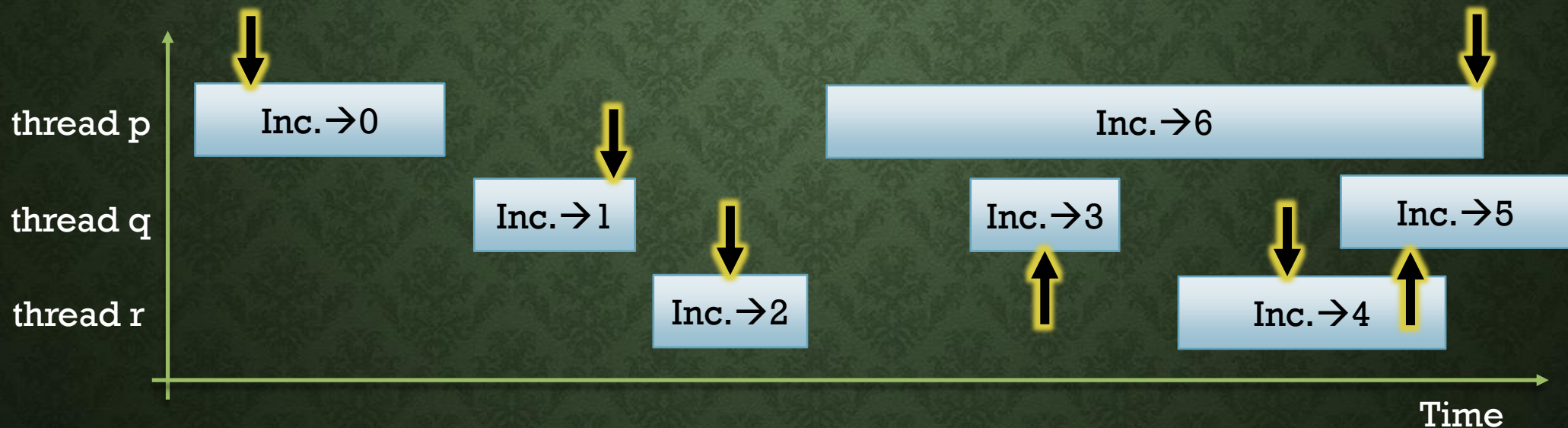# MULTICORE PROGRAMMING

Implementing a counter: what could go wrong?

**Lecture 2**

Trevor Brown

# RECALL: WHAT IT MEANS FOR AN **EXECUTION** TO BE **LINEARIZABLE**

- Must be <u>possible</u> to choose **linearization points** <u>during</u> each operation such that all operations return **the same values** that they would if they were executed **instantly** at their linearization points
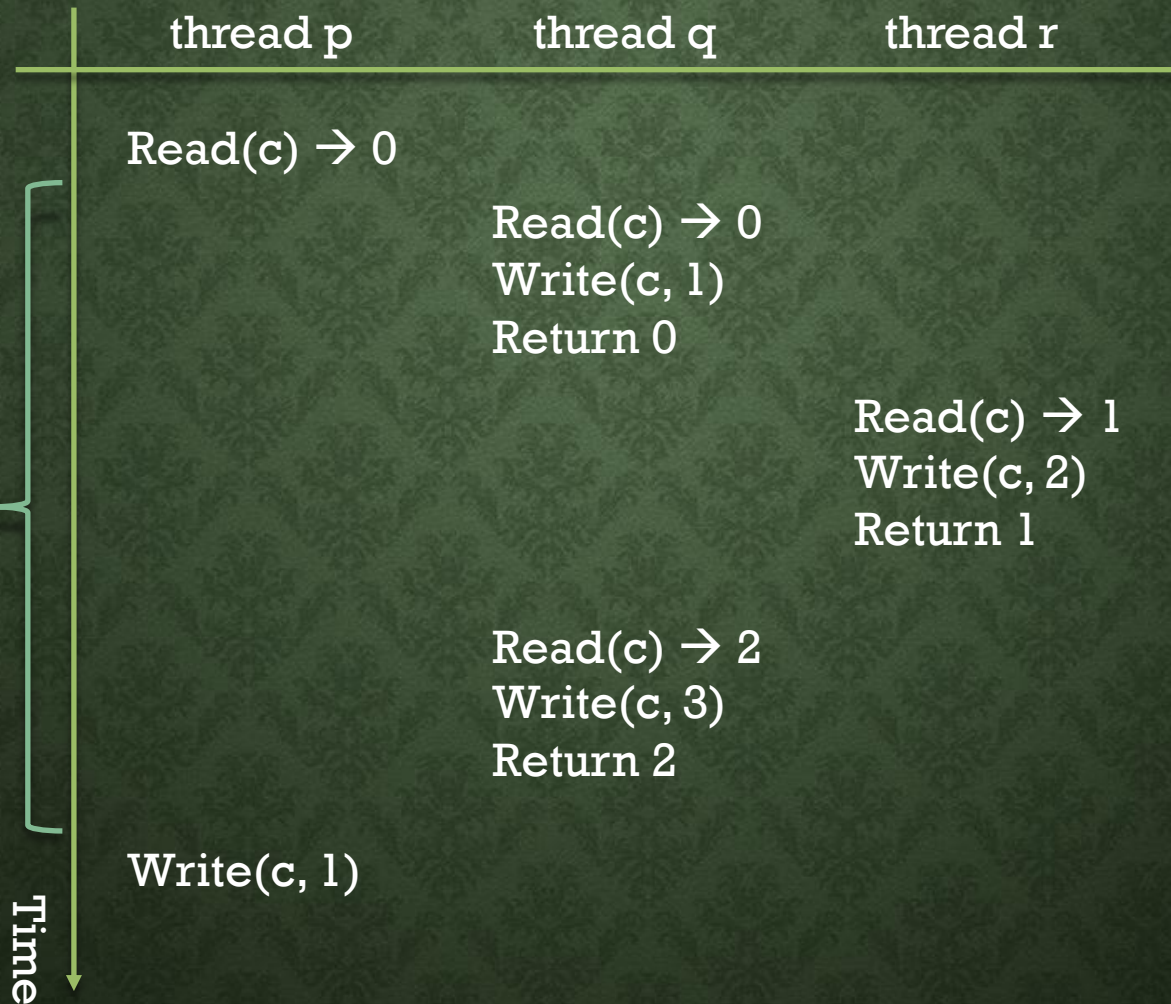


thread p — Inc.→0, Inc.→6

thread q — Inc.→1, Inc.→3, Inc.→5

thread r — Inc.→2, Inc.→4

Time

- <u>Non-linearizable</u> if it is **not possible** to pick such linearization points

# RECALL: WHAT IT MEANS FOR AN **OBJECT** TO BE **LINEARIZABLE**

- For **every possible execution E** of the object
  - E must be a linearizable **execution**

- An **object** is <u>non-linearizable</u>
  - if **any** execution of the object is non-linearizable
  - (i.e., the object can possibly behave "badly")

# RECALL: THIS EXECUTION OF THE
## NAÏVE COUNTER

| thread p | thread q | thread r |
|---|---|---|
| Read(c) → 0 | | |
| | Read(c) → 0<br>Write(c, 1)<br>Return 0 | |
| | | Read(c) → 1<br>Write(c, 2)<br>Return 1 |
| | Read(c) → 2<br>Write(c, 3)<br>Return 2 | |
| Write(c, 1) | | |

Any changes between this Read and Write are **overwritten**!

Time

# IMPLEMENTING A <u>LINEARIZABLE</u> COUNTER

- Intuition: **increment** must **atomically** Read and Write ("at the same time")

  - Otherwise a thread can always Write a "stale" value (overwriting "fresh" values)

- Need stronger tools!

  - How about using a **lock?**

# A LOCK / MUTEX

- **Guards** an object: allows only **one thread at a time** to access it

- Operations
  - Acquire / lock
    - Blocks until the calling thread has acquired the lock (then returns)
    - (Thread is then **allowed** to access the object)
  - Release / unlock
    - Releases the lock so other threads can acquire it
    - (Thread is **no longer** allowed to access the object)

```
Java: synchronized
C: pthread_spin_lock
C++: std::mutex
```
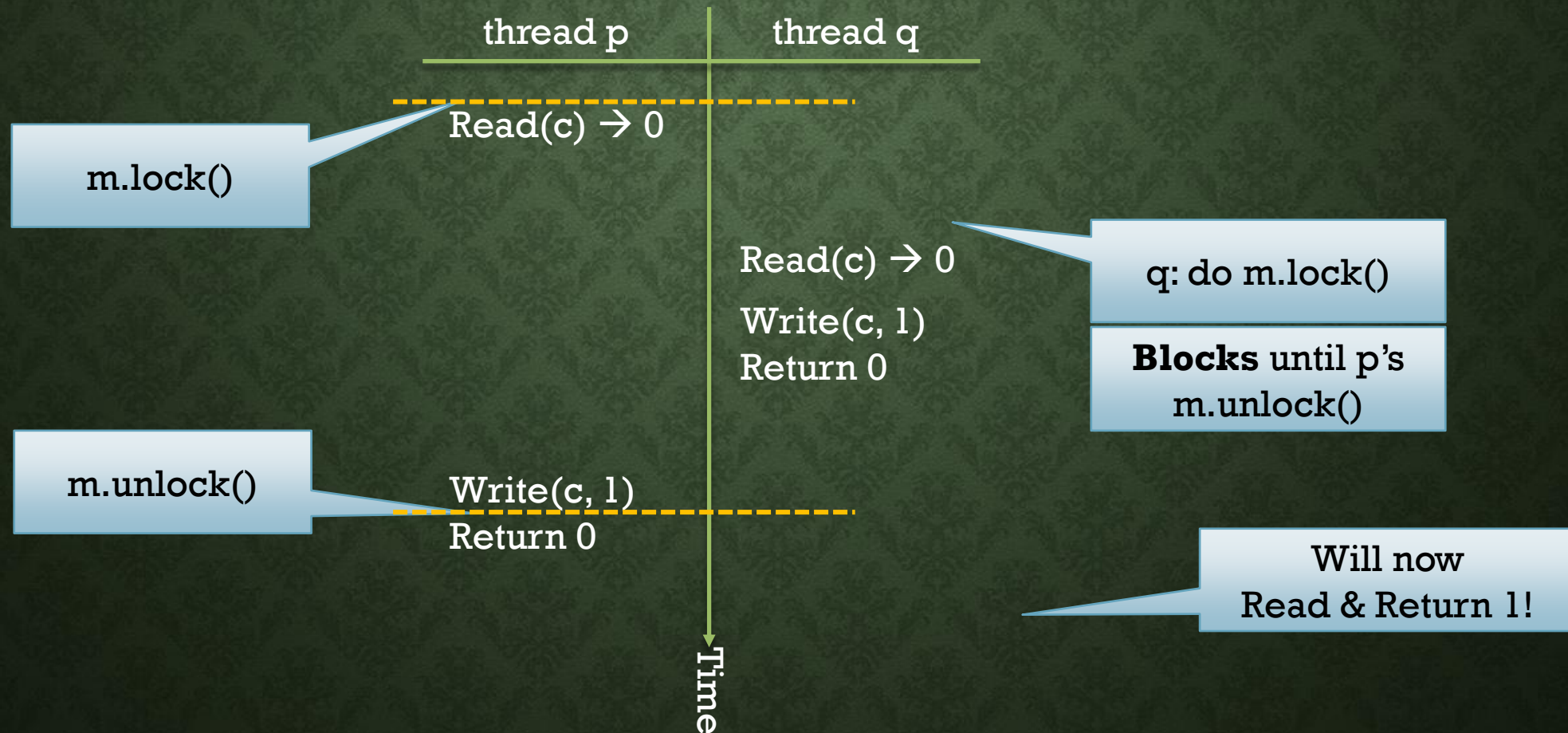
```cpp
36  class counter2 {
37  private:
38      std::mutex m;
39      int v;
40  public:
41      counter2() { v = 0; }
42      int increment(int threadID) {
43          m.lock();
44          auto ret = v++;
45          m.unlock();
46          return ret;
47      }
48      int get() {
49          m.lock();
50          auto ret = v;
51          m.unlock();
52          return ret;
53      }
54  };
```

This is really a **read** followed by a **write**…

How does this help?

Both **read** and **write** are done while holding the lock

# CAN THIS PROBLEM HAPPEN NOW?

thread p | thread q

Read(c) → 0

m.lock()

Read(c) → 0

Write(c, 1)
Return 0

q: do m.lock()

**Blocks** until p's
m.unlock()

m.unlock()

Write(c, 1)
Return 0

Will now
Read & Return 1!

Time

```
36  class counter2 {
37  private:
38      std::mutex m;
39      int v;
40  public:
41      counter2() { v = 0; }
42      int increment(int threadID) {
43          m.lock();
44          auto ret = v++;
45          m.unlock();
46          return ret;
47      }
48      int get() {
49          m.lock();
50          auto ret = v;
51          m.unlock();
52          return ret;
53      }
54  };
```

**Intuition behind why these linearization points work:**
to all threads, operations appear to happen instantly at these lines

**Linearize increment at the WRITE**
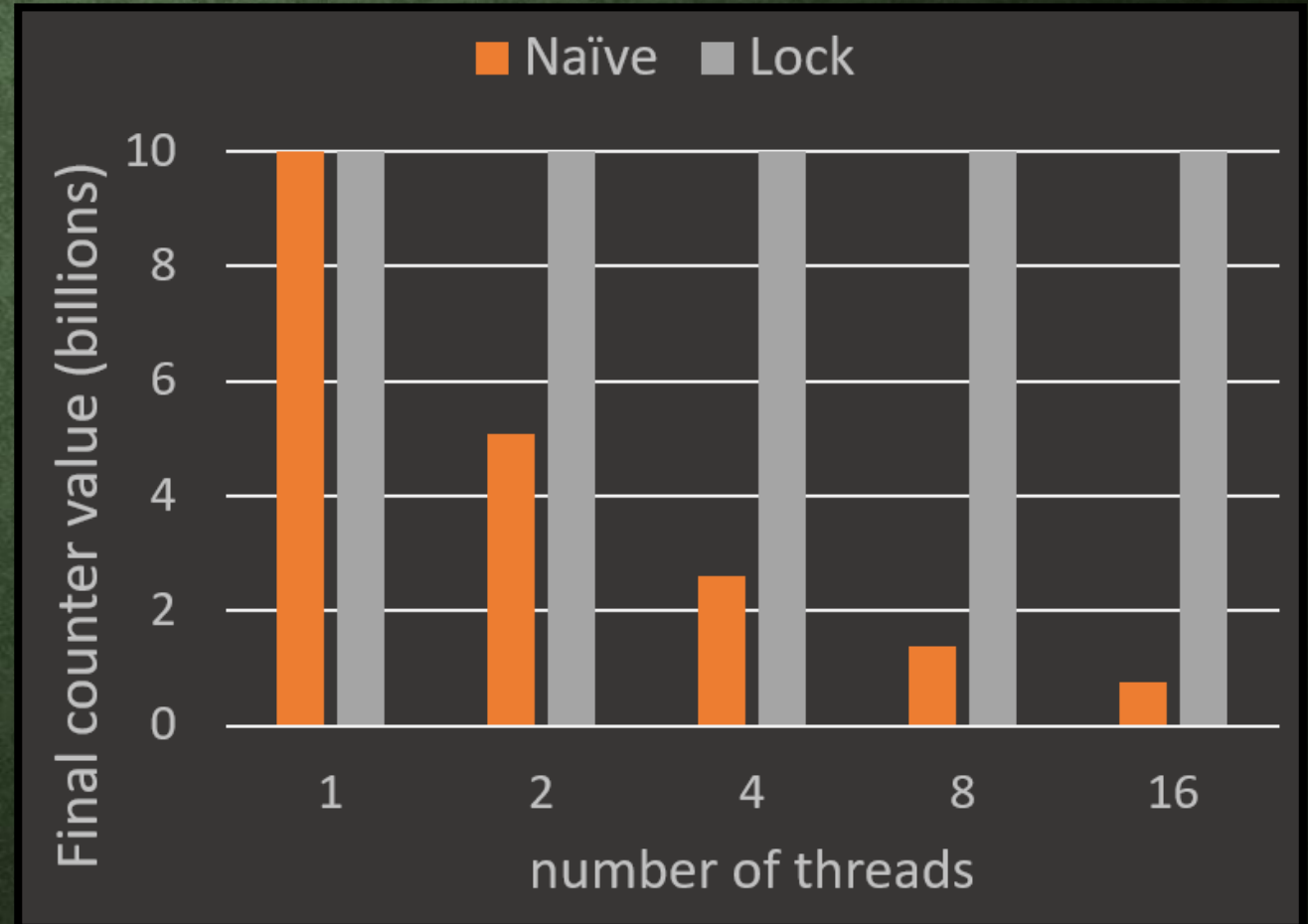(really, *any time* when the lock is held
would work)

(So, anything that happens
while the counter is locked
is **effectively atomic)**

(… **because**, from the perspective of other threads,
the counter cannot be accessed while it is locked)

**Linearize get at the READ of v**
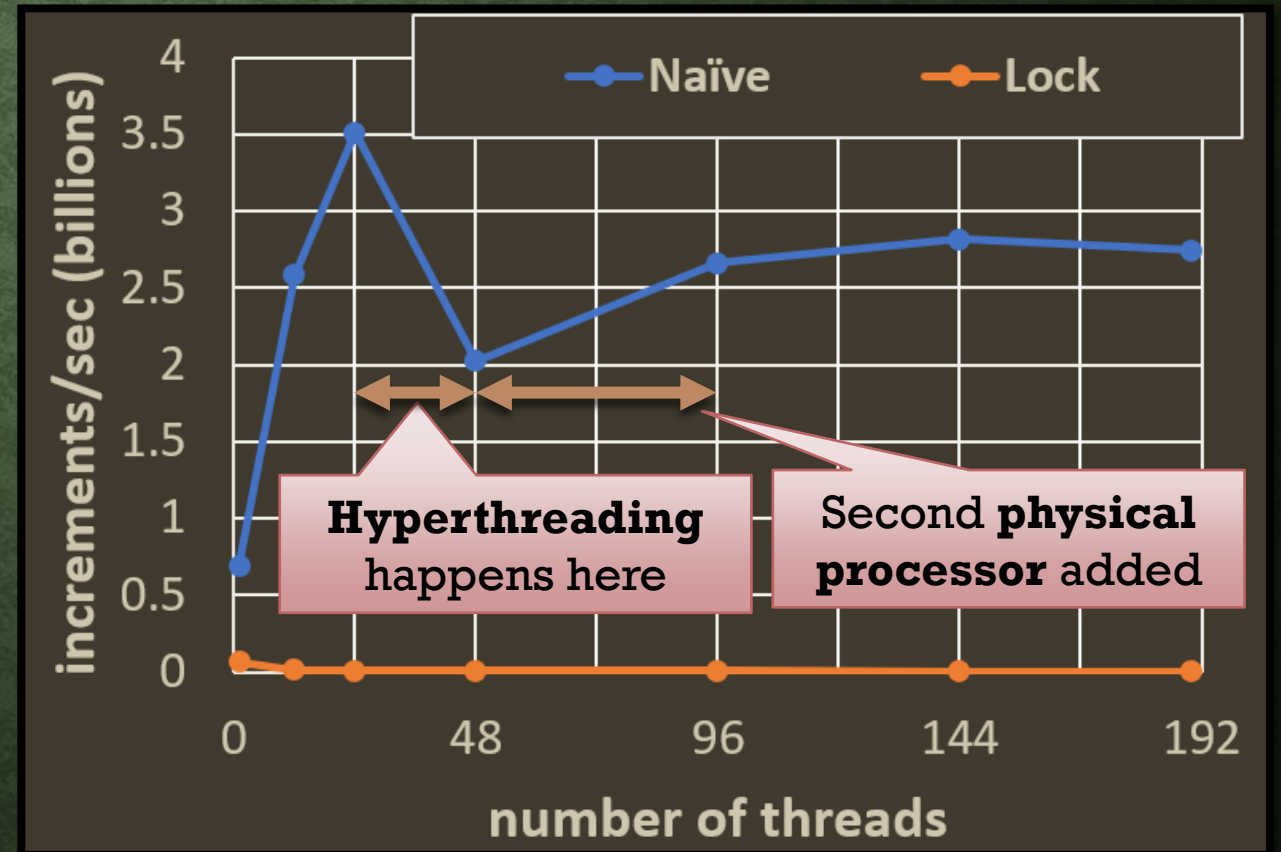
# OUTPUT AFTER ADDING A LOCK

- Same as the previous **"accuracy experiment"**
  - Comparing final counter value of **naïve** and **lock-based**

- Output is now **correct**!

- (Of course, this experiment is **not** a correctness proof)

- What about **performance**?

# PERFORMANCE COMPARISON

- Simple **timed** experiment

- Each data point = average of 5 trials

- In each trial, for 3 seconds,
  - threads repeatedly perform Increment,
  - and we measure increments/second

- What is the overhead of locking?
  - 10x slower with 1 thread
  - **450x slower** with 190 threads

- Is there a better tool than locking?

Machine with 4 physical processors
(each with 24 cores + hyperthreading)



Hyperthreading happens here

Second **physical processor** added

# FETCH AND ADD (FAA)

- Instruction implemented modern Intel and AMD systems:

  - **`lock xadd`**

- FAA(addr, val) does the following **atomically** (all at once)
  - old = Read(addr)
  - temp = old + val
  - Write(addr, temp)
  - Return old

# EASY FAA-BASED COUNTER
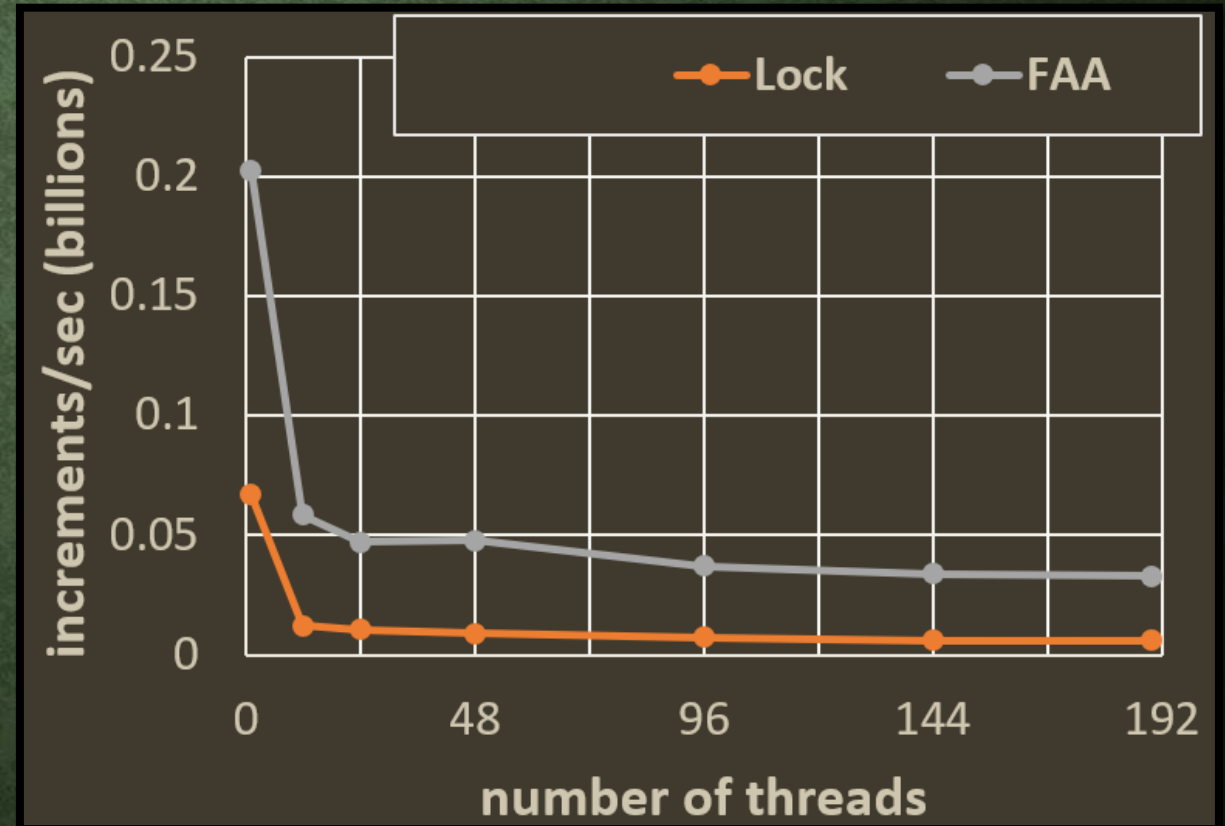
```cpp
57  class counter3 {
58  private:
59      atomic<int> v;
60  public:
61      counter3() : v(atomic<int>(0)) {}
62      int increment(int threadID) {
63          return v++;
64      }
65      int get() {
66          return v;
67      }
68  };
```

```asm
2          mov      DWORD PTR [rsp-4], 0
3          mov      eax, 1
4          lock xadd        DWORD PTR [rsp-4], eax
5          ret
```

Because **v** is **atomic<int>**, this is really a **FAA**!

# HOW DOES THIS PERFORM IN PRACTICE?

- Same timed experiment

- Excluding Naïve from the graph (to zoom in on Lock and FAA)

- Compared to Lock
  - FAA is up to **5.4x faster**

- Compared to Naïve (incorrect)
  - FAA is up to **83x slower**
  - (much better than Lock's **450x**)

# PROBLEM: TOO MUCH <u>CONTENTION</u>

- Accessing a single counter creates a **contention bottleneck**

- What if we **shard (partition)** the counter into multiple **sub counters**

  - Increment: pick one sub counter and increment it

  - What about Get?

    - Counter value is *distributed* over the sub counters

    - Trade-off

      - Single counter → slow Increment, fast Get

      - Sharded counter → fast Increment, slower/more complex Get?

    - We are going to **ignore** these complications and **only think about increment...**

- Each thread uses its own sub counter

- **No data sharing,** should scale perfectly

```cpp
73  class counter4 {
74  private:
75      atomic<int> data[MAX_THREADS];
76  public:
77      counter4() {
78          for (int threadID=0; threadID<MAX_THREADS; ++threadID)
79              new (&data[threadID]) atomic<int>(0);
80      }
81      int increment(int threadID) {
82          return data[threadID]++; // atomic
83      }
84      int get() {
85          int sum = 0;
86          for (int threadID=0; threadID<MAX_THREADS; ++threadID)
87              sum += data[threadID];
88          return sum;
89      }
90  };
```
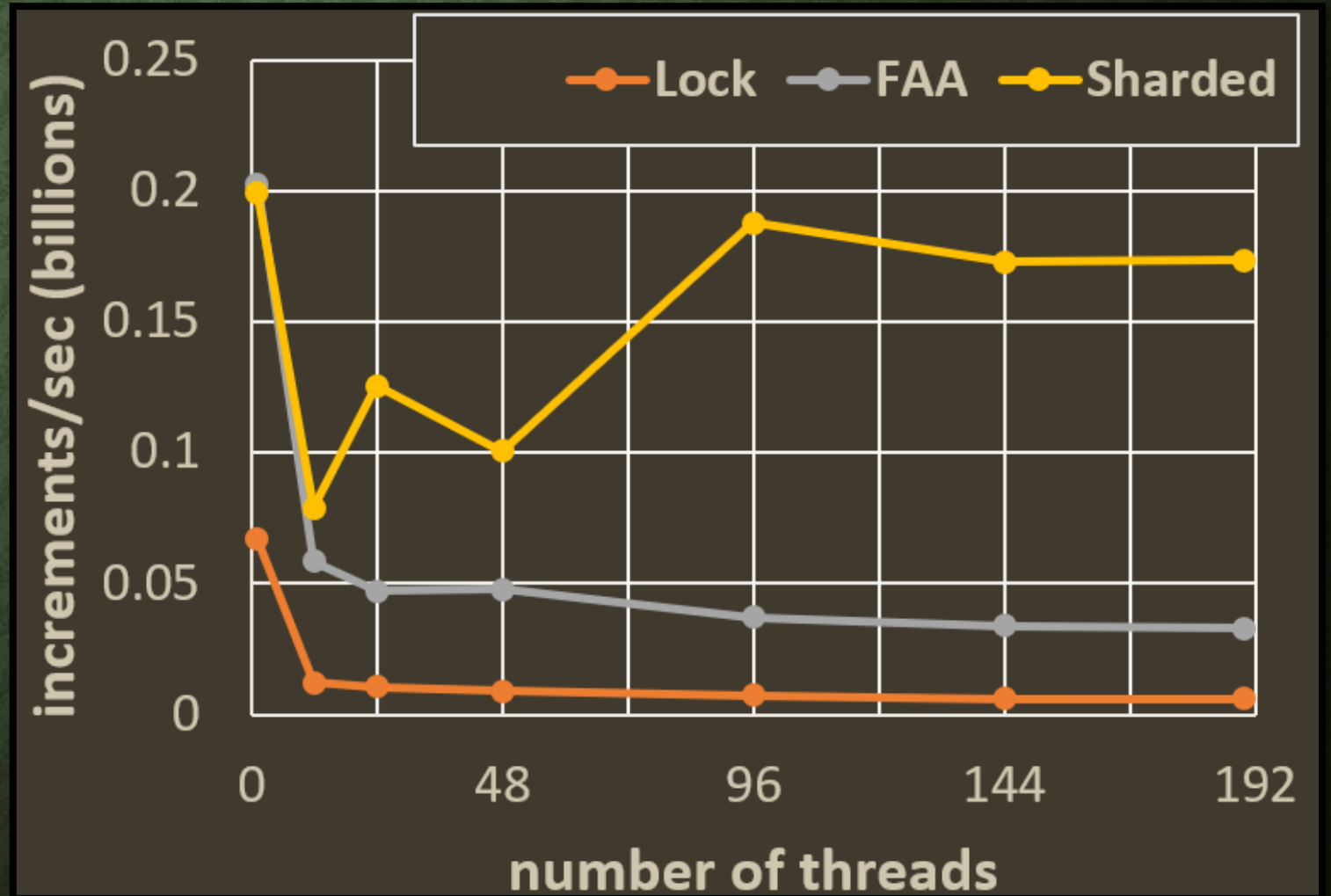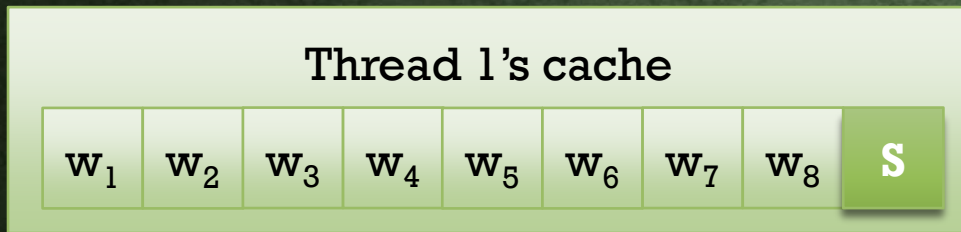
# HOW DOES THIS PERFORM?

- Same timed experiment
- Why is the scaling so poor?
  - **No shared data**, right?
- Answer: **cache coherence**

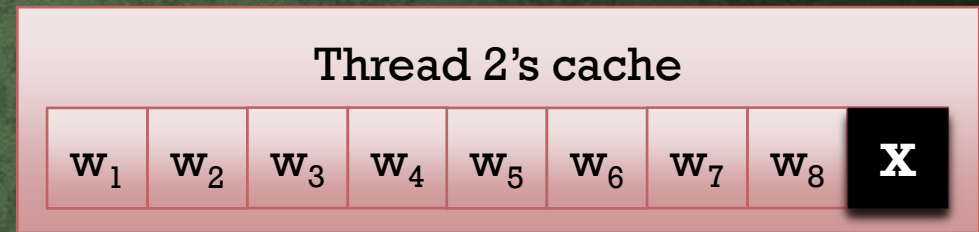# HOW CACHE COHERENCE WORKS

Thread 1's cache

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$ $w_7$ $w_8$ **S**

Thread 2's cache

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$ $w_7$ $w_8$ **X**

Thread 1 reads $w_2$

Thread 2 reads $w_7$

Cache line **invalidated** and **evicted**!

Thread 2 **writes** $w_7$

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$ $w_7$ $w_8$

64 byte (8 word) cache line

# MEMORY LAYOUT OF SUB COUNTERS

| c[0] | c[1] | c[2] | c[3] | c[4] | c[5] | c[6] | c[7] | c[8] | c[9] | c[10] | c[11] | c[12] | c[13] | c[14] | c[15] | ... |

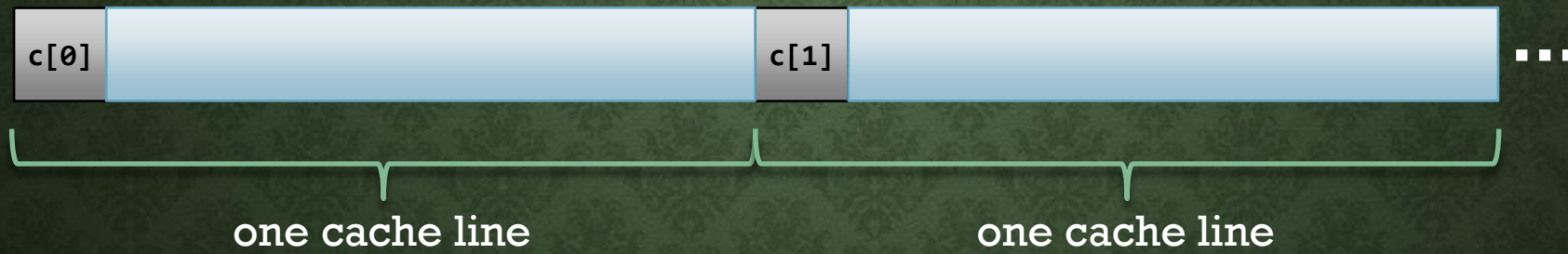one cache line

one cache line

Incrementing **one of these** will invalidate **all of them**
(causing huge contention)

This is called
**false sharing**

# SOLUTION: <u>PADDING</u>

- Add **empty space** to each sub counter
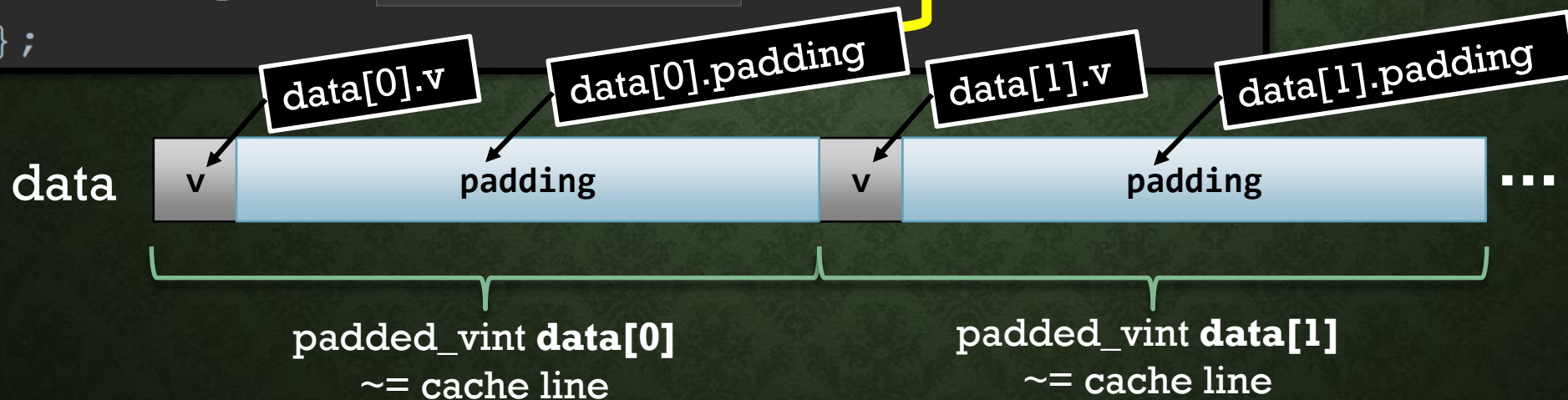    - To make it **cache line sized**



| c[0] | | c[1] | | ... |

one cache line       one cache line

```
93   class counter5 {
94   private:
95       struct padded_vint {
96           atomic<int> v;
97           char padding[64-sizeof(atomic<int>)];
98       };
99       padded_vint data[MAX_THREADS];
100  public:
101      counter5() {...4 lines }
105      int increment(int threadID) {..
108      int get() {...6 lines }
114  };
```
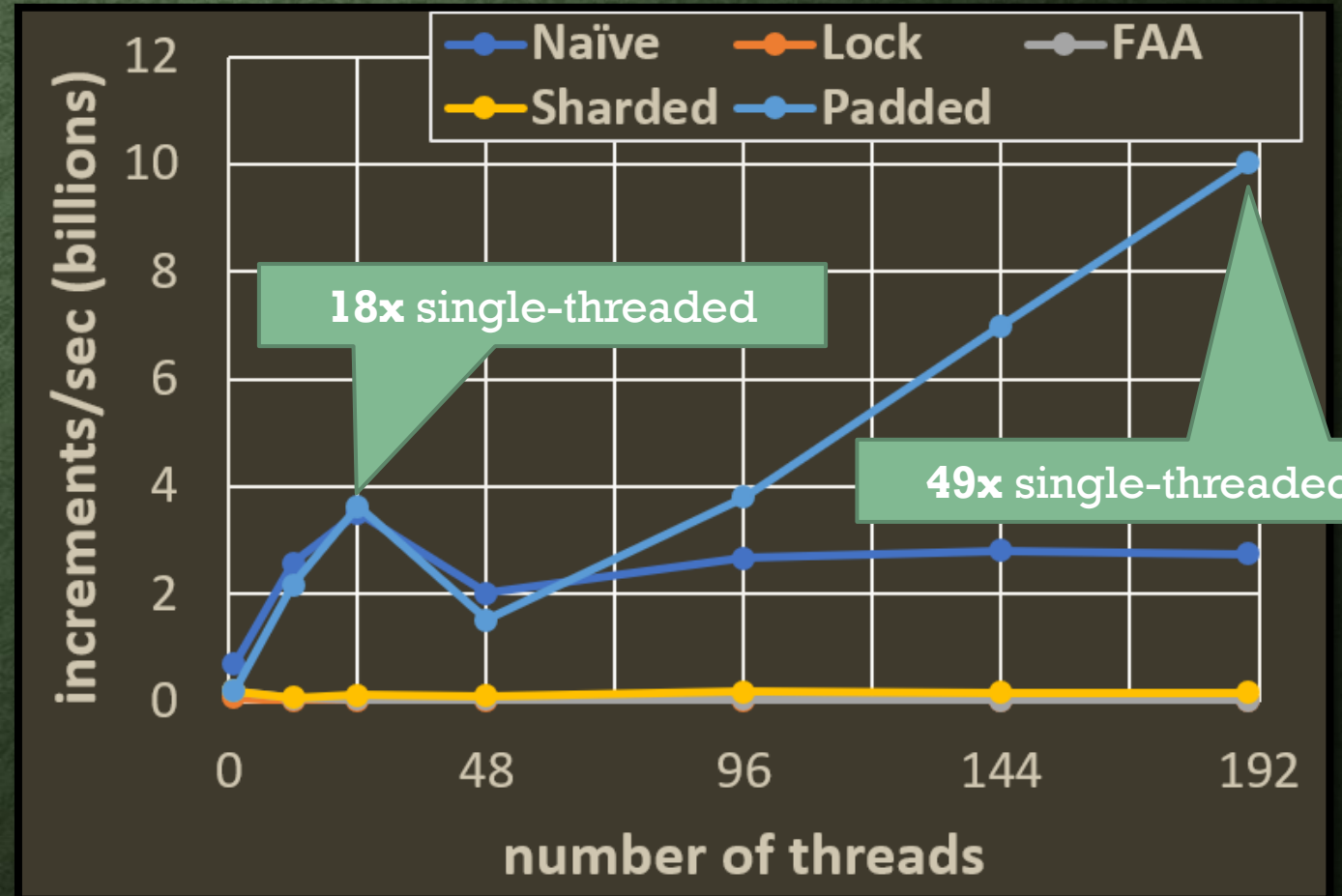
**This is where the magic happens**

**Identical to naïve sharded counter**

data[0].v  data[0].padding  data[1].v  data[1].padding

data  | v | padding | v | padding | ...

padded_vint **data[0]**
~= cache line

padded_vint **data[1]**
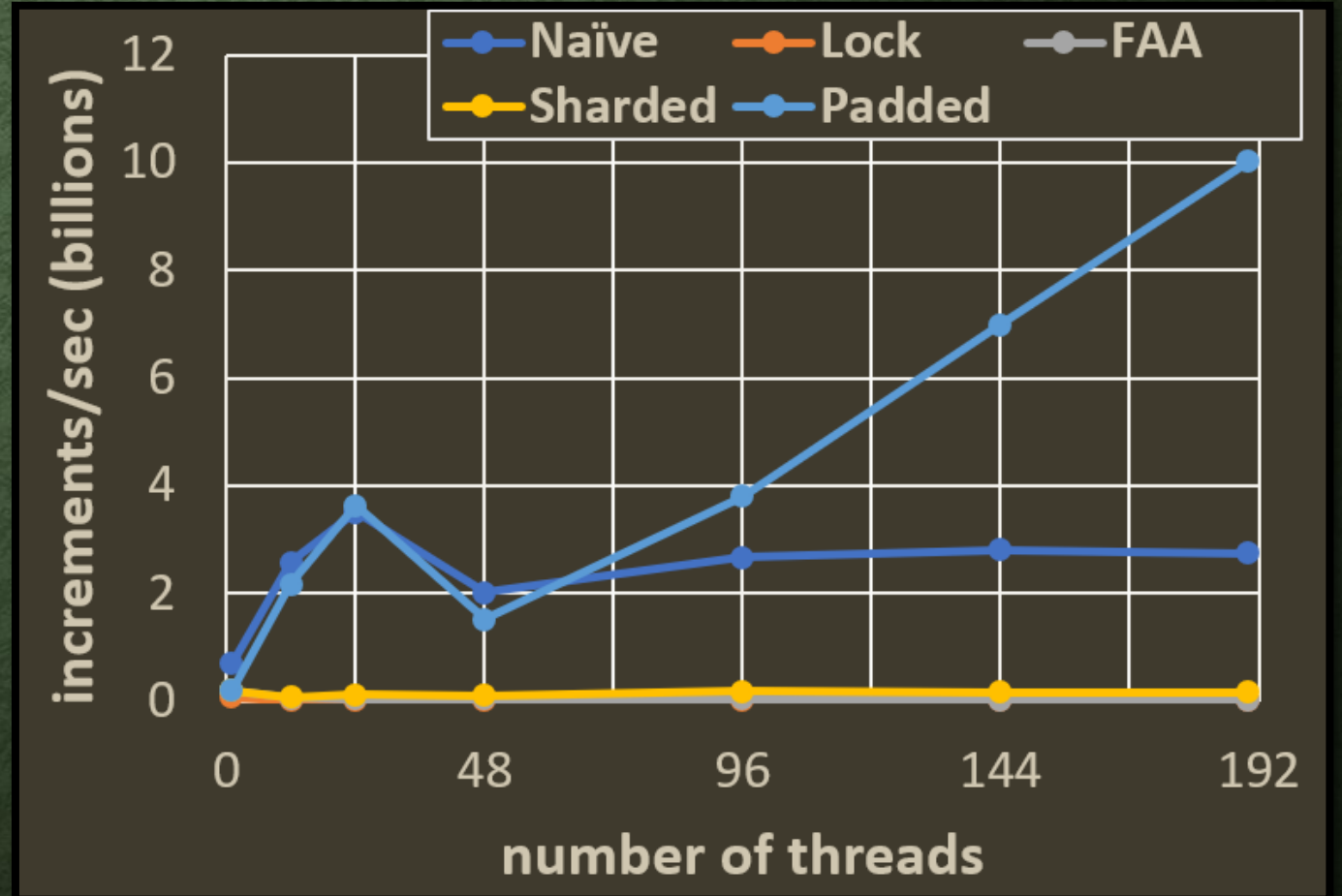~= cache line

# HOW DOES THIS PERFORM?

- Same experiment, but comparing **naïve sharding** with a **padded counter**

- **Pretty good scaling**
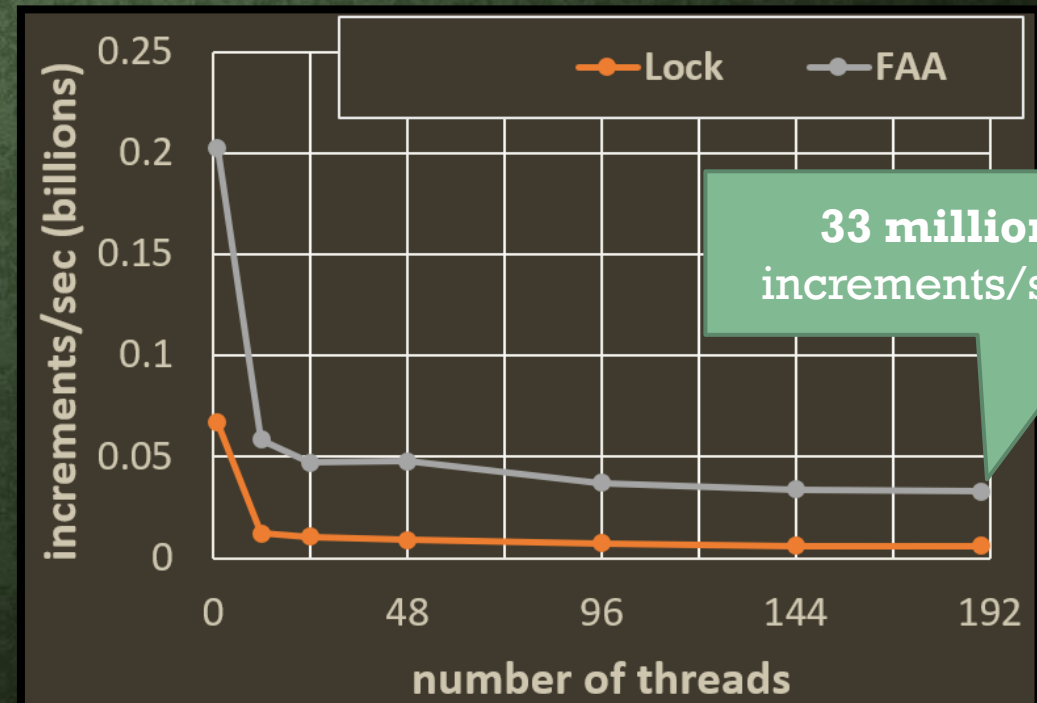  - **18x vs optimal 24x**
  - **49x vs optimal 190x**

# SPEED VS SIMPLICITY

- But… **reading** is hard!
  - Solving this will add complexity

- Simplicity is valuable!
  - Do we **need** a complex solution?
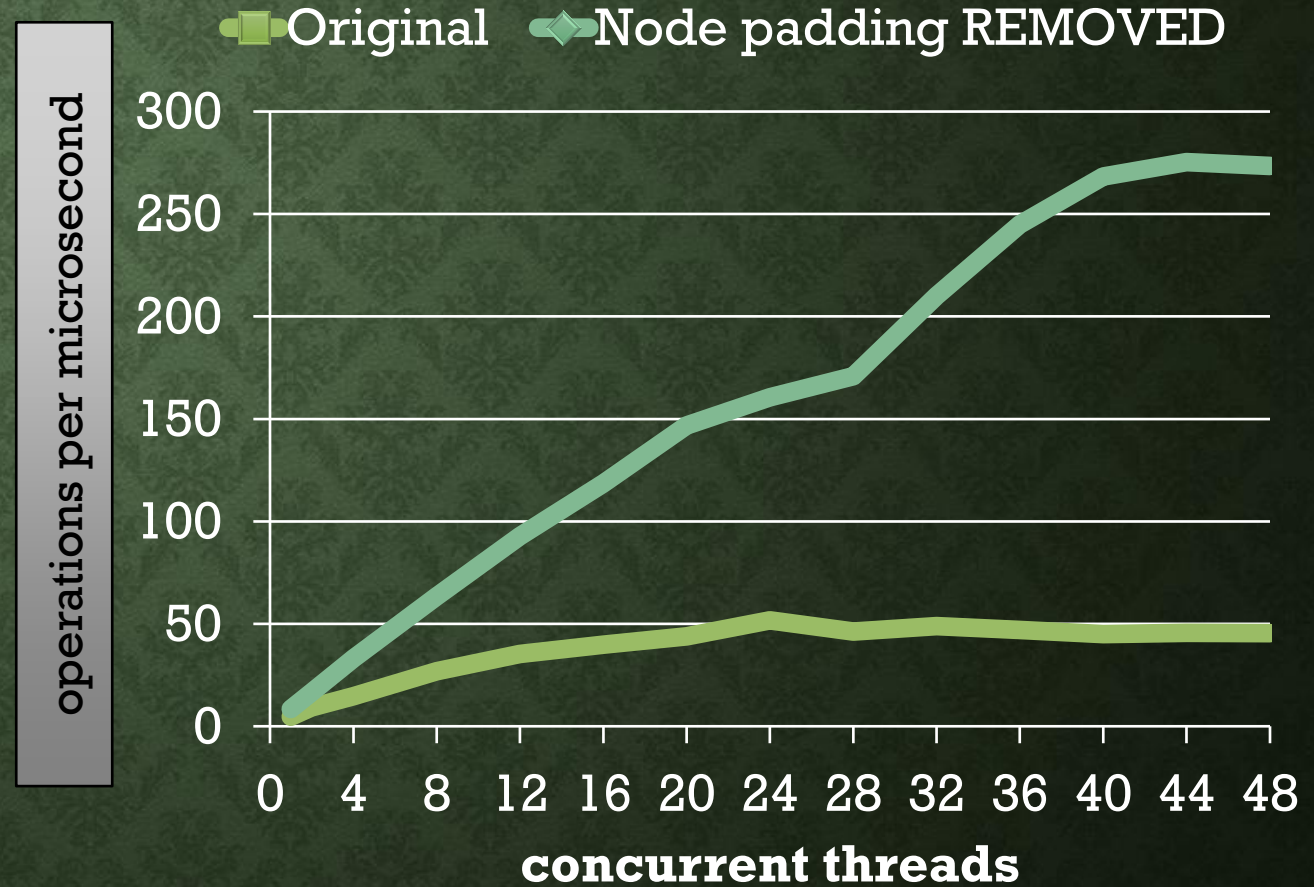  - Sometimes… but not always…

# CONSIDERING USE CASES: A FAA-COUNTER MIGHT BE <u>GOOD ENOUGH</u>

- FAA-based counter does not truly scale

- But, its absolute throughput might be **high <u>enough</u>** for your application

- Real applications do more than just increment a single counter
  - Avoid unnecessary optimization
  - Figure out **if it's a bottleneck** first



**33 million** increments/sec!

# A WORD OF WARNING: PADDING CAN <u>HURT</u>

- Union-find data structure

- Each 8b node was padded to 64b

- **Removing** padding → 5x faster!

- Why?
  - Many nodes, uniformly accessed
    → contention is rare
    → false sharing is rare
    → padding can't help much
  - Padding wastes space
    → 1/8<sup>th</sup> as many nodes in cache!

# WHEN TO PAD?

- When the number of objects being padded is O(# threads) for a small constant

- AND threads frequently **write** to these objects

- Try and see if it helps…

# SUMMARY

- **<u>Cache coherence</u>**, shared and exclusive modes, cache invalidations, contention

- Sharding (partitioning data to reduce contention)

- False sharing and padding (principle: when to pad)

- Locks, fetch-and-add

- Implementing linearizable counters

  - Lock-based counter

  - Fetch-and-add counter

  - Sharded counter

  - Padded sharded counter