# Hardware and its Habits

*Paul E. McKenney*
*IBM Linux Technology Center*
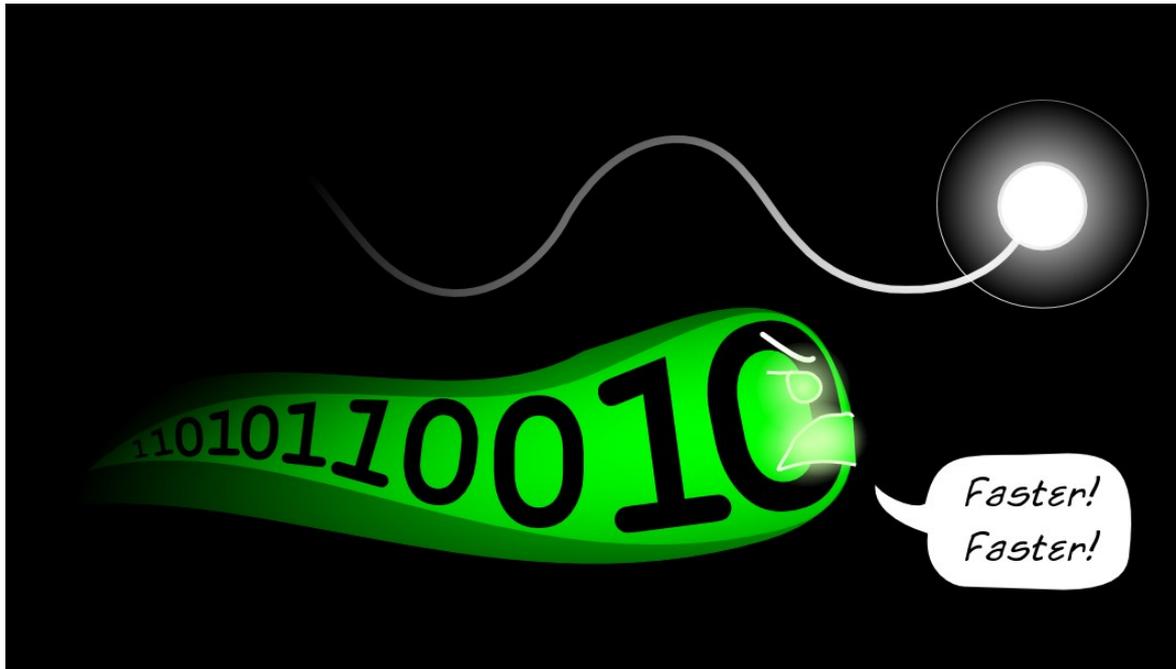
*http://www.rdrop.com/users/paulmck*
*https://www.kernel.org/pub/linux/kernel/people/paulmck/perfbook/perfbook.html*

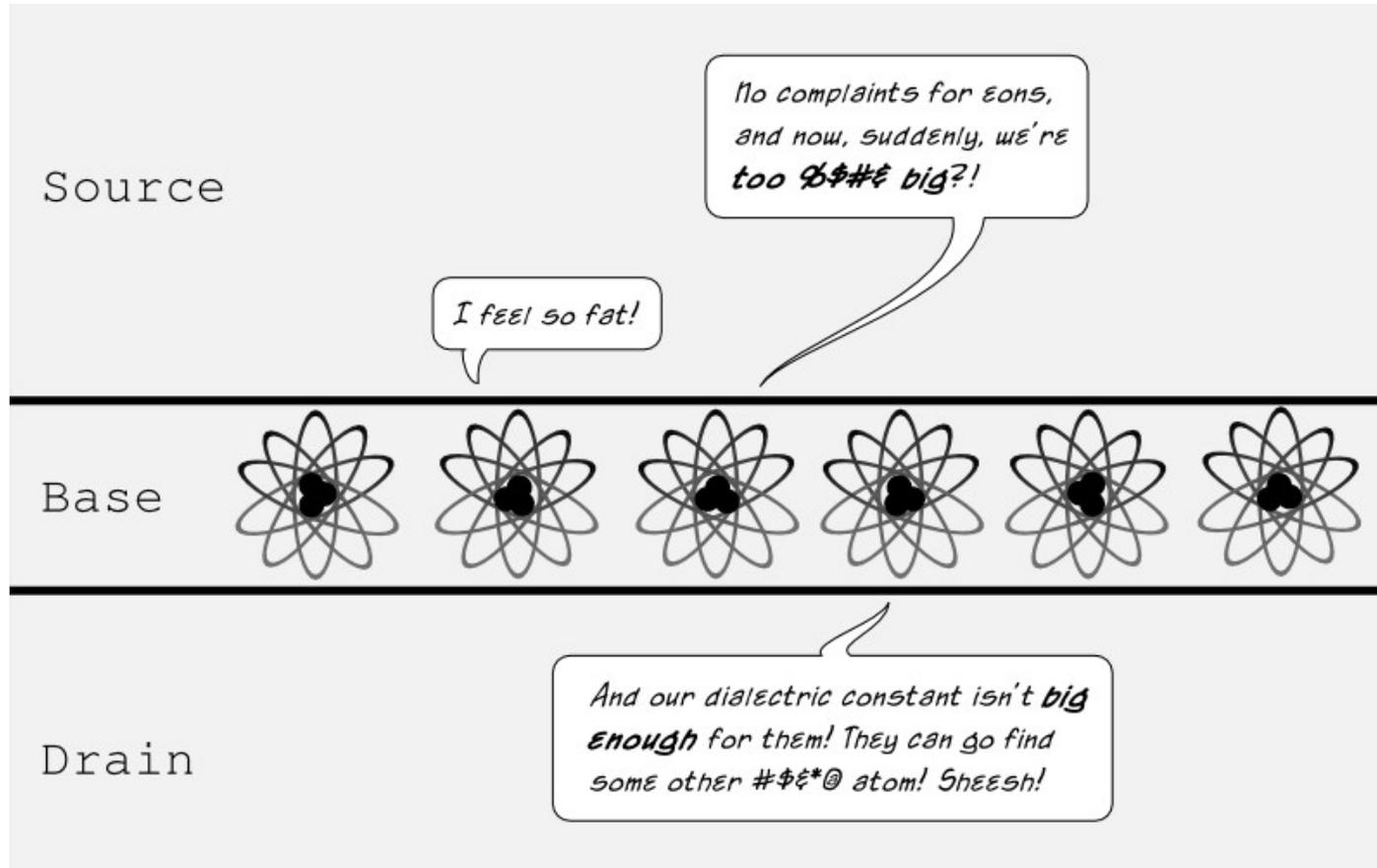# Premature Abstraction Is The Root of All Evil!!!

# Premature Abstraction

## Abstract away the finite speed of light???

# Premature Abstraction

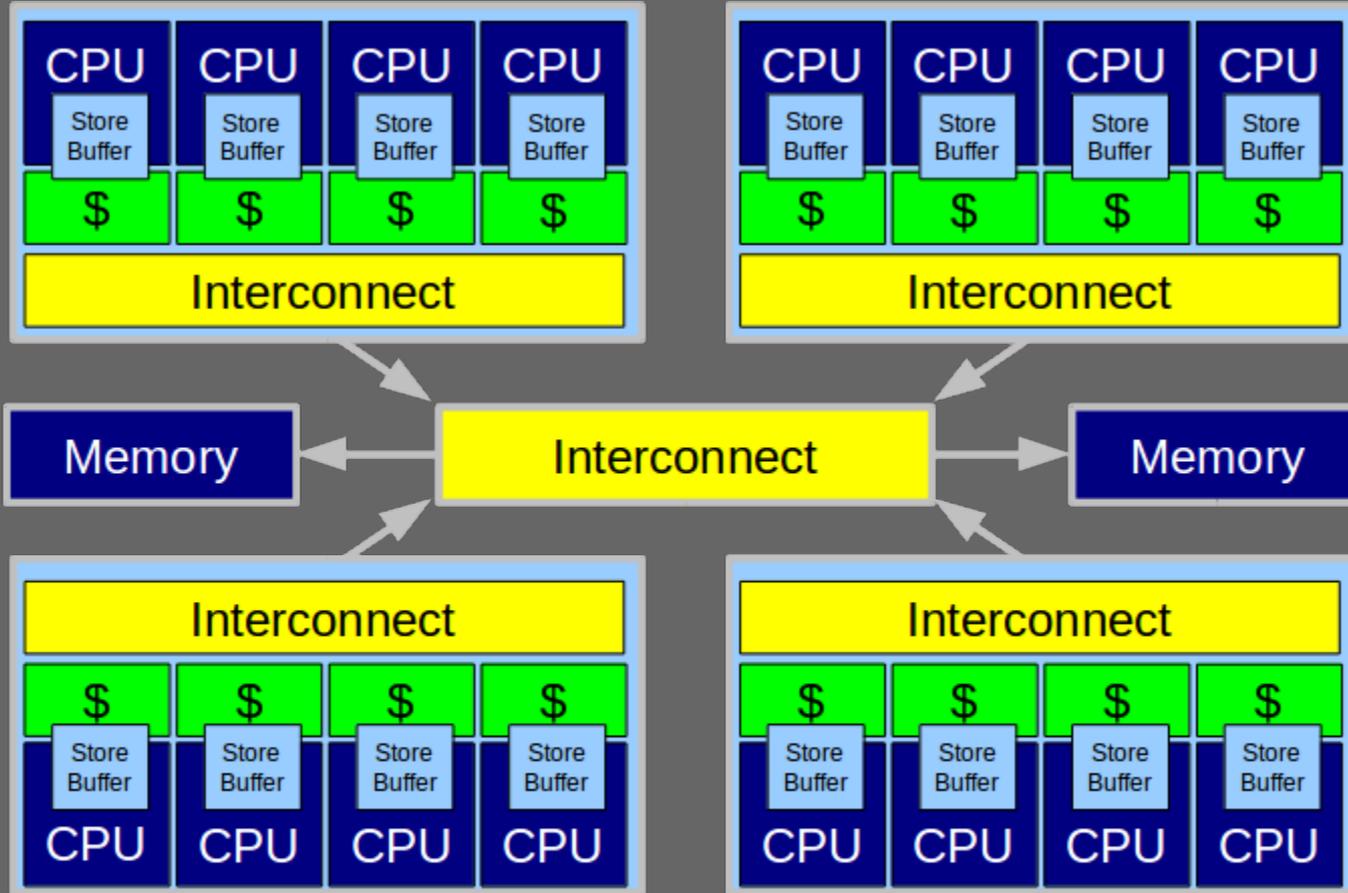## Abstract away the atomic nature of matter???

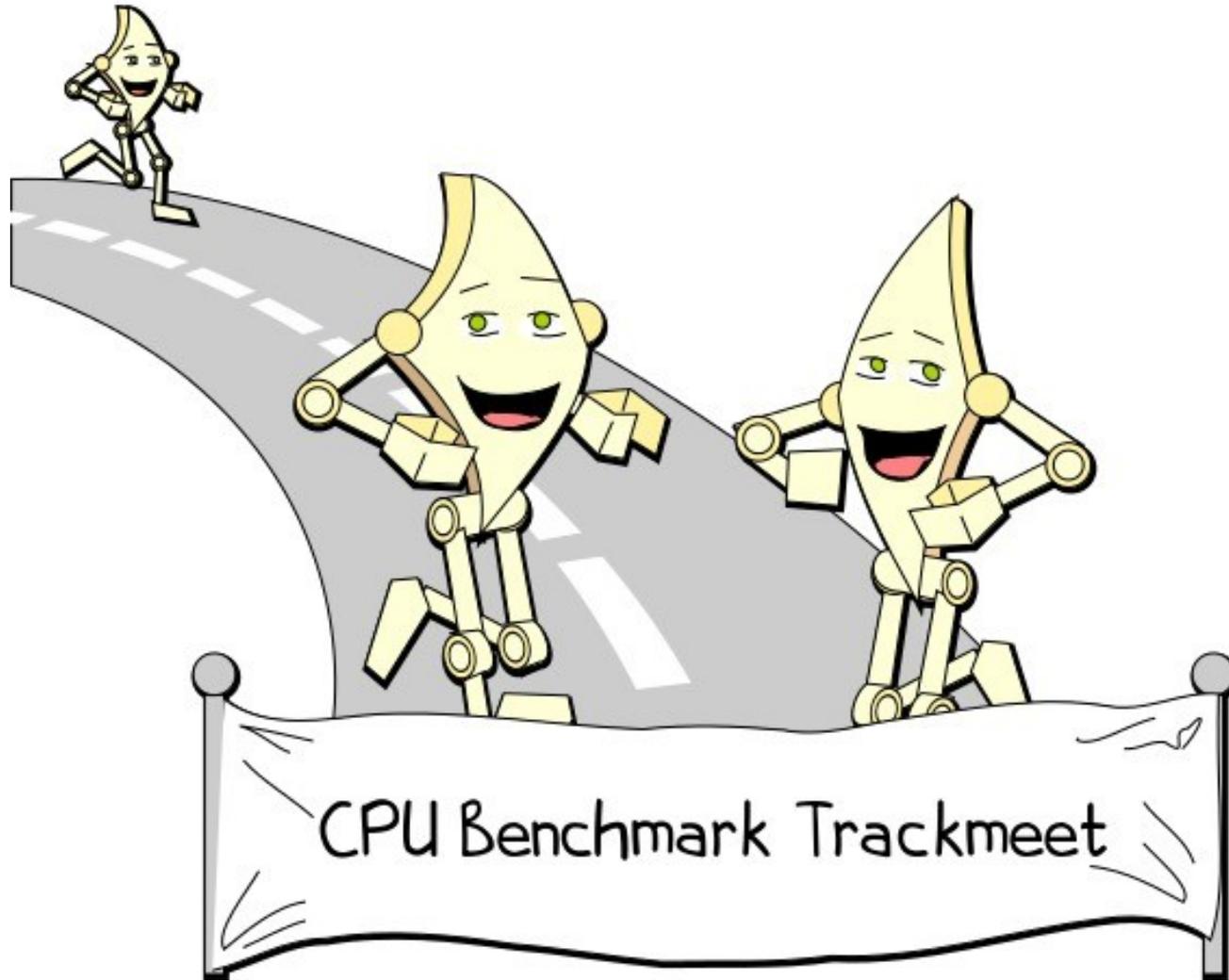# Who Cares About Laws of Physics?



Electrons move at 0.03C to 0.3C in transistors and, so lots of waiting. 3D???

# CPUs at Their Best
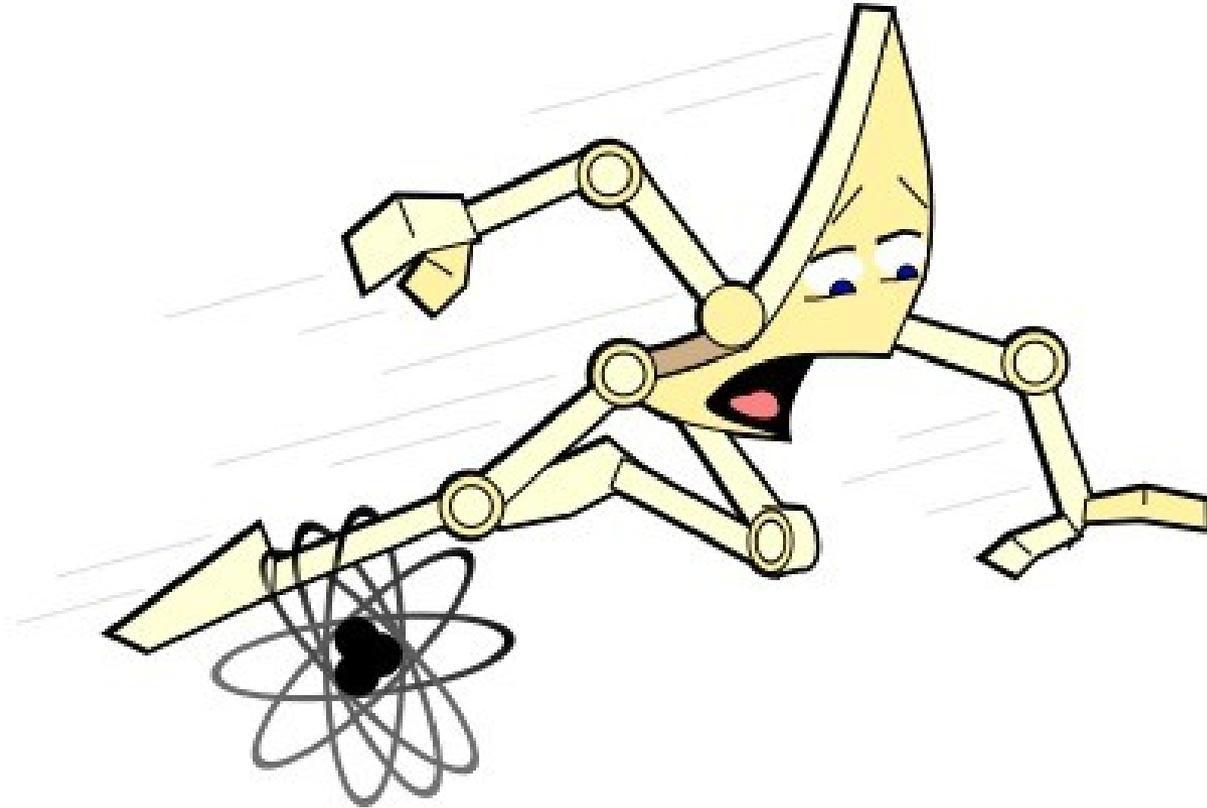


CPU Benchmark Trackmeet

# Mispredicted Branch

# Memory Reference

# RMW Atomic Instruction

# Memory Barrier (AKA Fence)

# Memory Ordering Not Created Equal



Exact relation depends strongly on hardware architecture and on compiler

# Cache Miss

# I/O Operation

# Costs of Operations

16-CPU 2.8GHz Intel X5550 (Nehalem) System

| Operation | Cost (ns) | Ratio |
|---|---|---|
| Clock period | 0.4 | 1 |
| "Best-case" CAS | 12.2 | 33.8 |
| Best-case lock | 25.6 | 71.2 |
| Single cache miss | 12.9 | 35.8 |
| CAS cache miss | 7.0 | 19.4 |
| Single cache miss (off-core) | 31.2 | 86.6 |
| CAS cache miss (off-core) | 31.2 | 86.5 |
| Single cache miss (off-socket) | 92.4 | 256.7 |
| CAS cache miss (off-socket) | 95.9 | 266.4 |

# Loophole in Laws of Physics



SOL RT @ 2GHz

7.5 centimeters

CPU — Store Buffer — $

Interconnect

Memory

Interconnect

Read-only data remains replicated in all caches

# In Real Estate:

Location, location, location!!!

# In Parallel Programming:

## Locality, locality, locality!!!

# But Can't The Hardware Help???

# Hardware Locality Optimizations

# Big Cachelines: The Good

Almost no additional  latency

```
r2 =  b ;
r3 =  c ;
r4 =  d ;
```

Long latency

r1 =  a ;

CPU 0

a b c d

a b c d

CPU 1

Time

# Big Cachelines: The Bad and Ugly

CPU 0

`a = 42;` Long latency `a = 256;` Long latency

`a b c d`

CPU 1

`a b c d` `a b c d`

`d = 1729;` Long latency

Time

False sharing!!!

# Big Cachelines: Alignment Directives

Immediate effect!!!

Long latency

`a = 42;`

`a = 256;`

CPU 0

a b

a b | c d

c d

CPU 1

`d = 1729;`

Immediate effect!!!

Time

# Prefetching Good!!!

r2 = b ;
r3 = c ;
r4 = d ;
r5 = e ;
r6 = f;

Almost no additional  latency

Long latency

r1 = a ;

CPU 0

a b c d    e f g h

a b c d    e f g h

CPU 1

Time

# Prefetching: Bad and Ugly!

# Store Buffers: The Good

Almost no additional latency

r2 = e ;
r3 = f;
r4 = g ;

Store to "a"
can complete

a = 42 ;

CPU 0

e f g h

a b c d

Hidden
latency

a b c d

CPU 1

Time

# Store Buffers: The Bad and Ugly

Almost no additional latency

```
r2 = e ;
r3 = f ;
r4 = g ;
a = 42 ;    r5 = a ; (42)
```

Store of 42 to "a" can complete

**CPU 0**

e f g h

a b c d

Hidden latency

"Store to 'a' happened before load from 'e'."

"No, store to 'a' happened **after** load from 'e'!!!!"

a b c d

**CPU 1**

**Time**

# Store Buffers: The Bad and Ugly

Almost no additional  latency

```
r2 = e;
r3 = f;
r4 = g;
r5 = a; (42)
```

a = 42;

Store of 42 to "a"
can complete

**CPU 0**

e f g h

Hidden
latency

a b c d

"Store to 'a' happened before
load from 'e'."

"No, store to 'a' happened
***after*** load from 'e'!!!!"

a b c d

**CPU 1**

**Time**

So we use ordering directives where needed!!!

# Portable Ordering: Memory Model

C++11 Memory Model: Ordering Portability

CPU Family 1
Toolchain

CPU Family 2
Toolchain

CPU Family 3
Toolchain

CPU Family 1
Hardware

CPU Family 2
Hardware

CPU Family 3
Hardware

# Big Caches & Speculative Execution

**Big caches**:
- Higher throughput when data fits in cache
- Higher latency for cache misses
- Energy inefficiency

**Speculative execution**:
- Hide latencies and memory misordering
- Higher worst-case latency, energy consumption

# Speculative Execution & Ordering

```
data = 42;

flag.store(1, mo_release)
```

```
while (!flag.load(mo_acquire))
    continue;

r1 = data;
```

# Maligned Workhorse: Locking

Stupidly simple lock acquisition:

```
while (lk.exchange(memory_order_acquire))
    continue;
```

Lock release:

```
lk.store(0, memory_order_release);
```

Works, but **horrible** high-contention behavior

# Maligned Workhorse: Locking

Somewhat less stupidly simple lock acquisition:

```
for (;;) {
    while (lk.exchange(1, memory_order_acquire))
        continue;
    if (!lk.load(memory_order_acquire))
        break;
}
```

Lock release:

```
lk.store(0, memory_order_release);
```

Works, but sub-optimal high-contention behavior

# Suboptimal Lock Handoff

CPU 0

unlock

lk=1

CPU 1

lock

lk=0    lk=1

Don't have
cacheline,
can't unlock!

CPU 2

Can't lock!

lk=1

Time

Lots of useless bus traffic after unlock time!!!
All CPUs spinning on same location...

# Queued Locks

| CPU 0<br>Holds Lock |
| :---: |

↓

| CPU 1<br>First Waiter |
| :---: |

↓

| CPU 2<br>Second Waiter |
| :---: |

**Handoff** ⟩

| CPU 1<br>Holds Lock |
| :---: |

↓

| CPU 2<br>First Waiter |
| :---: |

**Handoff** ⟩

| CPU 2<br>Holds Lock |
| :---: |

## At most two CPUs content for given location

http://www.cs.rochester.edu/~scott/papers/1991_TOCS_synch.pdf

# But Even Better...

Maintain low levels of contention!!!

Also, adaptive spin/sleep strategy
(for example, algorithms using futex
 system calls!)

# Reader-Writer Locking

Read-side parallelism!  What not to like?
Textbook reader-writer lock:

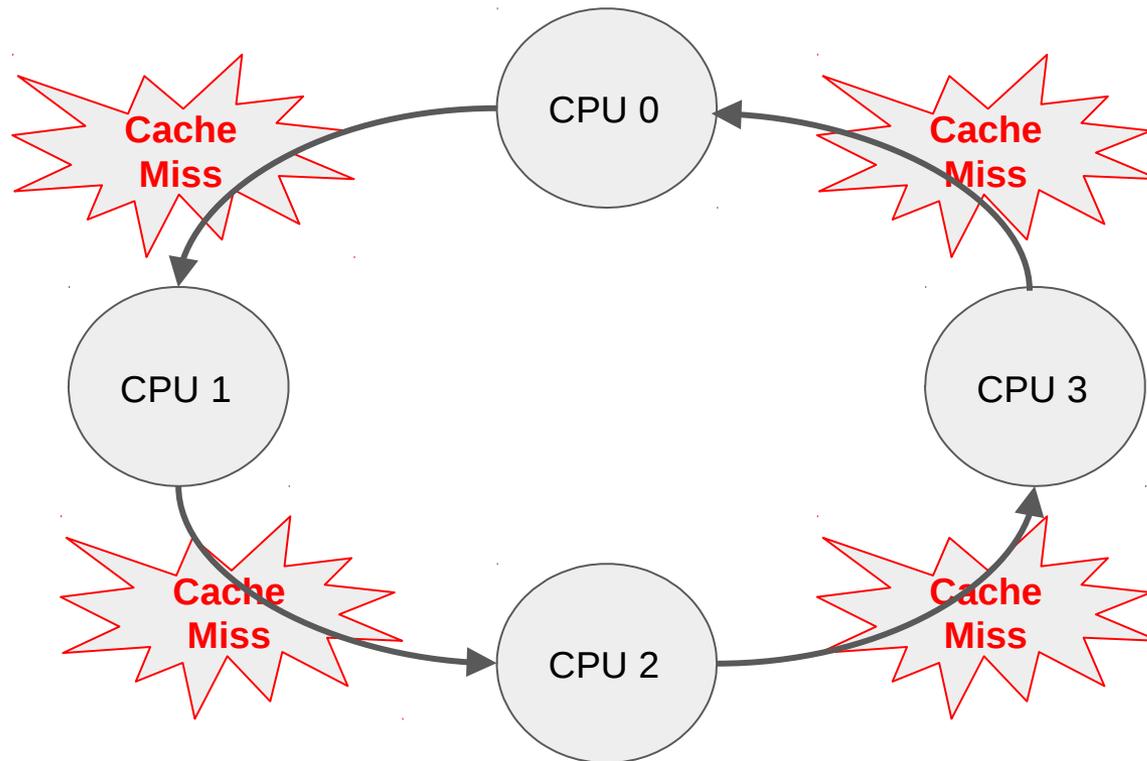**Single atomically manipulated machine word**

**Flag bits (fairness, etc.)**

**Number of Readers**

**Number of Writers**

# Reader-Writer Locking

Textbook read-only acquisition and release:



Orders of magnitude slowdown from cache misses!!!
Might as well have an exclusive lock unless **huge** critical sections...

# Software Must Help Hardware

# TLS Reader-Writer Lock
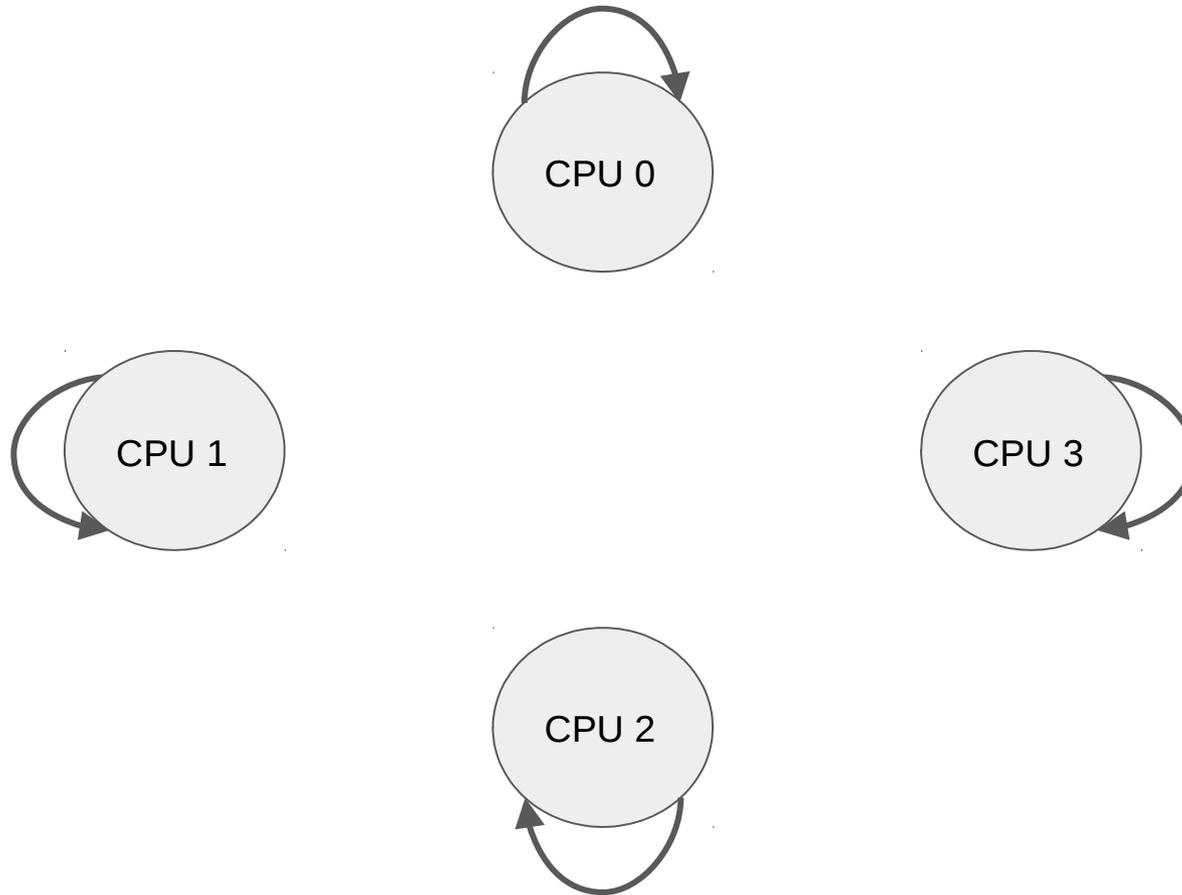
One exclusive lock per thread: TLS
Reading threads acquire and release
 own lock
Writing threads acquire and release **all**
 locks

Very fast readers, but very slow writers

# TLS Reader-Writer Lock

## TLS read-only acquisition and release:



Very fast readers!!!
Too bad about those poor writers...  Generality!!!

# Reader-Writer Locking

And now you understand one motivation for Paul's RCU and Maged's Hazard Pointers...

But that is another presentation.

Now for Maged's single-producer single-consumer buffer!!!