

Advanced OpenMP

Charles Boullaguet

`charles.boullaguet@lip6.fr`

October 18, 2022

atomic is not a Panacea

Example: sum of an array

```
int sum = 0;
for (int i = 0; i < n; i++)
    sum += A[i];
```

$$T = 5.95s \quad (n = 10^{10})$$

```
int sum = 0;
#pragma omp parallel for
for (int i = 0; i < n; i++)
    #pragma omp atomic
    sum += A[i];
```

```
int sum = 0;
#pragma omp parallel for reduction(+:sum)
for (int i = 0; i < n; i++)

    sum += A[i];
```

atomic is not a Panacea

Example: sum of an array

```
int sum = 0;
for (int i = 0; i < n; i++)
    sum += A[i];
```

$$T = 5.95s \quad (n = 10^{10})$$

```
int sum = 0;
#pragma omp parallel for
for (int i = 0; i < n; i++)
    #pragma omp atomic
    sum += A[i];
```

```
int sum = 0;
#pragma omp parallel for reduction(+:sum)
for (int i = 0; i < n; i++)

    sum += A[i];
```

$$T \geq 200s !!!$$

$$T = 0.46s (\times 12.9)$$

Histogram (E.g. `numpy.histogram`)

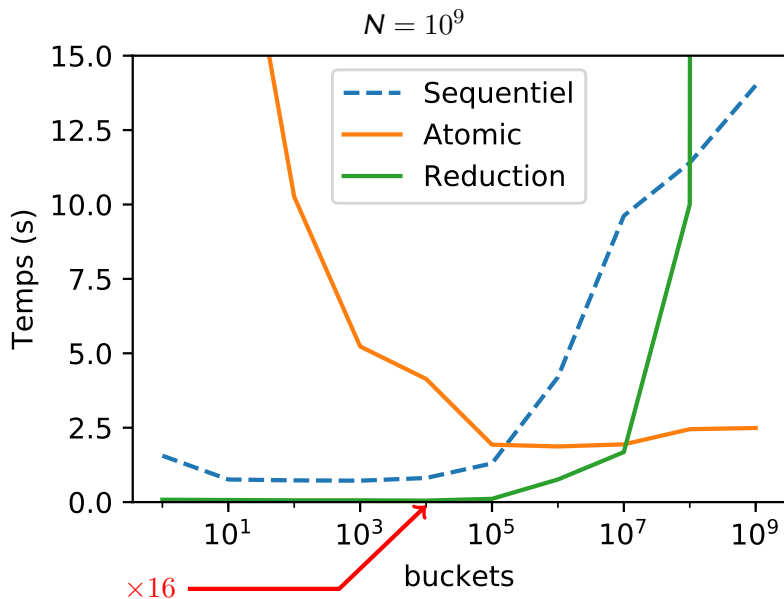
Plan A

```
void histogram(u32 A[], u64 n, u64 buckets, u64 H[])
{
    #pragma omp parallel for
    for (u64 i = 0; i < n; i++) {
        u64 x = (A[i] * buckets) >> 32;
        #pragma omp atomic
        H[x]++;
    }
}
```

Plan B

```
void histogram(u32 A[], u64 n, u64 buckets, u64 H[])
{
    #pragma omp parallel for reduction(+:H[0:buckets])
    for (u64 i = 0; i < n; i++) {
        u64 x = (A[i] * buckets) >> 32;
        H[x]++;
    }
}
```

Histogram



The Sad Truth



- ▶ Barrier → waiting
- ▶ Critical → serialization
- ▶ Atomic → slower than normal memory accesses

Synchronization → limits scalability

⇒ important role of data locality

You **must** follow the rule

You **must** follow the rule

Or else...

You **must** follow the rule

Or else...

You will be living

In a **world of PAIN**

Don't follow the rules?



<https://xkcd.com/292/>

Everyone inevitably goes through this

- ▶ Okay, but then if I avoid `i++`, it's okay, right?
- ▶ I read in my CPU doc that aligned reads/writes are atomic; if we stick to that, we'll be fine, right?

Everyone inevitably goes through this

- ▶ Okay, but then if I avoid `i++`, it's okay, right?
- ▶ I read in my CPU doc that aligned reads/writes are atomic; if we stick to that, we'll be fine, right?

Teenage crisis

Let's drop the golden rule

Everyone inevitably goes through this

- ▶ Okay, but then if I avoid `i++`, it's okay, right?
- ▶ I read in my CPU doc that aligned reads/writes are atomic; if we stick to that, we'll be fine, right?

Teenage crisis

Let's drop the golden rule

Guru switch →



Safety

World of PAIN

Peterson Lock (1981)

```
bool flag[2];
int turn;

void lock()
{
    int i = omp_get_thread_num();
    flag[i] = 1;                               // I'm interested
    turn = 1 - i;                               // you go first
    while (turn != i && flag[1-i]) {};          // wait 'til he's not interested
                                                // or its my turn
}

void unlock()
{
    int i = omp_get_thread_num();
    flag[i] = 0;                               // I'm not interested
}
```

CLAIMS

- ▶ Mutual exclusion
- ▶ Deadlock-free
- ▶ Starvation-free

LIVE DEMO

Guru problem #1: Compiler treachery

- ▶ We had:

```
turn = 1 - i;           // you go first
while (turn != i && flag[1-i]) {}; // wait
```

- ▶ The optimizer “knows” that `turn != i`
- ▶ This then becomes:

```
void lock()
{
    int i = omp_get_thread_num();
    flag[i] = true;           // I'm interested
    while (flag[1-i]) {};    // wait
}
```

- ▶ Both threads call `lock()` at the same time...
- ▶ ... `flag[0] == flag[1] == true` ...
- ▶ Deadlock.

Guru problem #2: ???

- ▶ Mutual exclusion failure
- ▶ Both threads present in the critical section simultaneously
- ▶ How is this possible?

Order Relations Between Concurrent Memory Accesses

► Memory Accesses

$W_i(x)a : T_i$ writes value a in variable x

$R_i(x)b : T_i$ reads variable x and gets value b

► Program Order:

$x \xrightarrow{po} y$: code states that x is done first and y after

► Extended Communication Order:

$W(x)a \xrightarrow{\text{rf}} R(x)a$: the read gets the written value

$W(x)a \xrightarrow{\text{mo}} W(x)b$: a is written before b

$R(x)a \xrightarrow{\text{rb}} W(x)b$: The read takes place before b is written in x
by definition, $\xrightarrow{\text{rb}} = (\xrightarrow{\text{rf}})^{-1}; \xrightarrow{\text{mo}}$

Order Relations Between Concurrent Memory Accesses

► Memory Accesses

$W_i(x)a$: T_i writes value a in variable x

$R_i(x)b$: T_i reads variable x and gets value b

► Program Order:

$x \xrightarrow{po} y$: code states that x is done first and y after

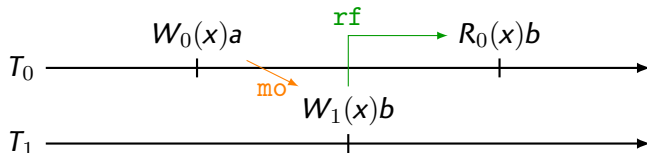
► Extended Communication Order:

$W(x)a \xrightarrow{rf} R(x)a$: the read gets the written value

$W(x)a \xrightarrow{mo} W(x)b$: a is written before b

$R(x)a \xrightarrow{rb} W(x)b$: The read takes place before b is written in x

by definition, $\xrightarrow{rb} = (\xrightarrow{rf})^{-1}$; \xrightarrow{mo}



Order Relations Between Concurrent Memory Accesses

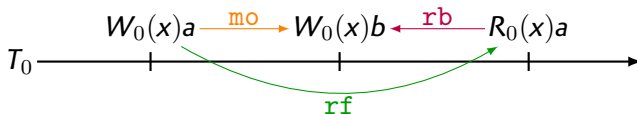
$W(x)a \xrightarrow{\text{rf}} R(x)a$: the read gets the written value

$W(x)a \xrightarrow{\text{mo}} W(x)b$: a is written before b

$R(x)a \xrightarrow{\text{rb}} W(x)b$: The read takes place before b is written in x
by definition, $\xrightarrow{\text{rb}} = (\xrightarrow{\text{rf}})^{-1}; \xrightarrow{\text{mo}}$

Intuition

- ▶ In principle, $\xrightarrow{\text{rf}}$, $\xrightarrow{\text{mo}}$, $\xrightarrow{\text{rb}}$ cannot “contradict” $\xrightarrow{\text{po}}$
- ▶ Cannot read “outdated” values ($\xrightarrow{\text{rb}}$ backwards)



Order Relations Between Concurrent Memory Accesses

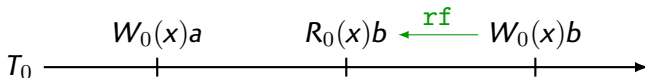
$W(x)a \xrightarrow{\text{rf}} R(x)a$: the read gets the written value

$W(x)a \xrightarrow{\text{mo}} W(x)b$: a is written before b

$R(x)a \xrightarrow{\text{rb}} W(x)b$: The read takes place before b is written in x
by definition, $\xrightarrow{\text{rb}} = (\xrightarrow{\text{rf}})^{-1}; \xrightarrow{\text{mo}}$

Intuition

- ▶ In principle, $\xrightarrow{\text{rf}}$, $\xrightarrow{\text{mo}}$, $\xrightarrow{\text{rb}}$ cannot “contradict” $\xrightarrow{\text{po}}$
- ▶ Cannot read “outdated” values ($\xrightarrow{\text{rb}}$ backwards)
- ▶ Cannot read “from the future” ($\xrightarrow{\text{rf}}$ backwards)



Order Relations Between Concurrent Memory Accesses

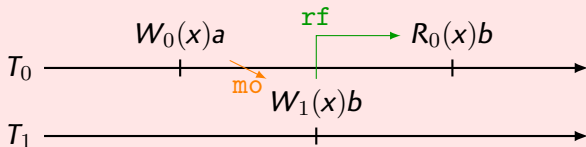
$W(x)a \xrightarrow{\text{rf}} R(x)a$: the read gets the written value

$W(x)a \xrightarrow{\text{mo}} W(x)b$: a is written before b

$R(x)a \xrightarrow{\text{rb}} W(x)b$: The read takes place before b is written in x
by definition, $\xrightarrow{\text{rb}} = (\xrightarrow{\text{rf}})^{-1}; \xrightarrow{\text{mo}}$

Sequentially : $\xrightarrow{\text{rf}} \subseteq \xrightarrow{po}$, $\xrightarrow{\text{mo}} \subseteq \xrightarrow{po}$ and $\xrightarrow{\text{rb}} \subseteq \xrightarrow{po}$

This is not true in parallel



Sequential Consistency

Intuition

Everything happens as if memory accesses were executed *sequentially* (in some order that we do not necessarily know).

Sequential Consistency

Intuition

Everything happens as if memory accesses were executed *sequentially* (in some order that we do not necessarily know).

Definition (Sequential Consistency)

A parallel system is **sequentially consistent** if, for each of the possible executions of the threads to which it can lead, we can build a **history** H :

- ▶ Totally ordered sequence
- ▶ Contains each memory access (only once)
- ▶ Compatible with the code of each thread (respects \xrightarrow{po})
- ▶ Reading x gives the last value written in x

Theorem

Sequential Consistency \iff no cycles with $\xrightarrow{po} \cup \xrightarrow{rf} \cup \xrightarrow{mo} \cup \xrightarrow{rb}$.

Peterson Lock (1981)

```
bool flag[2];
int turn;

void lock()
{
    int i = omp_get_thread_num();
    flag[i] = 1;                               // I'm interested
    turn = 1 - i;                               // you go first
    while (turn != i && flag[1-i]) {};          // wait 'til he's not interested
                                                // or its my turn
}

void unlock()
{
    int i = omp_get_thread_num();
    flag[i] = 0;                               // I'm not interested
}
```

CLAIMS

- ▶ Mutual exclusion
- ▶ Deadlock-free
- ▶ Starvation-free

Correctness of the *Peterson Lock*

Sequential Consistency

Everything happens as if memory accesses were executed *sequentially* (in some order that we do not necessarily know).

Theorem

If the memory is sequentially consistent, then the Peterson lock guarantees mutual exclusion.

Proof: ...

► Suppose not: both T_0 and T_1 call `lock()`, enter the critical section

► Initial state: `flag[0] = false; flag[1] = false;`

► Code requires that:

$T_0 : W_0(\text{flag}[0]) \text{ true} \xrightarrow{po} W_0(\text{turn})1 \xrightarrow{po} R_0(\text{turn})? \xrightarrow{po} R_0(\text{flag}[1])? \xrightarrow{po} CS_0$

$T_1 : W_1(\text{flag}[1]) \text{ true} \xrightarrow{po} W_1(\text{turn})0 \xrightarrow{po} R_1(\text{turn})? \xrightarrow{po} R_1(\text{flag}[0])? \xrightarrow{po} CS_1$

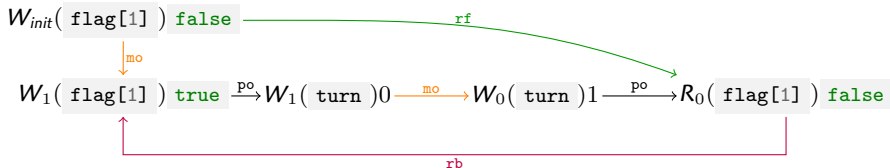
► Suppose that T_0 writes `turn` last:

$W_1(\text{turn})0 \xrightarrow{mo} W_0(\text{turn})1.$

► T_0 exits the loop, and `turn == 1`, then necessarily:

$W_0(\text{turn})1 \xrightarrow{po} R_0(\text{turn})1 \xrightarrow{po} R_0(\text{flag}[1]) \text{ false}.$

► Put all this together:



► Cycle \Rightarrow non-SC \Rightarrow Contradiction!

Guru problem #2: **lack of *sequential consistency***

- ▶ Proof: sequential consistency \Rightarrow mutual exclusion
- ▶ Observation : ~~mutual exclusion~~
- ▶ *Ergo*:



- ▶ My laptop is not sequentially consistent...

Guru problem #2: **lack of *sequential consistency***

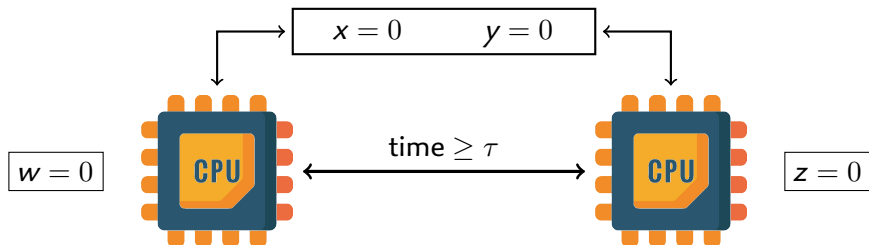
- ▶ Proof: sequential consistency \Rightarrow mutual exclusion
- ▶ Observation : ~~mutual exclusion~~
- ▶ *Ergo*:



- ▶ My laptop is not sequentially consistent...
- ▶ Yours is not either!

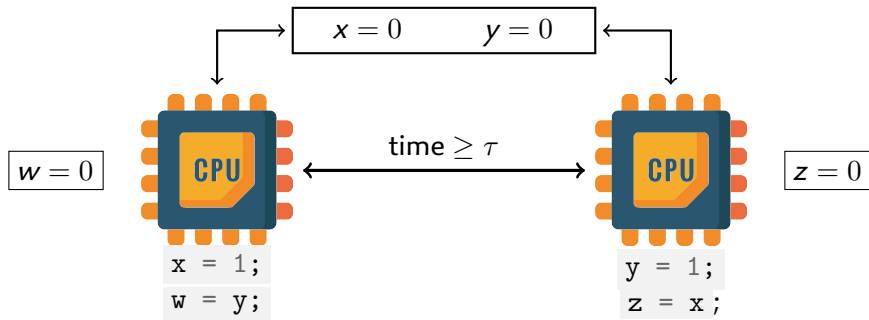
CPUs are not sequentially consistent!

Sequential Consistency is costly



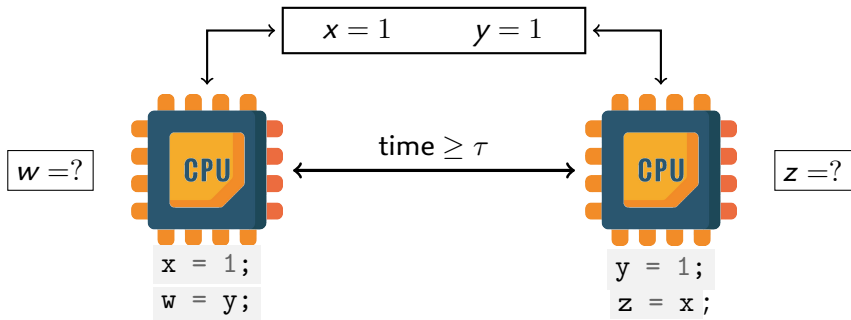
CPUs are not sequentially consistent!

Sequential Consistency is **costly**



CPUs are not sequentially consistent!

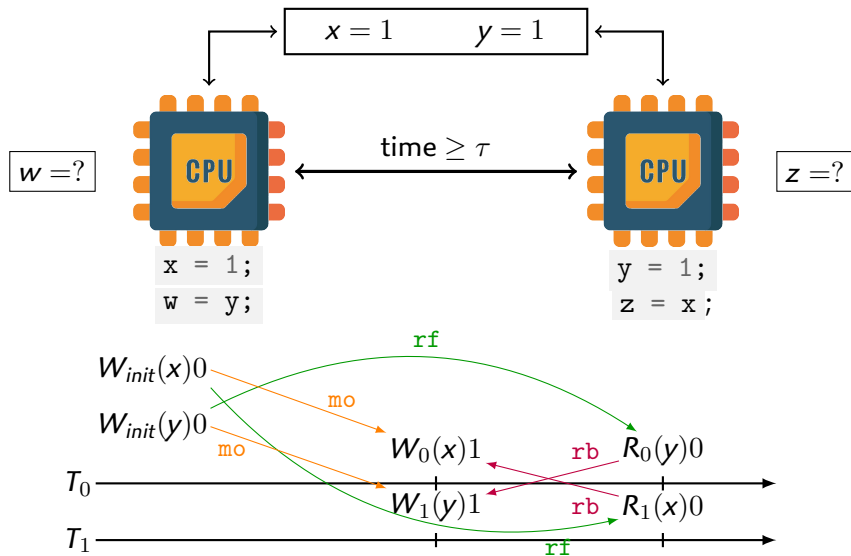
Sequential Consistency is **costly**



Sequential Consistency $\Rightarrow (w, z) \neq (0, 0)$

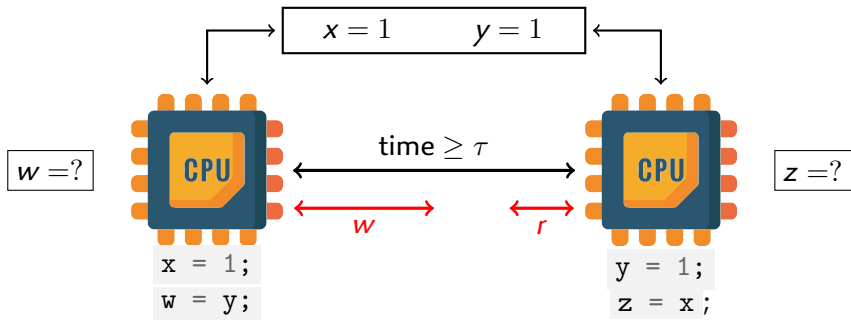
CPUs are not sequentially consistent!

Sequential Consistency is **costly**



CPUs are not sequentially consistent!

Sequential Consistency is *costly*

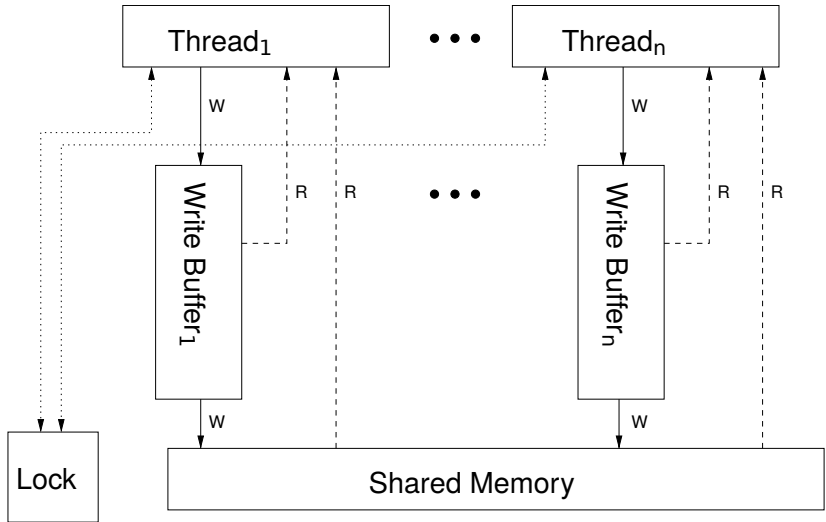


r : min. time to do a read
 w : min. time to do a write

Sequential Consistency $\implies r + w \geq \tau$

One must read a write from the other...

Store Buffering



(image : A Tutorial Introduction to the ARM and POWER Relaxed Memory Models)

Architectures with *Total Store Ordering*

x86 and SPARC

Each **hardware** thread has a **Store Buffer**

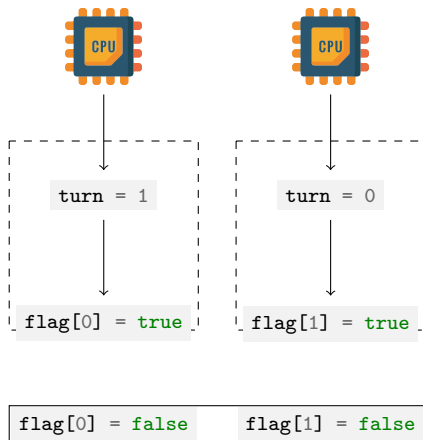
- ▶ My writes are enqueued in my *Store Buffer*
- ▶ FIFO queue
- ▶ I see the effects of my writes immediately
- ▶ My *Store Buffer* **will** be “flushed” to the memory... later
- ▶ At that time:
 - ▶ *All the other threads “see” my writes*
 - ▶ They see them in the same order (the order I did them)

Peterson Lock Failure in the presence of Store Buffering

Peterson Lock

```
bool flag[2];
int victim;

void lock()
{
    int i = omp_get_thread_num();
    flag[i] = true;
    turn = 1 - i;
    while (turn != i && flag[1-i])
        {};
}
```

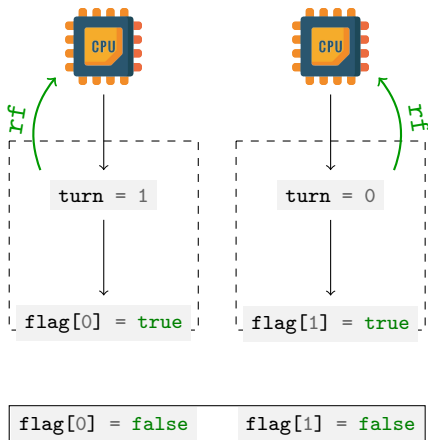


Peterson Lock Failure in the presence of Store Buffering

Peterson Lock

```
bool flag[2];
int victim;

void lock()
{
    int i = omp_get_thread_num();
    flag[i] = true;
    turn = 1 - i;
    while (turn != i && flag[1-i])
        {};
}
```

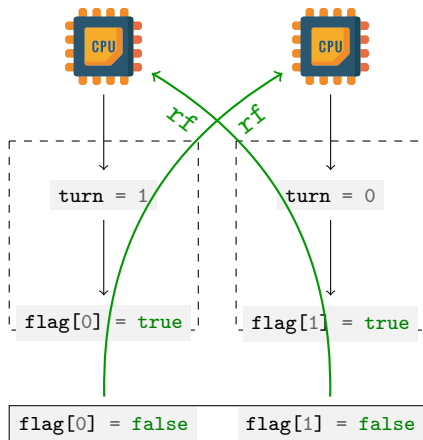


Peterson Lock Failure in the presence of Store Buffering

Peterson Lock

```
bool flag[2];
int victim;

void lock()
{
    int i = omp_get_thread_num();
    flag[i] = true;
    turn = 1 - i;
    while (turn != i && flag[1-i])
        {};
}
```

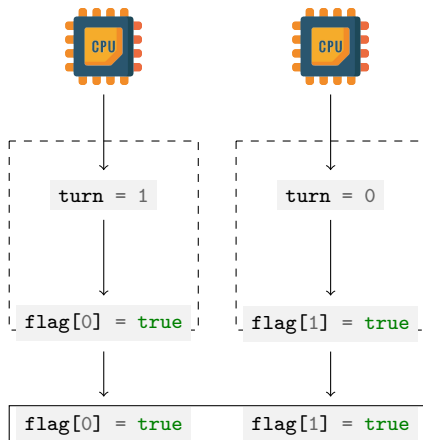


Peterson Lock Failure in the presence of Store Buffering

Peterson Lock

```
bool flag[2];
int victim;

void lock()
{
    int i = omp_get_thread_num();
    flag[i] = true;
    turn = 1 - i;
    while (turn != i && flag[1-i])
        {};
}
```



Message Passing

Producer

```
void produce(void *payload)
{
    msg = payload;
    flag = 1;
}
```

Consumer

```
void * receive()
{
    while (!flag) {}; // wait
    return msg;
}
```

Message Passing

Producer

```
void produce(void *payload)
{
    msg = payload;
    flag = 1;
}
```

Consumer

```
void * receive()
{
    while (!flag) {}; // wait
    return msg;
}
```



- ▶ Works with *Total Store Ordering*
- ▶ x86, SPARC, ...

Message Passing

Producer

```
void produce(void *payload)
{
    msg = payload;
    flag = 1;
}
```

Consumer

```
void * receive()
{
    while (!flag) {}; // wait
    return msg;
}
```



- ▶ Fails on ARM, POWER, ...
- ▶ Threads do **not necessarily** see the writes in the same order!

Message Passing

Producer

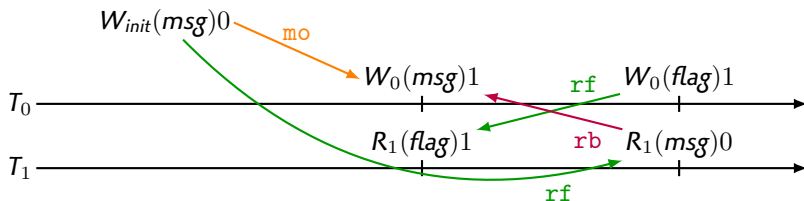
```
void produce(void *payload)
{
    msg = payload;
    flag = 1;
}
```

Consumer

```
void * receive()
{
    while (!flag) {}; // wait
    return msg;
}
```



- ▶ Fails on ARM, POWER, ...
- ▶ Threads do **not necessarily** see the writes in the same order!



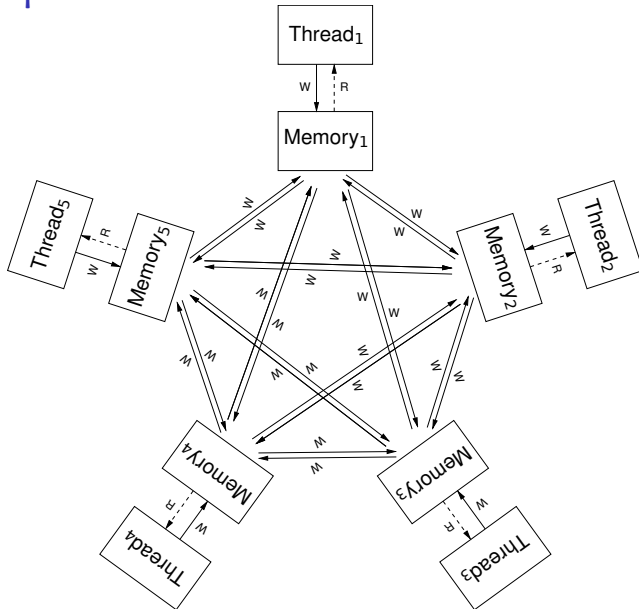
The OpenMP Memory Model

Threads have access to the same memory, but...

Each thread has a **private temporary view** of the memory

- ▶ Not necessarily always synchronized
 - ▶ A read **may** come from the private temporary view
 - ▶ A write **may** stay inside the private temporary view
 - ▶ (Implicit) synchronizations on:
 - ▶ `#pragma omp barrier`
 - ▶ Exit of `#pragma omp for/sections/single`
 - ▶ Entry/exit of `#pragma omp parallel/critical/atomic`
 - ▶ *Task Scheduling Points*
 - ▶ (Explicit) synchronizations with `#pragma omp flush`
-
- ▶ Various reasons, including compilers and the hardware itself
 - ▶ Atomic operations are **sequentially consistent**

Visual Representation



(image : A Tutorial Introduction to the ARM and POWER Relaxed Memory Models)

Quizz

Initially, $x = y = 0$.

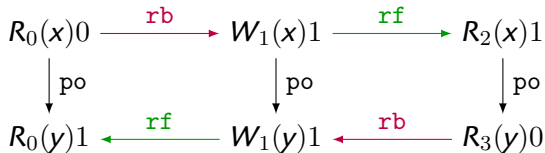
 $R_0(x)0$ $\downarrow \text{po}$ $R_0(y)1$ $W_1(x)1$ $\downarrow \text{po}$ $W_1(y)1$ $R_2(x)1$ $\downarrow \text{po}$ $R_3(y)0$

Possible?

1. No, it is contradictory
2. No, because it is not sequentially consistent
3. Yes, possible on ARM and POWER CPUs but not on x86
4. Yes, it is sequentially consistent

Quizz

Initially, $x = y = 0$.

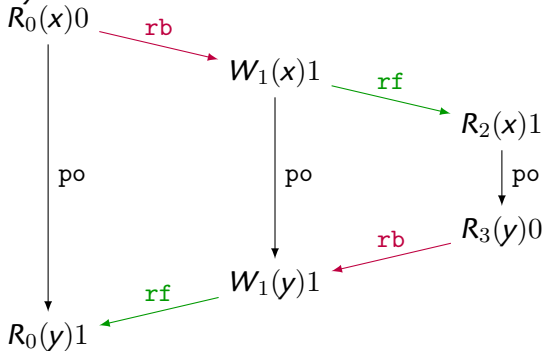


Possible?

1. No, it is contradictory
2. No, because it is not sequentially consistent
3. Yes, possible on ARM and POWER CPUs but not on x86
4. Yes, it is sequentially consistent

Quizz

Initially, $x = y = 0$.



Possible?

1. No, it is contradictory
2. No, because it is not sequentially consistent
3. Yes, possible on ARM and POWER CPUs but not on x86
4. Yes, it is sequentially consistent

Quizz (harder)

Initially, $x = y = 0$.

$R_0(y)1$

$\downarrow \text{po}$

$R_0(x)0$

$W_1(x)1$

$\downarrow \text{po}$

$W_1(y)1$

$R_2(x)1$

$\downarrow \text{po}$

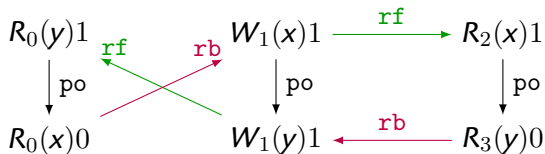
$R_3(y)0$

Possible ?

1. No, it is contradictory
2. No, because it is not sequentially consistent
3. Yes, possible on ARM and POWER but not on x86
4. Yes, it is sequentially consistent

Quizz (harder)

Initially, $x = y = 0$.

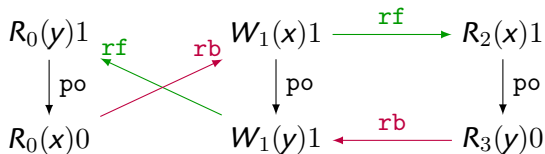


Possible ?

1. No, it is contradictory
2. No, because it is not sequentially consistent
3. Yes, possible on ARM and POWER but not on x86
4. Yes, it is sequentially consistent

Quizz (harder)

Initially, $x = y = 0$.



Possible ?

1. No, it is contradictory
2. No, because it is not sequentially consistent
3. Yes, possible on ARM and POWER but not on x86
4. Yes, it is sequentially consistent

Everyone inevitably goes through this

- ▶ Okay, but then if I avoid `i++`, it's okay, right?
- ▶ I read in my CPU doc that aligned reads/writes are atomic; if we stick to that, we'll be fine, right?

Teenage crisis

Let's drop the golden rule

Guru switch →



Safety

World of PAIN

Golden rule of multi-thread programming

ALL potentially conflicting accesses* to **shared** variables **MUST** be protected (atomic, critical, ...).

* at least one of them is a write

Case Study: Bucket Sort

```
// Initialization
for (int i = 0; i < M; i++) {
    C[i] = 0;
}

// Histogram
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    C[bucket]++;
}

// Prefix-sum
int s = 0;
for (int i = 0; i < M; i++) {
    P[i] = s;
    s += C[i];
}

// Dispatch
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    B[P[bucket]] = A[i];
    P[bucket]++;
}
```



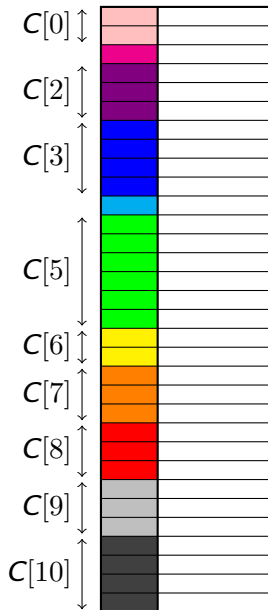
Case Study: Bucket Sort

```
// Initialization
for (int i = 0; i < M; i++) {
    C[i] = 0;
}

// Histogram
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    C[bucket]++;
}

// Prefix-sum
int s = 0;
for (int i = 0; i < M; i++) {
    P[i] = s;
    s += C[i];
}

// Dispatch
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    B[P[bucket]] = A[i];
    P[bucket]++;
}
```



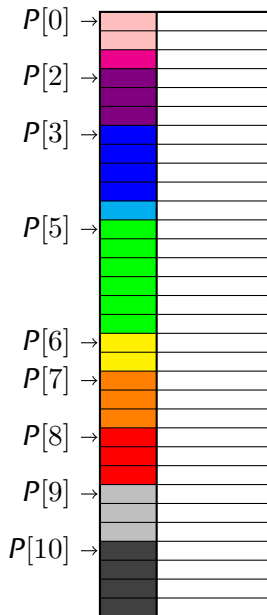
Case Study: Bucket Sort

```
// Initialization
for (int i = 0; i < M; i++) {
    C[i] = 0;
}

// Histogram
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    C[bucket]++;
}

// Prefix-sum
int s = 0;
for (int i = 0; i < M; i++) {
    P[i] = s;
    s += C[i];
}

// Dispatch
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    B[P[bucket]] = A[i];
    P[bucket]++;
}
```



Case Study: Bucket Sort

```
// Initialization
```

```
for (int i = 0; i < M; i++) {  
    C[i] = 0;  
}
```

```
// Histogram
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    C[bucket]++;  
}
```

```
// Prefix-sum
```

```
int s = 0;  
for (int i = 0; i < M; i++) {  
    P[i] = s;  
    s += C[i];  
}
```

```
// Dispatch
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    B[P[bucket]] = A[i];  
    P[bucket]++;  
}
```



$P[0]$	→		
$P[2]$	→		
$P[3]$	→		
$P[5]$	→		
$P[6]$	→		
$P[7]$	→		
$P[8]$	→		
$P[9]$	→		
$P[10]$	→		

Case Study: Bucket Sort

```
// Initialization
```

```
for (int i = 0; i < M; i++) {  
    C[i] = 0;  
}
```

```
// Histogram
```

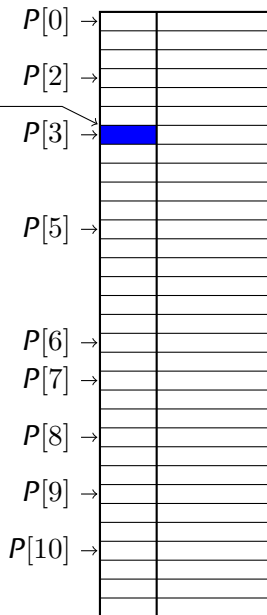
```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    C[bucket]++;  
}
```

```
// Prefix-sum
```

```
int s = 0;  
for (int i = 0; i < M; i++) {  
    P[i] = s;  
    s += C[i];  
}
```

```
// Dispatch
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    B[P[bucket]] = A[i];  
    P[bucket]++;  
}
```



Case Study: Bucket Sort


```
// Initialization
for (int i = 0; i < M; i++) {
    C[i] = 0;
}

// Histogram
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    C[bucket]++;
}

// Prefix-sum
int s = 0;
for (int i = 0; i < M; i++) {
    P[i] = s;
    s += C[i];
}

// Dispatch
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    B[P[bucket]] = A[i];
    P[bucket]++;
}
```

	3
	0
	7
	1
	10
	2
	5
	2
	10
	5
	2
	3
	8
	7
	9
	10
	3
	6
	4
	6
	10
	0
	5
	3
	5
	9
	5
	7
	8
	9
	8
	5

$P[0]$	→		
$P[2]$	→		
$P[3]$	→		
$P[5]$	→		
$P[6]$	→		
$P[7]$	→		
$P[8]$	→		
$P[9]$	→		
$P[10]$	→		

Case Study: Bucket Sort

```
// Initialization
```

```
for (int i = 0; i < M; i++) {  
    C[i] = 0;  
}
```

```
// Histogram
```

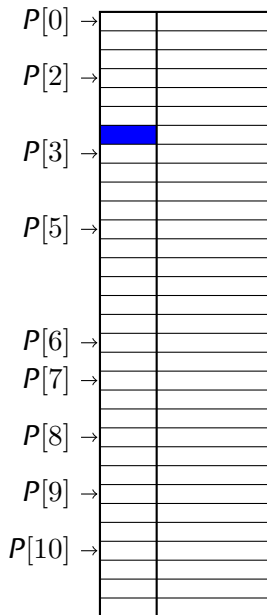
```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    C[bucket]++;  
}
```

```
// Prefix-sum
```

```
int s = 0;  
for (int i = 0; i < M; i++) {  
    P[i] = s;  
    s += C[i];  
}
```

```
// Dispatch
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    B[P[bucket]] = A[i];  
    P[bucket]++;  
}
```



Case Study: Bucket Sort

```
// Initialization
```

```
for (int i = 0; i < M; i++) {  
    C[i] = 0;  
}
```

```
// Histogram
```

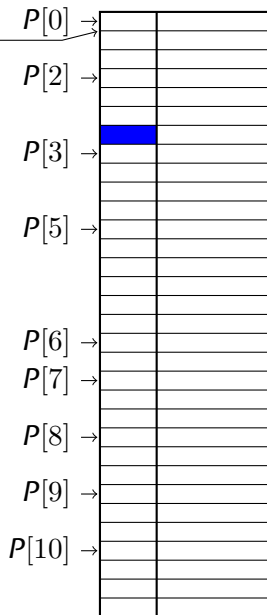
```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    C[bucket]++;  
}
```

```
// Prefix-sum
```

```
int s = 0;  
for (int i = 0; i < M; i++) {  
    P[i] = s;  
    s += C[i];  
}
```

```
// Dispatch
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    B[P[bucket]] = A[i];  
    P[bucket]++;  
}
```



Case Study: Bucket Sort

```
// Initialization
for (int i = 0; i < M; i++) {
    C[i] = 0;
}

// Histogram
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    C[bucket]++;
}

// Prefix-sum
int s = 0;
for (int i = 0; i < M; i++) {
    P[i] = s;
    s += C[i];
}

// Dispatch
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    B[P[bucket]] = A[i];
    P[bucket]++;
}
```

	3
	0
	7
	1
	10
	2
	5
	2
	10
	5
	2
	3
	8
	7
	9
	10
	3
	6
	4
	6
	10
	0
	5
	3
	5
	9
	5
	7
	8
	9
	8
	5

$P[0]$	→		
$P[2]$	→		
$P[3]$	→		
$P[5]$	→		
$P[6]$	→		
$P[7]$	→		
$P[8]$	→		
$P[9]$	→		
$P[10]$	→		

Case Study: Bucket Sort

```
// Initialization
```

```
for (int i = 0; i < M; i++) {  
    C[i] = 0;  
}
```

```
// Histogram
```

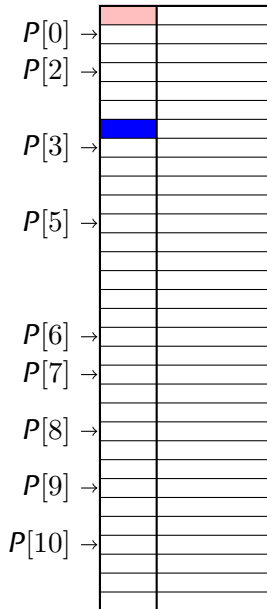
```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    C[bucket]++;  
}
```

```
// Prefix-sum
```

```
int s = 0;  
for (int i = 0; i < M; i++) {  
    P[i] = s;  
    s += C[i];  
}
```

```
// Dispatch
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    B[P[bucket]] = A[i];  
    P[bucket]++;  
}
```



Case Study: Bucket Sort

```
// Initialization
```

```
for (int i = 0; i < M; i++) {  
    C[i] = 0;  
}
```

```
// Histogram
```

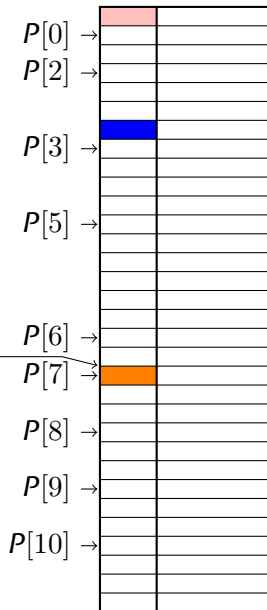
```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    C[bucket]++;  
}
```

```
// Prefix-sum
```

```
int s = 0;  
for (int i = 0; i < M; i++) {  
    P[i] = s;  
    s += C[i];  
}
```

```
// Dispatch
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    B[P[bucket]] = A[i];  
    P[bucket]++;  
}
```



Case Study: Bucket Sort




```
// Initialization
for (int i = 0; i < M; i++) {
    C[i] = 0;
}

// Histogram
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    C[bucket]++;
}

// Prefix-sum
int s = 0;
for (int i = 0; i < M; i++) {
    P[i] = s;
    s += C[i];
}

// Dispatch
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    B[P[bucket]] = A[i];
    P[bucket]++;
}
```

	3
	0
	7
	1
	10
	2
	5
	2
	10
	5
	2
	3
	8
	7
	9
	10
	3
	6
	4
	6
	10
	0
	5
	3
	5
	9
	5
	7
	8
	9
	8
	5

$P[0]$	→		
$P[2]$	→		
$P[3]$	→		
$P[5]$	→		
$P[6]$	→		
$P[7]$	→		
$P[8]$	→		
$P[9]$	→		
$P[10]$	→		

Case Study: Bucket Sort

Direct Naïve Parallelization

```
// Counting
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    C[bucket]++;  
}
```

```
// Prefix-sum
```

```
int s = 0;  
for (int i = 0; i < M; i++) {  
    P[i] = s;  
    s += C[i];  
}
```

```
// Dispatch
```

```
for (int i = 0; i < N; i++) {  
    int bucket = f(A[i]);  
    int ptr;  
  
    ptr = P[bucket]++;  
    B[ptr] = A[i];  
}
```



Case Study: Bucket Sort

Direct Naïve Parallelization

```
// Counting
#pragma omp parallel for reduction(+:C[0:M])
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    C[bucket]++;
}

// Prefix-sum (sequential)
int s = 0;
for (int i = 0; i < M; i++) {
    P[i] = s;
    s += C[i];
}

// Dispatch
#pragma omp parallel for
for (int i = 0; i < N; i++) {
    int bucket = f(A[i]);
    int ptr;
    #pragma omp atomic capture
    ptr = P[bucket]++;
    B[ptr] = A[i];
}
```

T_0		3
		0
		7
		1
		10
		2
		5
		2
		10
		5
T_1		2
		3
		8
		7
		9
		10
		3
		6
		4
		6
T_2		10
		0
		5
		3
		5
		9
		5
		7
		8
		9
T_3		8
		5
		5

General Principle n°1: **reorganize**

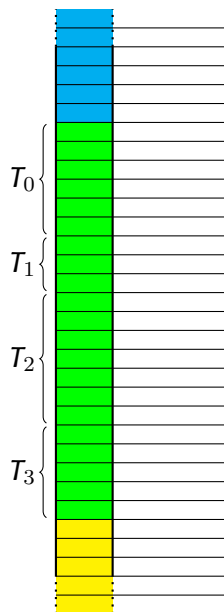


- ▶ Do a (tiny) bit of extra computation...

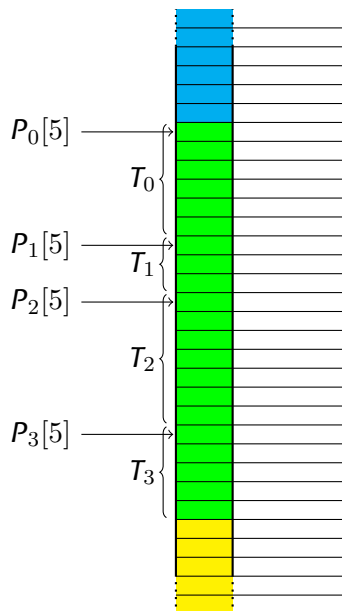


- ▶ ... to completely eliminate conflicts

Case Study: Bucket Sort

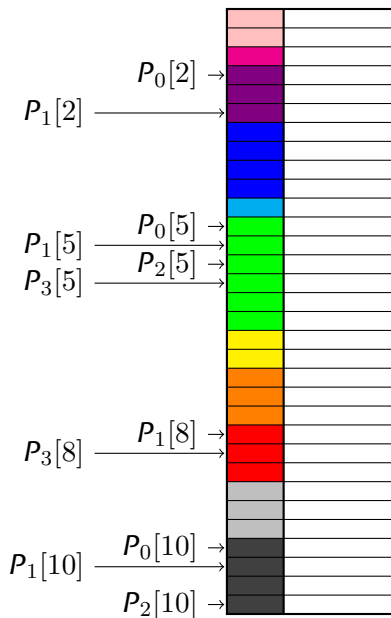


Case Study: Bucket Sort



Case Study: Bucket Sort

T_0		3
		0
		7
		1
		10
		2
		5
		2
T_1		10
		5
		2
		3
		8
		7
		9
		10
T_2		3
		6
		4
		6
		10
		0
		5
		3
T_3		5
		9
		5
		7
		8
		9
		8
		5



Case Study: Bucket Sort

T_0		3
		0
		7
		1
		10
		2
		5
		2
T_1		10
		5
		2
		3
		8
		7
		9
		10
T_2		3
		6
		4
		6
		10
		0
		5
		3
T_3		5
		9
		5
		7
		8
		9
		8
		5

$C_i[x] = \text{\#items of kind } x \text{ seen by thread } i$

		Buckets											
		C	0	1	2	3	4	5	6	7	8	9	10
Threads	0	1	1	2	1		1		1				1
	1			1	1		1		1	1	1	2	
	2	1			2	1	1	2					1
	3						3		1	2	2		

Case Study: Bucket Sort

T_0		3
		0
		7
		1
		10
		2
T_1		5
		2
		10
		5
		3
		8
T_2		7
		9
		10
		3
		6
		4
T_3		6
		10
		0
		5
		3
		5
T_4		9
		5
		7
		8
		5

$$C_i[x] = \begin{array}{l} \text{\#items of kind } x \\ \text{seen by thread } i \end{array}$$

The diagram illustrates a parallel bucket sort algorithm. It shows a grid of 4 threads (rows) and 11 buckets (columns). Red arrows indicate the distribution of elements from threads to buckets. The final row shows the size of each bucket after distribution.

	Buckets										
	0	1	2	3	4	5	6	7	8	9	10
0	1	1	2	1	1	1	1	1	1	1	1
1			1	1		1		1	1	1	2
2	1			2	1	1	2				1
3						3		1	2	2	
Size of each bucket	2	1	3	4	1	6	2	3	3	3	4

Case Study: Bucket Sort

T_0	3
	0
	7
	1
	10
	2
	5
T_1	2
	10
	5
	2
	3
	8
	7
T_2	9
	10
	3
	6
	4
	6
	10
T_3	0
	5
	3
	5
	9
	5
	7
	8
	9
	8
	5
	5

$C_i[x] = \text{\#items of kind } x \text{ seen by thread } i$

		Buckets											
		C	0	1	2	3	4	5	6	7	8	9	10
Threads	0												
	1												
	2												
	3												
			2	1	3	4	1	6	2	3	3	3	4

Prefix-Sum

Size of each bucket

Case Study: Bucket Sort



$C_i[x] = \text{\#items of kind } x \text{ seen by threads } < i$

		Buckets											
		C	0	1	2	3	4	5	6	7	8	9	10
Threads	0	0	0	0	0	0	0	0	0	0	0	0	0
	1	1	1	2	1	0	1	0	1	0	0	1	
	2	1	1	3	2	0	1	0	2	1	1	3	
	3	2	1	3	4	1	3	2	2	1	1	4	
		2	1	3	4	1	6	2	3	3	3	4	

Size of each bucket

Case Study: Bucket Sort

T_0		3
		0
		7
		1
		10
		2
		5
T_1		2
		10
		5
		2
		3
		8
		7
T_2		9
		10
		3
		6
		4
		6
		10
T_3		0
		5
		3
		5
		9
		5
		7
		8
		9
		8
		5
		5

$C_i[x] = \text{\#items of kind } x$
seen by threads $< i$

		Buckets											
		C	0	1	2	3	4	5	6	7	8	9	10
Threads	0	0	0	0	0	0	0	0	0	0	0	0	0
	1	1	1	2	1	0	1	0	1	0	0	1	
	2	1	1	3	2	0	1	0	2	1	1	3	
	3	2	1	3	4	1	3	2	2	1	1	4	
		2	1	3	4	1	6	2	3	3	3	4	

Prefix-Sum

Case Study: Bucket Sort

T_0		3
		0
		7
		1
		10
		2
		5
		2
T_1		10
		5
		2
		3
		8
		7
		9
		10
T_2		3
		6
		4
		6
		10
		0
		5
		3
T_3		5
		9
		5
		7
		8
		9
		8
		5

$C_i[x] = \text{\#items of kind } x$
seen by threads $< i$

		Buckets											
		C	0	1	2	3	4	5	6	7	8	9	10
Threads	0	0	0	0	0	0	0	0	0	0	0	0	0
	1	1	1	2	1	0	1	0	1	0	0	1	
	2	1	1	3	2	0	1	0	2	1	1	3	
	3	2	1	3	4	1	3	2	2	1	1	4	
		0	2	3	6	10	11	17	19	22	25	28	

Start indice of each bucket
(#items in previous buckets)

Case Study: Bucket Sort

T_0		3
		0
		7
		1
		10
		2
		5
T_1		2
		10
		5
		2
		3
		8
		7
T_2		9
		10
		3
		6
		4
		6
		10
T_3		0
		5
		3
		5
		9
		5
		7
		8
		9
		8
		5
		5

$$C_i[x] = \begin{array}{l} \text{\#items of kind } x \\ \text{seen by threads } < i \end{array}$$

Buckets												
	C	0	1	2	3	4	5	6	7	8	9	10
Threads	0	0	0	0	0	0	0	0	0	0	0	0
	1	1	1	2	1	0	1	0	1	0	0	1
	2	1	1	3	2	0	1	0	2	1	1	3
	3	2	1	3	4	1	3	2	2	1	1	4
		0	2	3	6	10	11	17	19	22	25	28

Sum

Start indice of each bucket
(#items in previous buckets)

Case Study: Bucket Sort



$C_i[x]$ = Start indice of kind x
for thread i

		Buckets											
		C	0	1	2	3	4	5	6	7	8	9	10
Threads	0	0	2	3	6	10	11	17	19	22	25	28	
	1	1	3	5	7	10	12	17	20	22	25	29	
	2	1	3	6	8	10	12	17	21	23	26	31	
	3	2	3	6	10	11	14	19	21	23	26	32	
		0	2	3	6	10	11	17	19	22	25	28	

Start indice of each bucket
(#items in previous buckets)

Case Study: Bucket Sort

```
int C[T][M], S[M];

#pragma omp parallel
{
    int t = omp_get_thread_num();

    // Counting
    for (int i = 0; i < M; i++)
        C[t][i] = 0;

    #pragma omp for schedule(static)
    for (int i = 0; i < N; i++) {
        int bucket = f(A[i]);
        C[t][bucket]++;

        // <<COMPUTE POINTERS>> ----->
    }

    // Dispatch
    #pragma omp for schedule(static)
    for (int i = 0; i < N; i++) {
        int bucket = f(A[i]);
        int ptr = C[t][bucket]++;
        B[ptr] = A[i];
    }
}
```

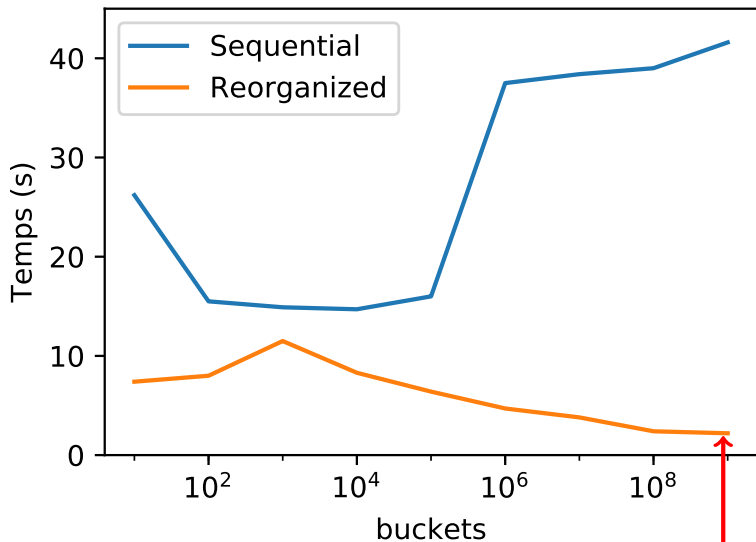
```
// sum (columns)
#pragma omp for
for (int i = 0; i < M; i++) {
    S[i] = 0;
    for (int j = 0; j < T; j++)
        S[i] += C[j][i];
}

// horizontal prefix-sum (sequential)
#pragma omp single
{
    int s = 0;
    for (int i = 0; i < M; i++) {
        int t = S[i];
        S[i] = s;
        s += t;
    }
}

// prefix-sum (columns)
#pragma omp for
for (int i = 0; i < M; i++) {
    int s = S[i];
    for (int j = 0; j < T; j++) {
        int t = C[j][i];
        C[j][i] = s;
        s += t;
    }
}
```

Case Study: Bucket Sort

$$N = 10^{10}$$



$\times 19$

An array of integers to sort?

Pro Tip

Practical parallel sorting algorithm

- ▶ *Parallel Bucket Sort* on the 8 most significant bits
- ▶ For $0 \leq i < 2^8$, do (in parallel) :
 - ▶ Sort the i -th Bucket (using a normal sequential sort)

Parallel Transactions

Read $A[i_1], A[i_2], \dots \rightarrow$ **Compute** \rightarrow **Write** $A[k_1], A[k_2], \dots$

Obstacle to “atomic” execution:

- ▶ The read data was modified before the writes
- ▶ Result of the computation is “outdated”

Pessimistic approach (*“Ask for Permission”*)

- ▶ “Lock” read data
- ▶ Read/Acquire locks \rightarrow Compute \rightarrow Write \rightarrow Release locks
 - ▶ Prevent **potential** modification by other threads
- ▶ Assume that conflict **WILL** take place
- ▶ Useless overhead in the absence of conflict

Parallel Transactions

Read $A[i_1], A[i_2], \dots \rightarrow$ **Compute** \rightarrow **Write** $A[k_1], A[k_2], \dots$

Obstacle to “atomic” execution:

- ▶ The read data was modified before the writes
- ▶ Result of the computation is “outdated”

Optimistic Approach (*“Shoot First, Ask Questions Later”*)

- ▶ Read (**without precaution!!!**) \rightarrow Compute \rightarrow **Commit**:
 - ▶ Check freshness of read data,
 - ▶ If OK (=unmodified), write; otherwise, restart from the beginning
- ▶ Assume that conflict **WILL NOT** take place
- ▶ Lost work in case of conflict

General Principle: **Analyze Conflict Frequency**



- ▶ Take the risk to waste a little bit of computation...



- ▶ ... To reduce the cost of handling conflicts

Generic Technique: *Versioning*

- ▶ Shared data structure, with a **version number**
- ▶ $v = 0$ initially (**odd** during writes)

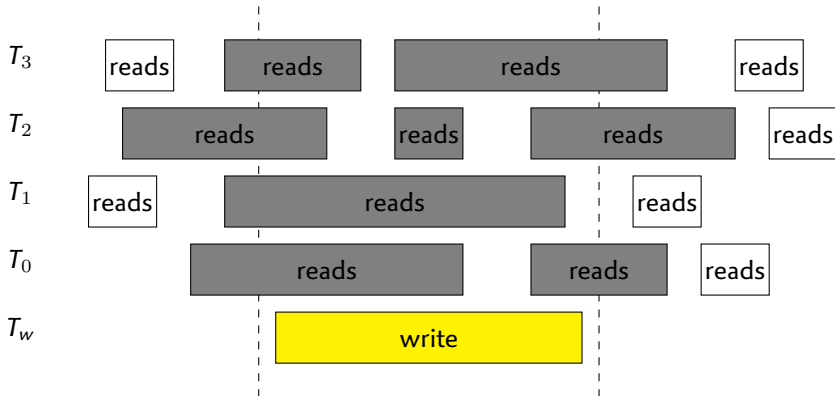
Writer

1. Enter critical section; increment v
2. Do the writes
3. increment v ; exit critical section

Reader

1. $v_{before} \leftarrow v$
2. Do the reads
3. $v_{after} \leftarrow v$
4. If v_{before} is odd or $v_{before} \neq v_{after}$, retry

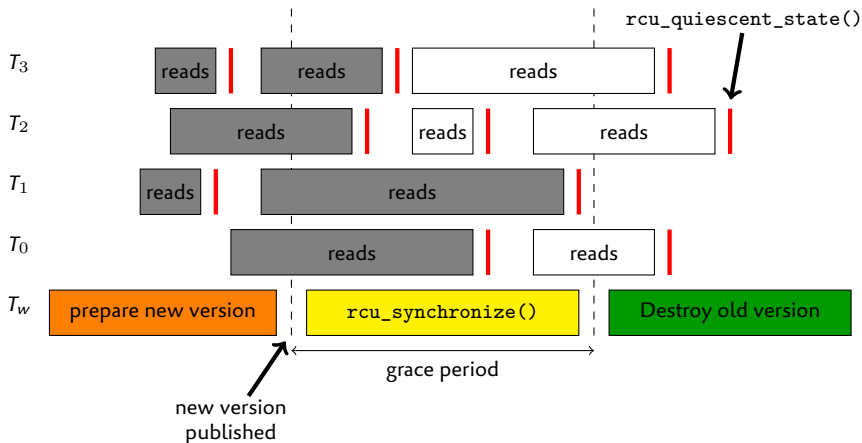
Generic Technique: *Versioning*



- Writers have priority over readers

Generic Technique: *Read-Copy-Update*

- ▶ Shared data structure, readers access via a pointer
- ▶ Readers grab the pointer, read, call `rcu_quiescent_state()`
- ▶ Writers make new copy, atomically update the pointer
- ▶ **Cannot** release old copy right now (readers still read it)



```
int tid = omp_get_num_thread();
int N = omp_get_max_threads();
int gc = 1;
int rc[N]; // {0, ..., 0}
```

```
void rcu_quiescent_state()
{
    #pragma omp atomic write
    rc[tid] = gc;
}
```

```
void rcu_thread_offline()
{
    #pragma omp atomic write
    rc[tid] = 0;
}
```

```
void rcu_thread_online()
{
    #pragma omp atomic write
    rc[tid] = gc;
}
```

```
void rcu_synchronize()
{
    bool was_online = (rc[tid] > 0);
    if (was_online)
        rcu_thread_offline();
    #pragma omp critical
    {
        #pragma omp atomic update
        gc++;
        for (int i = 0; i < N; i++)
            for (;;) {
                int s;
                #pragma omp atomic read
                s = rc[i];
                if (s == 0 || s == gc)
                    break;
            }
    }
    if (was_online)
        rcu_thread_online();
}
```

Transactional Memory

- ▶ Similar problems in database servers
- ▶ Many concurrent transactions management techniques

(very) modern CPUs: transactional memory

```
#include <immintrin.h>
unsigned int status = _xbegin();
if (status == _XBEGIN_STARTED) {
    // Access shared data ...
    if (problem) // give up ?
        _xabort(0);
    // Access more shared data ...
    _xend();
    /* <----- Success !!! */
} else { /* <--- Failure */
    if (status & _XABORT_EXPLICIT)
        ...
    if (status & _XABORT_CONFLICT)
        ...
    if (status & _XABORT_CAPACITY)
        ...
}
```

- ▶ `_xbegin()` starts a transaction
 - ▶ Returns `_XBEGIN_STARTED`
 - ▶ Flush the cache...
- ▶ `_xend()` attempts to “commit”
 - ▶ OK → execution continues
- ▶ `_xabort(cst)` aborts transaction
- ▶ **In case of failure:**
 - ▶ Returns after `_xbegin()`
 - ▶ Error code (conflict, resources, ...)
- ▶ Still not a panacea
 - ▶ Non-negligible cost
 - ▶ False positives, ...
- ▶ See also TinySTM library