Alan Stern, Rowland Institute at Harvard
Paul E. McKenney, Meta Platforms Kernel Team
Michael Wong, YetiWare Inc.
Maged Michael, Category Labs
Gonzalo Brito, NVIDIA
Kangrejos, Copenhagen, Denmark, September 8, 2024

# Lifetime-End Pointer Zap & How to Avoid OOTA Without Really Trying

# Overview

This is just an overview, not a replacement for the papers themselves

- P2414R10 "Pointer lifetime-end zap proposed solutions"
  - https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2025/p2414r10.pdf
- P3347R5 Invalid/Prospective Pointer Operations
  - https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2025/p3347r5.pdf
  - Based on Davis Herring's P2434R4 "Nondeterministic pointer provenance"
    - https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2025/p2434r4.html
- P3790R1 "Pointer lifetime-end zap proposed solutions: Bag-of-bits pointer class"
  - https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2025/p3790r1.pdf
- P3692R2 "How to Avoid OOTA Without Really Trying"
  - https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2025/p3692r2.pdf

# Overview

- Lifetime-end pointer zap

- Out-of-thin-air (OOTA) cycles

- Where are we on OOTA?

- Leverage restrictions:
  - Real computer systems
  - Speculate properly or not at all
  - Existing restrictions for volatile atomics
  - No invention or repurposing of atomic loads
  - Tooling looks at object code

- Future directions

# Lifetime-End Pointer Zap

# Problem Restatement (C11, 1/2)

```c
struct node_t* _Atomic top;

void list_push(value_t v)
{
  struct node_t *newnode = (struct node_t *) malloc(sizeof(*newnode));
  Struct node_t *next = atomic_load(&top);

  set_value(newnode, v);
  do {
    set_next(newnode, next);
    // newnode's next pointer may have become invalid
  } while (!atomic_compare_exchange_weak(&top, &next, newnode));
}
```

```c
void list_pop_all()
{
  struct node_t *p = atomic_exchange(&top, NULL);

  while (p) {
    struct node_t *next = p->next;

    foo(p);
    free(p);
    p = next;
  }
}
```

# Problem Illustration (C11)
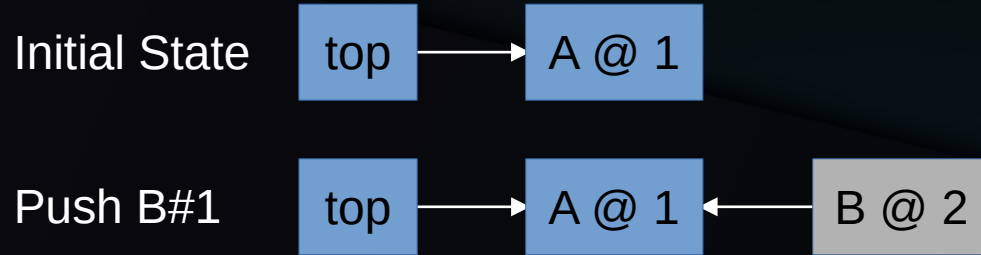
Freelist

Initial State    [ top ] → [ A @ 1 ]
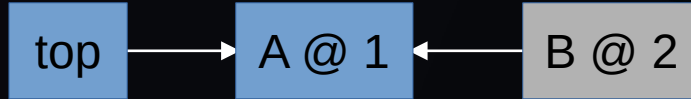
# Problem Illustration (C11)

Freelist

Initial State    | top | → | A @ 1 |

Push B#1    | top | → | A @ 1 | ← | B @ 2 |

# Problem Illustration (C11)

# Problem Illustration (C11)

# Problem Illustration (C11)



Freelist

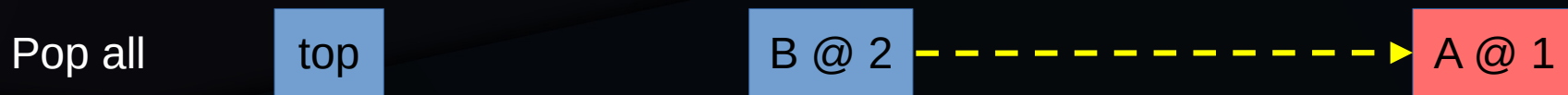| | | |
|---|---|---|
| Initial State | top → | A @ 1 |
| Push B#1 | top → A @ 1 ← | B @ 2 |
| Pop all | top | B @ 2 ⇢ A @ 1 |
| Push C | top → C @ 1 ⇠ | B @ 2 |
| Push B#2 | top → B @ 2 ⇢ | C @ 1 |

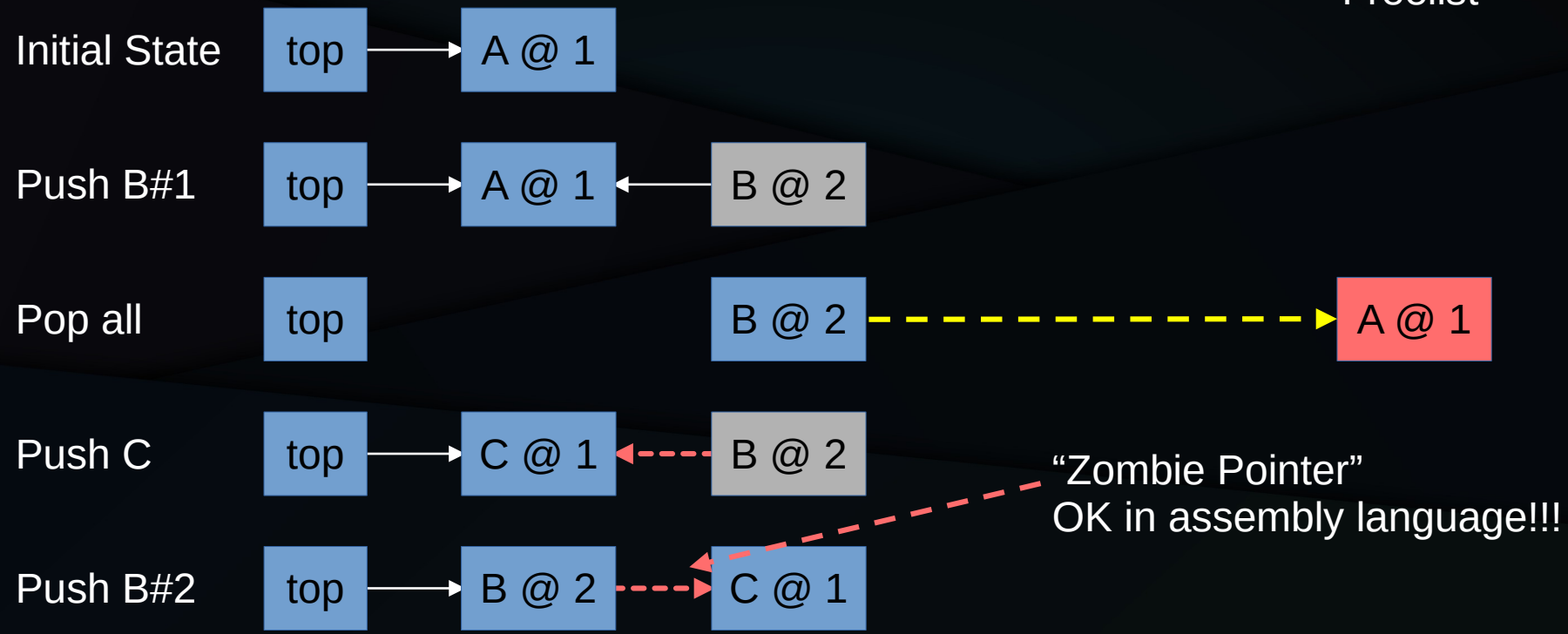# Problem Illustration (C11)

# Problem Illustration (C11)



Freelist

| | | |
|---|---|---|
| Initial State | top → A @ 1 | |
| Push B#1 | top → A @ 1 ← B @ 2 | |
| Pop all | top    B @ 2 ⇢ | A @ 1 |
| Push C | top → C @ 1 ⇠ B @ 2 | |
| Push B#2 | top → B @ 2 ⇢ C @ 1 | |

"Zombie Pointer"
OK in assembly language!!!

# Problem Illustration (C11)

Freelist

| Initial State | top → A @ 1 |
| Push B#1 | top → A @ 1 ← B @ 2 |
| Pop all | top ⇢ A @ 1 |
| Push C | C @ 1 ⇠ B @ 2 |
| B#2 | top → B @ 2 ⇢ C @ 1 |

"Zombie Pointer"
OK in assembly language!!!

LIFO stack with pop-all is ABA tolerant

# This is Real and Isn't Going Away

- LIFO stack described by Treiber in 1986
  - Written in IBM BAL, avoiding issues with compilers
- LIFO stack alluded to in early 1970s
- LIFO stack implemented in Rust library
  - Though with `pop()`, not `pop_all()`.
- Used in heavily production in many languages
  - Often open-coded, often inadvertently reinvented

# OK, OK, What is New Since 2024???

# C and C++: Pointer Provenance

- Pointers contain bits and also "provenance"
  - Compiler may assume that pointers from two different calls to the allocator are unequal
  - Some provenance might be part of pointer value (ARM MTE)
- Provenance may be erased
  - Conversion to integer, I/O, optimization frontiers
- Davis Herring C++ proposal (P2434R4) provides "angelic provenance", **but now limited**

# C++: Angelic Provenance

- Davis Herring P2434R4 ("Nondeterministic pointer provenance") restricts provenance restoration
  - Conversion from integer, I/O, optimization frontiers
  - At which point, the compiler must choose provenance (if any) that allows the program to be well-formed
    - **But compiler need not consider objects where provenance restoration happens-before the beginning of that object's storage duration**

# C++: Angelic Provenance

- Davis Herring P2434R4 ("Nondeterministic pointer provenance") restricts provenance restoration
  - Conversion from integer, I/O, optimization barriers
  - At which point, the compiler can choose provenance (if any) that allows the program to be well-formed
    - **But compiler could not consider objects where provenance restoration happens-before the beginning of that object's storage duration**
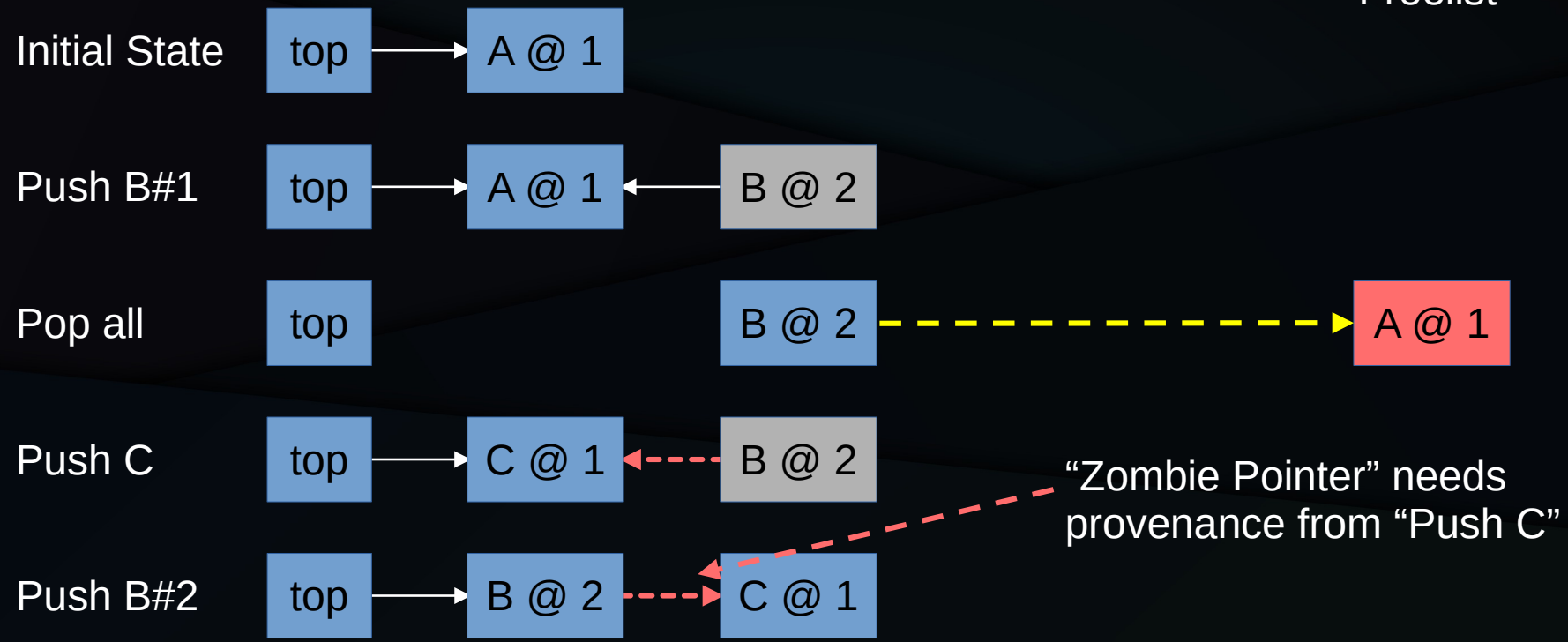
Not enough for LIFO stack…

# Problem Illustration (C11)



Freelist

Initial State — top → A @ 1

Push B#1 — top → A @ 1 ← B @ 2

Pop all — top    B @ 2 ⇢ A @ 1

Push C — top → C @ 1 ⇠ B @ 2

Push B#2 — top → B @ 2 ⇢ C @ 1

"Zombie Pointer" needs provenance from "Push C"

# Problem Illustration (C11)

| | | |
|---|---|---|
| Initial State | top → A @ 1 | |
| Push B#1 | top → A @ 1 ← B @ 2 | |
| Pop all | top | A @ 1 |
| Push C | B @ 2 | |
| Push B#2 | top → B @ 2 ⇢ C @ 1 | |

Freelist

"Zombie Pointer" needs provenance from "Push C"

**No conversion, I/O, or optimization frontier, but useful definitions**

# What Else Is Needed?

- P2414R10 ("Pointer lifetime-end zap proposed solutions"): Provenance restoration results from:
  - Conversions to/from atomic<T *>
    - Including old pointer referenced by successful CAS operations
  - Volatile accesses involving pointers
- P3347R5 ("Pointer lifetime-end zap proposed solutions: Tighten IDB for invalid and prospective pointers")
  - Glvalue-to-rvalue conversions from invalid pointers must produce representation values consistent with those of the lvalue
- P3790R1 ("Pointer lifetime-end zap proposed solutions: bag-of-bits pointer class"): Provenance restoration results from:
  - `ptr_bits<T>` **(But now internal representation not visible to user per IBM System i)**
  - `launder_ptr_bits()` "identity" function

# What Else Is Needed?

- P2414R10 ("Pointer lifetime-end zap proposed solutions"): Provenance restoration results from:
  - Conversions to/from atomic<T *>
    - Including old pointer referenced by successful CAS operations
  - Volatile accesses involving pointers
- P3347R5 ("Pointer lifetime-end zap prop̶o̶s̶e̶d̶ ̶s̶o̶l̶u̶t̶i̶o̶n̶s̶: Tighten IDB for invalid and prospective pointers")
  - Glvalue-to-rvalue conversi̶o̶n̶ ̶̶̶u̶ pointers must produce representation values consistent with tho̶s̶e̶

~~Not all that much!!!~~

- P3790R1 ("Pointer li̶f̶etime-end zap proposed solutions: bag-of-bits pointer class"): Provenance restoration results from:
  - `ptr_bits<T>` **(But now internal representation not visible to user per IBM System i)**
  - `launder_ptr_bits()` "identity" function

# Status in C++ Committee

- All progressing through C++ committee:
  - P2414R10 "Pointer lifetime-end zap proposed solutions"
  - P3347R5 "Pointer lifetime-end zap proposed solutions: Tighten IDB for invalid and prospective pointers"
  - P3790R1 "Pointer lifetime-end zap proposed solutions: bag-of-bits pointer class "
  - Davis Herring's P2434R4 "Nondeterministic pointer provenance"
- No guarantees, but best progress thus far

# Pointer-Zap Discussion

# OOTA Cycles

# Proposed Change to C++ Standard

# Proposed Change to C++ Standard

- P3692R2 ("How to Avoid OOTA Without Really Trying"):
  - After N5008 33.5.4p8 ([atomics.order])33.5.4p8 ([atomics.order]):
    - "Implementations should ensure that no "out-of-thin-air" values are computed that circularly depend on their own computation. [ Note 6 *… example …* ]"
  - Add the following:
    - "Compiler-based implementations whose binaries run on conventional hardware are guaranteed not to compute out-of-thin-air values in programs that are free of undefined behavior, as long as they restrict themselves to thread- at-a-time analysis and and treat non-volatile atomic accesses as if they were volatile, except that, when permitted by the as-if rule, they may omit accesses, merge accesses to the same object, or reorder accesses to different objects."
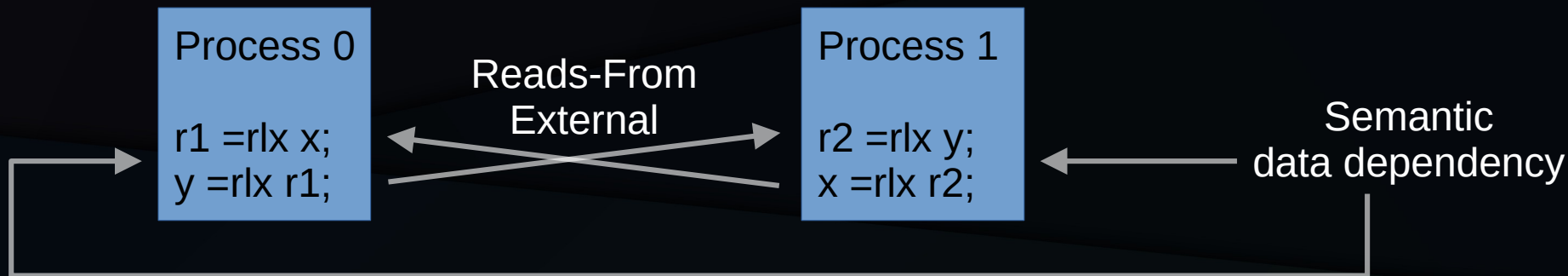
# Proposed Change to C++ Standard

- P3692R2 ("How to Avoid OOTA Without Really T̶r̶y̶i̶n̶g̶"):
  - After N5008 33.5.4p8 ([atomics.order])33.5̶.̶4̶p̶8̶ ̶(̶[̶a̶t̶o̶m̶i̶c̶s̶.̶o̶r̶der]):
    - "Implementations should ensure that n̶o̶ ̶v̶a̶l̶u̶es are computed that circularly depend on their own co̶m̶p̶u̶t̶a̶t̶i̶o̶n̶.̶ ̶[̶ example … ]"
  - Add the following:
    - "Compiler-bas̶e̶d̶ ̶i̶m̶p̶l̶e̶m̶e̶n̶t̶a̶t̶i̶ons run on conventional hardware are guarantee̶d̶ ̶t̶o̶ ̶a̶v̶o̶i̶d̶ ̶O̶O̶T̶A̶ ̶v̶a̶lues in programs that are free of und̶e̶f̶i̶n̶e̶d̶ ̶b̶e̶h̶a̶v̶i̶o̶r̶ ̶a̶s̶ ̶l̶o̶ng as t̶h̶e̶y̶ ̶restrict themselves to thread- at-a-time analysis and̶ ̶t̶r̶e̶a̶t̶ ̶n̶o̶n̶-̶volatile atomic accesses as if they were volatile, except that, when p̶e̶r̶m̶i̶t̶t̶ed by the as-if rule, they may omit accesses, merge accesses to the same object, or reorder accesses to different objects."

# OOTA Cycles: Background

# OOTA Cycles

- Self-satisfying load-buffering cycle, x==y==42

Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Process 1

r2 =rlx y;
x =rlx r2;

Semantic
data dependency

Relaxed load or store denoted by "=rlx"

# OOTA Cycles

- Self-satisfying load-buffering cycle, x==y==42

Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Proces...

...lx y;
x =rlx r2;

Semantic
data dependency

**Why "external" in reads-from external?**

```
r1 =rlx X;
Y =rlx r1;
r2 =rlx Y;
Z =rlx r2;
```

rfi

```
r1 =rlx X;              r1 =rlx X;

Y =rlx r1;              Z =rlx r1;

r2 =rlx Y;              Y =rlx r1;

Z =rlx r2;              r2 = r1;
```
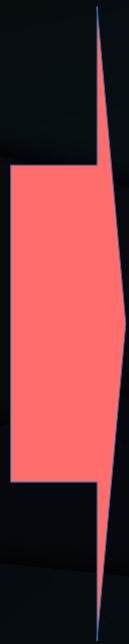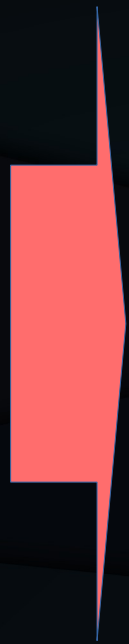
# OOTA Cycles: Reads-From Internal

```
r1 =rlx X;
Y =rlx r1;
r2 =rlx Y;
Z =rlx r2;
```

```
r1 =rlx X;
Z =rlx r1;
Y =rlx r1;
r2 = r1;
```

Compiler eliminated the read from Y so that
the store to Z can now occur before the store to Y

```
r1 =rlx X;              r1 =rlx X;
Y =rlx r1;              Z =rlx
r2 =rlx Y;              Y    lx r1;
Z =rlx r2;              r2 = r1;
```

...piler eliminated the read from Y so that
...e store to Z can now occur before the store to Y

**Hence "external" in reads-from external**

See Section 2.1 ("OOTA: rf vs. rfe") of P3692R2

# OOTA Cycles, Original Diagram

- Self-satisfying load-buffering cycle, x==y==42



Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Process 1

r2 =rlx y;
x =rlx r2;

Semantic
data dependency

# OOTA Cycles, Original Diagram

- Self-satisfying load-buffering cycle, x==y==42

Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Process ...

... y;
x =rlx r2;

Semantic
data dependency

Not seen in "real life"

# OOTA Cycles, Original Diagram

- Self-satisfying load-buffering cycle, x==y==42

Process 0

Reads-From
External

r1 =rlx x;
y =rlx r1;

Process 1

r2 =rlx y;
x =rlx r2;

Semantic
data dependency

**Not seen in "real life"**

**Why???**

# Where Are We on OOTA?

# Where Are We on OOTA? (TL;DR)

Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Process 1

r2 =rlx y;
x =rlx r2;

# Where Are We on OOTA? (TL;DR)



Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Process 1

r2 =rlx y;
x =rlx r2;

Reads take time

# Where Are We on OOTA? (TL;DR)



Instruction execution takes time

Process 0

r1 =rlx x;
y =rlx r1;

Reads-From External

Process 1

r2 =rlx y;
x =rlx r2;

Reads take time

Instruction execution takes time

# Where Are We on OOTA? (TL;DR)

Instruction execution takes time

**Process 0**

r1 =rlx x;
y =rlx r1;

**Reads-From External**

**Process 1**

r2 =rlx y;
x =rlx r2;

Instruction execution takes time

**Reads take time**

To form an OOTA cycle, at least one step must go backwards in time!!!

# Where Are We on OOTA? (TL;DR)

Instruction execution takes time

Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
Exter...

...rlx y;
x =rlx r2;

Instruction execution takes time

...take time

**OOTA cycle cannot form (on real compiler-based systems)**

To form an OOTA cycle, at least one step must go backwards in time!!!

# Where Are We on OOTA?

- Generalized "OOTA Cycle" (Section 2.2.2)

- Fundamental property of semantic dependency (Sections 5.3 and 6.1)

- Demonstrate OOTA-freedom under restrictions (Sections 6.2-6.4 for demonstration, 4.4 for restrictions)
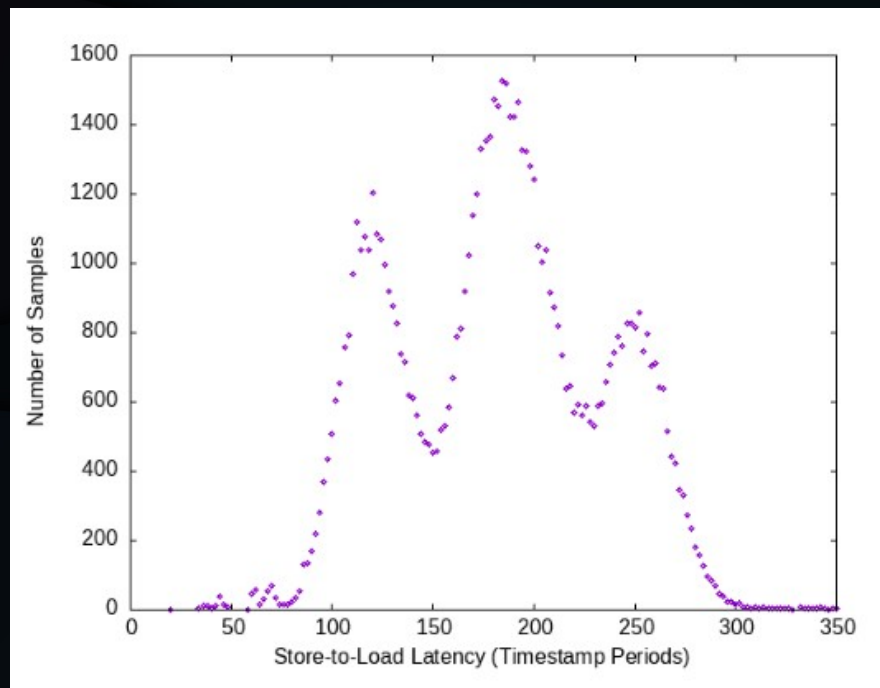
# Leverage Restrictions

# Real Computer Systems

# Real Computer Systems: Store-to-Load

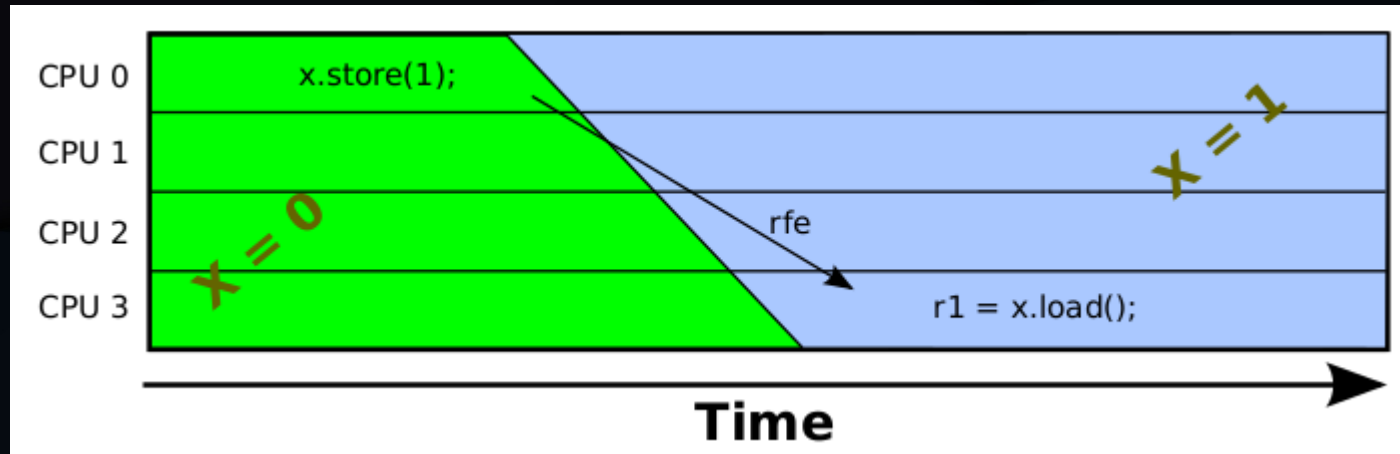- Store-to-load links are temporal*

\* The event that is logically first must happen before the other event in real-world time
Dual-socket Intel(R) Xeon(R) Gold 6138 CPUs @ 2.00 GHz, 80 hardware threads total: Measure beginning of store to end of load

# Real Computer Systems: Store-to-Load

- Store-to-load links are temporal: HW view

# Real Computer Systems: Store-to-Store

- Store-to-store links are atemporal*

* The event which is logically first can happen after the other event in real-world time
Dual-socket Intel(R) Xeon(R) Gold 6138 CPUs @ 2.00 GHz, 80 hardware threads total: Measure beginning of winning store to end of store
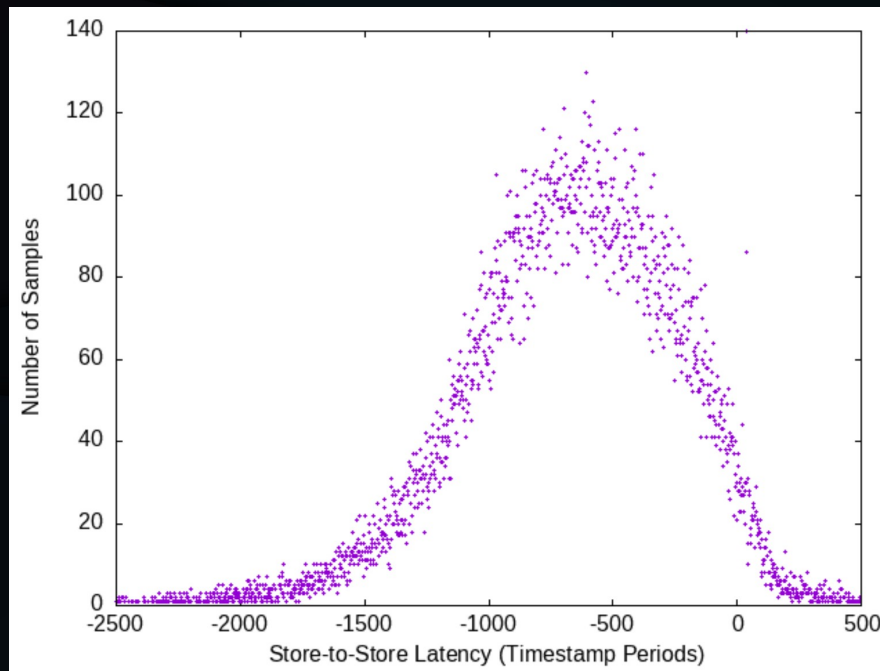
# Real Computer Systems: Store-to-Store

- Store-to-store links are atemporal: HW view

"co" is "modification order" in the C++ memory model

# Real Computer Systems: Load-to-Store

- Load-to-store links are atemporal

Dual-socket Intel(R) Xeon(R) Gold 6138 CPUs @ 2.00 GHz, 80 hardware threads total: Measure beginning of load to end of store

# Real Computer Systems: Load-to-Store

- Load-to-store links are atemporal: HW view

"fr" is "from reads", which connects a read to a write that happened too late to affect the value loaded

# Real Computer Systems: Summary

- Load-to-store links: Atemporal

- Store-to-store links: Atemporal

- Store-to-load links: Temporal
  - And thus have ordering properties on the cheap

See Appendix A ("Interthread Communications") of wg21.link/P3064R2 ("How to Avoid OOTA Without Really Trying")

# Speculate Properly or Not At All

# Speculate Properly or Not At All

```
r1 =speculate$_X$ 2;

r2 = somefunc(r1);

Y = r2;
```

atemporal!!!

```
X =rlx 1;
```

# Speculate Properly or Not At All

$$r1 = speculate_X\ 2;$$

$$r2 = somefunc(r1);$$

$$Y = r2;$$

*Also improper!!!*

~~ate____al!!!~~

$$X =rlx\ 1;$$

Don't just guess!  Guess and then check!!!

# Speculate Properly or Not At All

```
r1 =speculateₓ 2;
r2 = somefunc(r1);
```
~~Y = r2;~~
```
r3 =rlx X; // 1, not 2!
if (r1 != r3)
    r2 = somefunc(r3);
Y = r2;
```

X =rlx 1; —— **temporal!!!** ——→

# Speculate Properly or Not At All

```
r1 =speculateₓ 2;
r2 = somefun       );
Y = r2;
r2            1, not 2!
       r3)
       r2 = somefunc(r3);
Y = r2;
```

X =rlx 1;  ——— temporal!!! ———▶

**Speculation must be checked against the value from an actual load!!!**

See Section 1.1 ("Brief OOTA Overview") and Section 5.2 ("Instruction Ordering") of P3692R2

# Existing Restrictions on Volatile Atomics

# Existing Restrictions on Volatile Atomics

- Compiler may not:
  - Reorder accesses
  - Invent, duplicate, or repurpose accesses
  - Merge or fuse accesses
  - Omit accesses
- Relax restrictions for non-volatile atomics?

# No Atomic-Load Invention/Repurposing

# No Atomic-Load Invention

- Guaranteed perfect square for small X:

```
int r0 =rlx x;
int r1 = r0 * r0 + 2 * r0 + 1;
```

- But not if atomic loads are invented!!!

```
int r0 =rlx x;
int invented =rlx x;
int r1 = r0 * r0 + 2 * invented + 1;
```

# No Atomic-Load Invention

- Guaranteed no future square for small X:

```
int r0 = x x;
int r = * r0 + 2 r0 + 1;
```

- But not atomic ds are invnted!!!

```
int r =rlx x;
int inted =rlx
int r1 = * invented + 1;
```

# No Atomic-Load Repurposing

- Guaranteed perfect square for small X:

```
r2 =rlx x;
do_something(r2); // No synchronization or stores to x
int r0 =rlx x;
int r1 = r0 * r0 + 2 * r0 + 1;
```

- But not if atomic loads are repurposed!!!

```
r2 =rlx x;
do_something(r2); // No synchronization or stores to x
int r0 =rlx x;
int r1 = r0 * r0 + 2 * r2 + 1;
```

# No Atomic-Load Repurposing

- Guaranteed perf... ...re for small X:

```
r2 =r1...
do_something(r2); // No ...chronization or stores to x
in...0 ... x;
i... r1 = r... r0 + 2 * r0 ...;
```

- But... if atomic lo... ...re repurp...d!!!

```
r... ...lx x;
do... ...ething(r2); // ...chronization or stores to x
int ...lx x;
int r1 ... ... r2 + 1;
```

# Instead, Merge the Atomic Loads

- Guaranteed perfect square for small X:

```
r2 =rlx x;
do_something(r2); // No synchronization or stores to x
int r0 =rlx x;
int r1 = r0 * r0 + 2 * r0 + 1;
```

- And that guarantee is maintained for merged loads:

```
r0 =rlx x;
do_something(r0); // No synchronization or stores to x
int r1 = r0 * r0 + 2 * r0 + 1;
```

# Instead, Merge the Atomic Loads

- Guaranteed perfect square for small X:

  ```
  r2 =rlx x;
  do_something(r2); // No synchronizati          to x
  int r0 =rlx x;
  int r1 = r0 * r0 + 2 * r0
  ```

- And that guarantee is n            ged loads:

  ```
  r0 =rlx x;
  do_so              No synchronization or stores to x
             0 + 2 * r0 + 1;
  ```

**If do_something() contains synchronization, then must keep both atomic loads**

# Atomic Loads and Memory Ordering

```
r1 =rlx X;                          X =rlx 1;
          sdep?
r2 =rlx Y;

Z =rlx (r1 == r2);
```

Note: X, Y, and Z boolean and initially zero

See Section 4 ("C++ Compilers")of P3692R2

# Atomic Loads and Memory Ordering
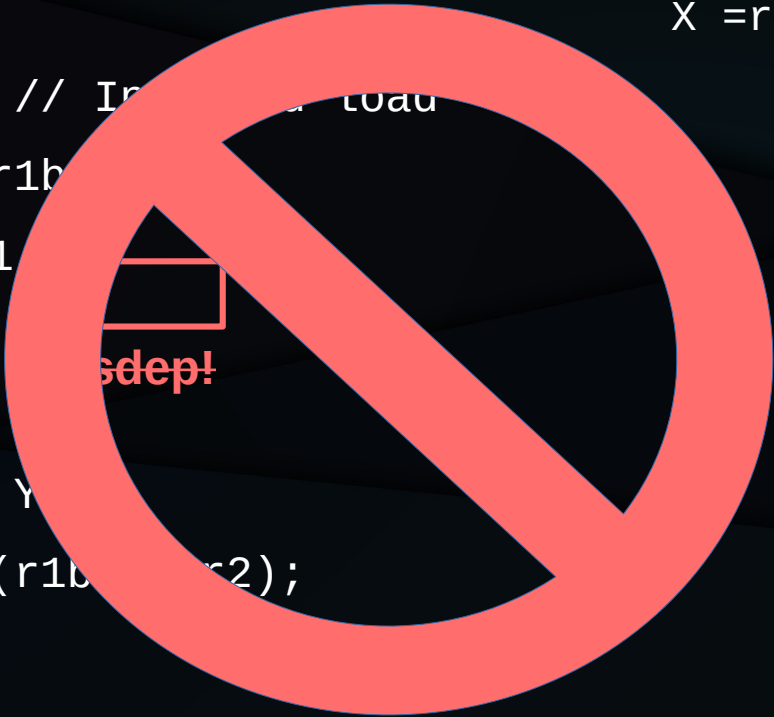
```
r1a =rlx X;

r1b =rlx X; // Invented load

If (r1a != r1b) {

    Z =rlx 1;

    r2 =rlx Y;

} else {

    r2 =rlx Y;

    Z =rlx (r1b == r2);

}
```

**sdep!**

```
X =rlx 1;
```

Note: X, Y, and Z boolean and initially zero

Inventing atomic load likely also invents hundreds-of-cycles cache miss!!!

# Atomic Loads and Memory Ordering

```
r1a =rlx X;                              X =rlx 1;

r1b =rlx X; // I      load

If (r1a != r1b

    Z =rlx 1

    r2 =rlx          sdep!

} else {

    r2 =rlx Y

    Z =rlx (r1b    r2);

}                                 Note: X, Y, and Z boolean and initially zero
```

See Section 4.3 ("Inventing Atomic Loads Can Cause Errors and Destroy Semantic Dependencies") of P3692R2
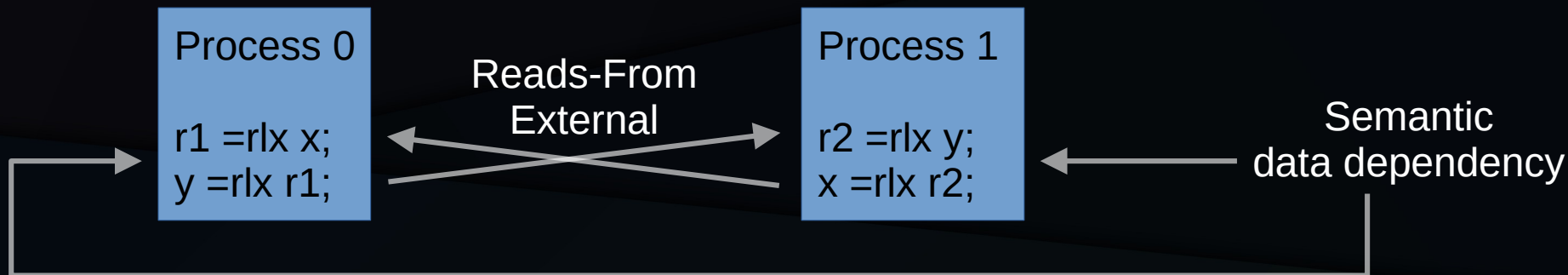
# Non-Volatile Atomics Optimizations?

- Looking only at relaxed operations:
  - **Reorder loads/stores from/to different objects**
  - **Merge back-to-back loads to same object**
  - **Drop loads whose values are unused**
  - **Discard first of back-to-back stores to same object**
  - **Fuse loads from (or stores to) adjacent objects if this results in a machine-word-sized/aligned access**
  - **But no invented, duplicated, or repurposed loads!!!**

# Tooling Looks at Object Code

See Section 7.3 ("Semantic Dependencies and Tooling") and Appendix C ("But What About Tooling?") of wg21.link/P3064R2

# OOTA Cycles, Original Diagram

- Self-satisfying load-buffering cycle, x==y==42



Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Process 1

r2 =rlx y;
x =rlx r2;

Semantic
data dependency

# OOTA Cycles, Original Diagram

- Self-satisfying load-buffering cycle, x==y==42

**Temporal**

Process 0

r1 =rlx x;
y =rlx r1;

Reads-From
External

Process 1

r2 =rlx y;
x =rlx r2;

Semantic
data dependency

**Temporal?
Semantic dependency?**

- Self-satisfying load-buffering cycle ........ =42

**Temporal**

Process 0

r1 =rlx x;
y =rlx

Re...

...x r2;

Semantic
data dependency

**Temporal?
Semantic dependency?**

*If each step in an OOTA cycle is temporal, then that cycle cannot happen in the real world because no step could happen first!*

# Semantic Dependencies are Tricky

- At source-code level, semantic dependencies:
    - Are not strict functions of source code (Section 2.2.1)
    - Can be many-to-one (Section 2.2.2)
    - Depend on partially defined executions (Section 2.2.8-9)
    - Depend on compilers and their users (Section 2.2.8 & 4.1)
- Current paper assumes local analysis (no global cross-thread optimizations)

# Semantic Dependencies in Code?

- Semantic dependencies are temporal:
  - Instructions take time to execute
  - Speculation must be checked against actual load

# Semantic Dependencies in Code?

- Semantic dependencies are temporal:
  - Instructions take time to execute
  - Speculation must be checked against actual load
- Compiler optimizations break dependencies:
  - But HW memory models respect dependencies
  - Thus look at object code (seL4 verification approach)
  - Also look at other compiler-produced artifacts

See Sections 4-6, P3692R2

# Semantic Dependencies in Code?

- Semantic dependencies are temporal:
  - Instructions take time to execute
  - Speculation must be chec~~ked~~
- Compiler optimiz~~ations~~ ~~depen~~dencies:
  - But HW ~~reo~~ ~~rdering~~ dependencies
  - T~~he~~ ~~se~~ (seL4 verification approach)
  - ~~consider~~ compiler-produced artifacts

**If compiler optimizes dependency away, it was not semantic. Otherwise, executing dependency's code will take time.**

See Sections 4-6, P3692R2

# Where Are We on OOTA? (Reprise)

- Generalized "OOTA Cycle" (Section 2.2.2)

- Fundamental property of semantic dependency (Sections 5.3 and 6.1)

- Demonstrate OOTA-freedom under restrictions (Sections 6.2 and 6.3 for demonstration, 4.4 for restrictions)
  - The main restriction is: No invented, duplicated, or repurposed atomic loads

# Future Directions

- From compilers to (some) JITs, interpreters, and link-time optimizations (LTO)

- Compilers doing (some) global analysis given volatile atomics

- Identify absolute semantic dependencies inherent in source code

- Non-shared-memory communication

# Discussion