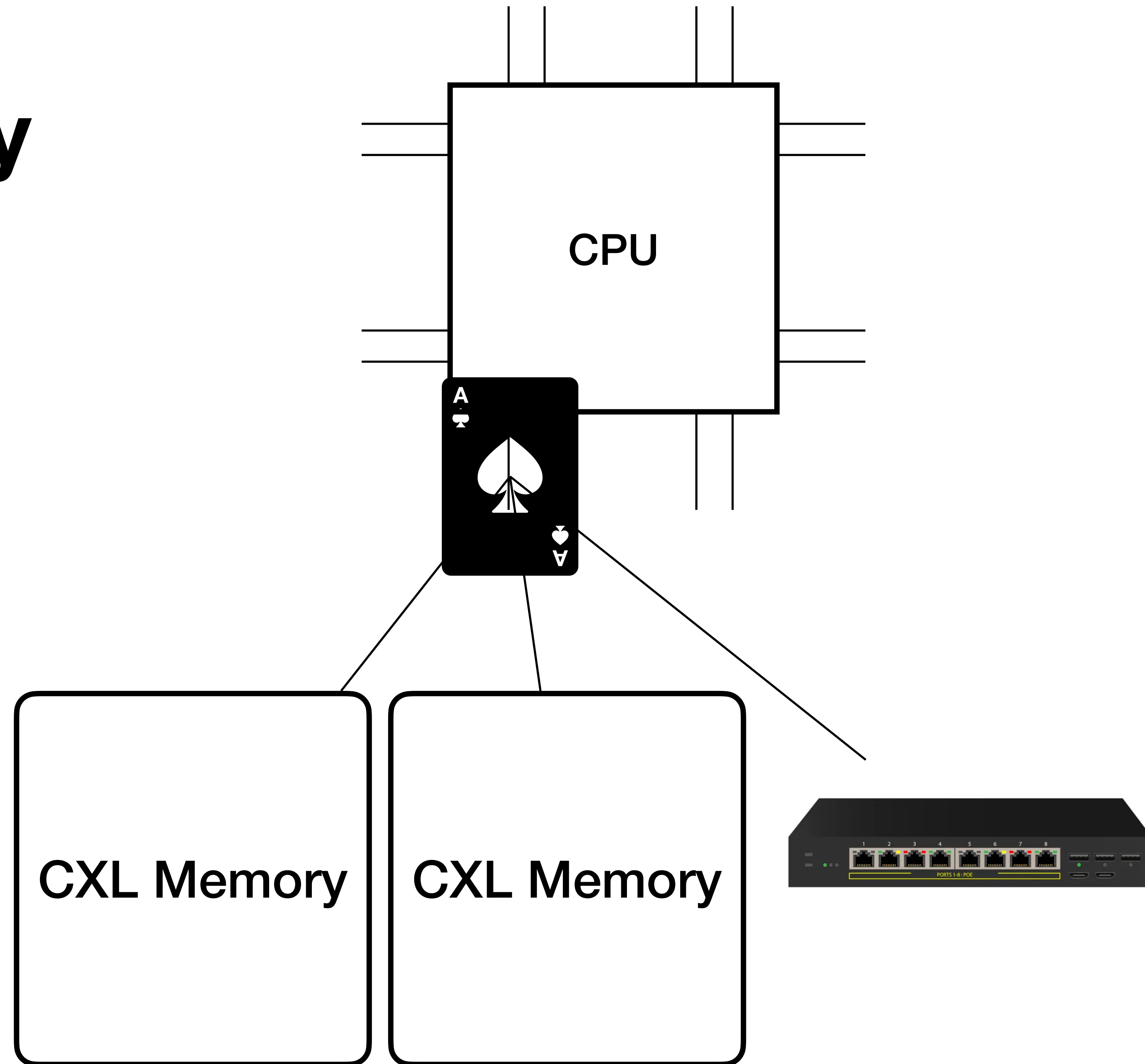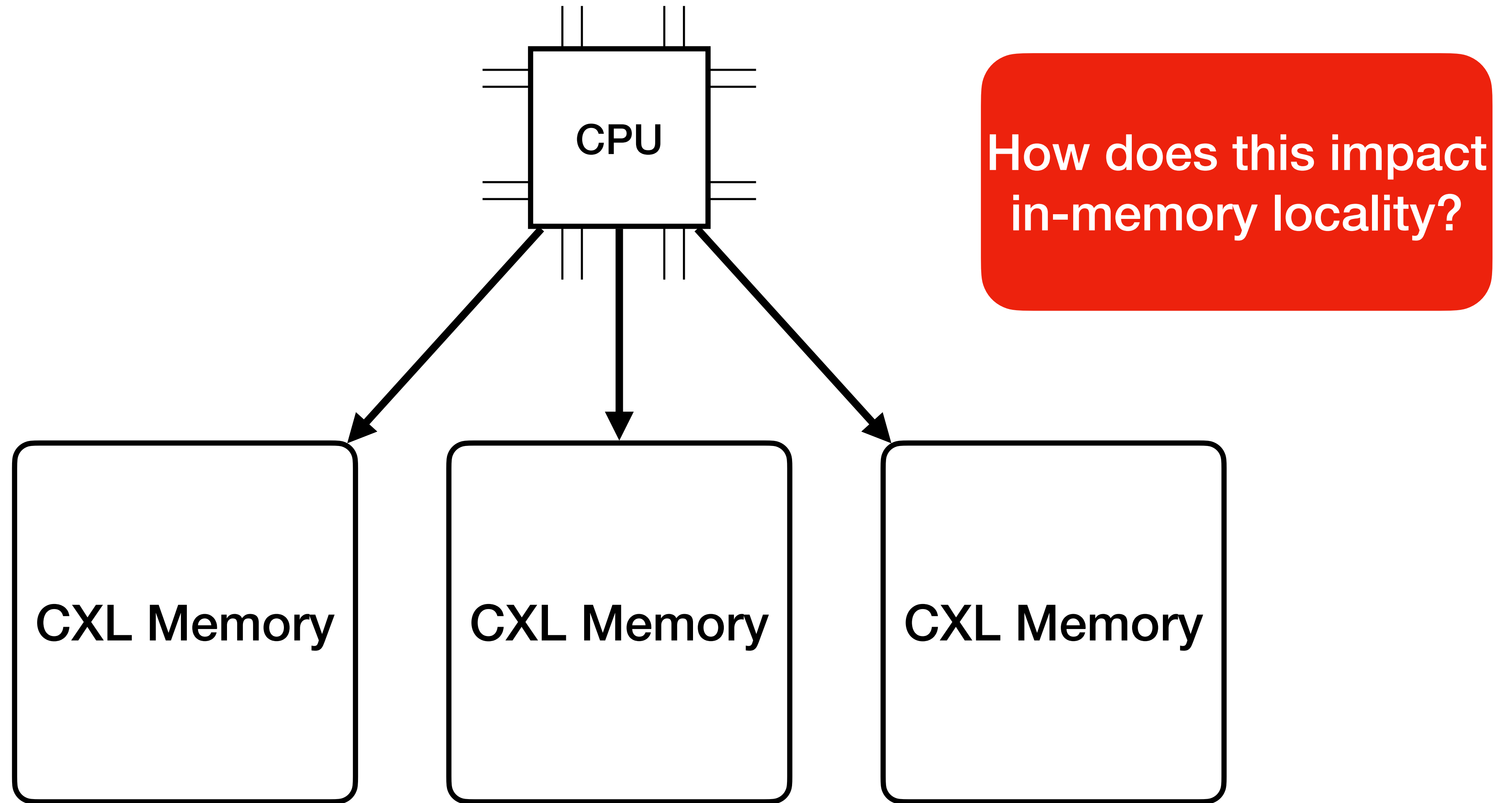
CXL Memory

CXL Memory

CXL Memory

Li, et al. *Pond: CXL-based memory pooling systems for cloud platforms*

# ✔️ Scalable Capacity

- Cheaper than buying new devices

- Doesn't require inter-device coordination

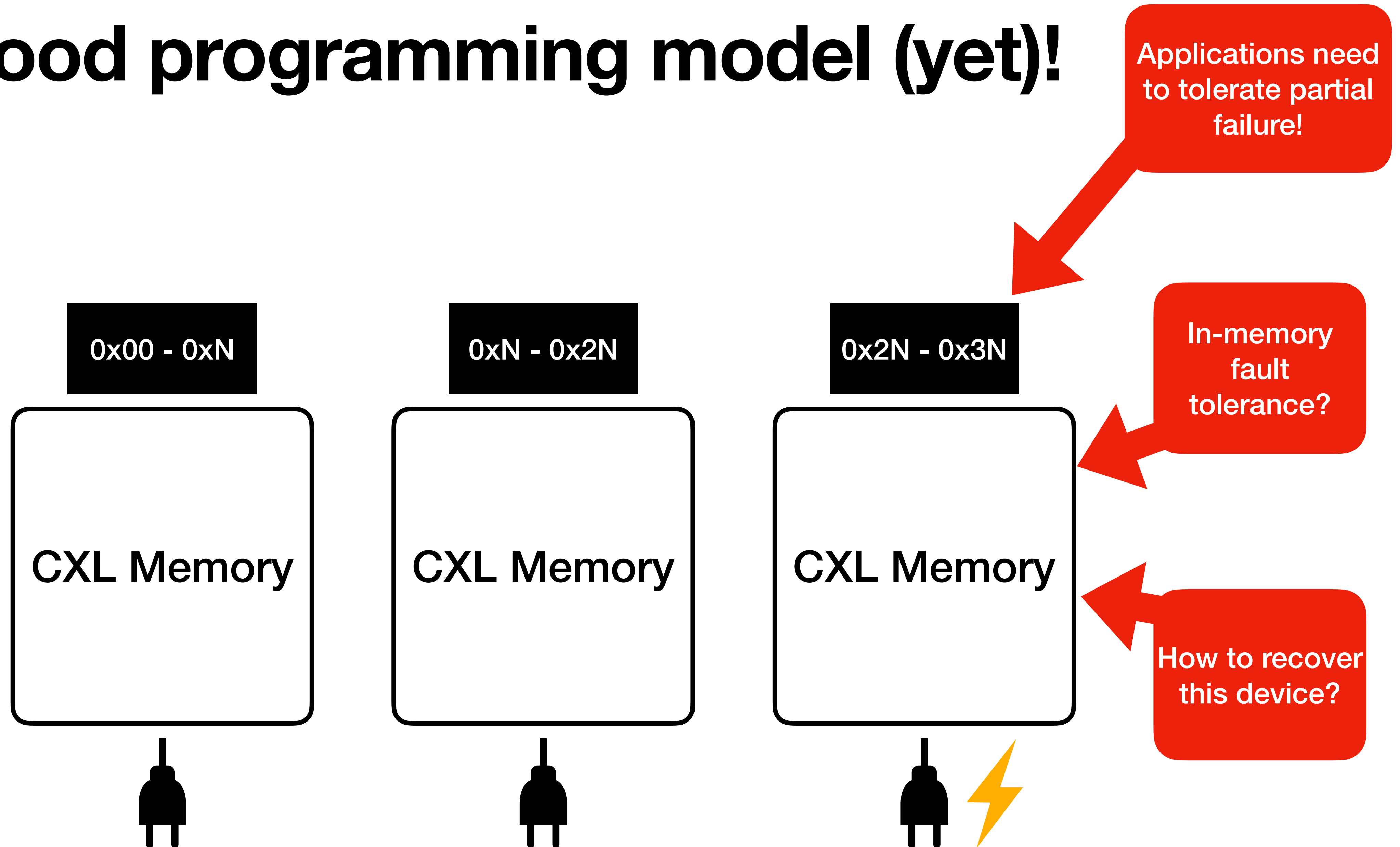- Extra hardware for computer architects to mess around with :-)

**CPU**

**CXL Memory**  **CXL Memory**

✅ **Increased Potential Bandwidth**

CPU

How does this impact in-memory locality?

CXL Memory    CXL Memory    CXL Memory

# ❌ No good programming model (yet)!

**Applications need to tolerate partial failure!**

**Is this shared memory consistency or a distributed system?**

| 0x00 - 0xN | 0xN - 0x2N | 0x2N - 0x3N |
|---|---|---|
| CXL Memory | CXL Memory | CXL Memory |

**In-memory fault tolerance?**

**How to recover this device?**