# Memory Consistency

## You can't ignore it, you should fake it

Rekai González Alberquilla
*AMD UK*

**16th May 2024**

**AMD**
together we advance_

# Outline

AMD
together we advance_

# Conventions
**And disclaimers, and such**

- `Bold monospace` is used for code or linux commands.
- The pairs *load instruction* - *read operation*, *store instruction* - *write operation* are used interchangeably throughout the presentation.

**AMD**
together we advance_

# Recomended bibliography

- Bruce Jacob, *"The Memory System. You can't avoid it, You can't ignore it, you can't fake it"*.
- Daniel Sorin, Mark Hill, David Wood, *"A Primer on Memory Consistency and Cache Coherence"*.

Both on Springer Synthesis Lectures on Computer Architecture: Affordable, probably available in your library, and offered in open-access every so often.

- Sarita Adve, Kourosh Garachorloo, *"Shared Memory Consistency Models: A Tutorial"*. IEEE Computer'96, Vol. 29, Issue 12.
- Shen and Lipasti, *"Modern Processor Design: Fundamentals of Superscalar Processors"*. McGraw-Hill, 2003.

AMD
together we advance_

# Preamble

This talk assumes

- You are familiar with Cache memories.
- You are familiar with Cache coherence.
- You are familiar with Out-of-order execution and superscalar processors.
  - Branch prediction is orthogonal.
- All data in C/C++ is `volatile`-qualified.
- All memory reads/writes are 8-bytes long, aligned.

# Preamble

This talk assumes

- You are familiar with Cache memories.
- You are familiar with Cache coherence.
- You are familiar with Out-of-order execution and superscalar processors.
  - Branch prediction is orthogonal.
- All data in C/C++ is `volatile`-qualified.
- All memory reads/writes are 8-bytes long, aligned.

Consistency vs Coherence, often put together, two totally different beasts.

- Coherence
  - Transparent to the programmer.
  - Only HW designers and low level programmers care about it.
- Consistency
  - Part of the architecture (implication on instruction semantics).
  - Has an impact beyond HW and low level SW.

# Definitions

- Parallel: Two events, A, and B, are parallel iff there is no order relationship between them.
- Static instruction: A sequence of bytes that adheres an element of the ISA. Defined by the value of the bytes.
- Dynamic instruction: An instance of a static instruction, defined by the static instruction as well as the execution context.

AMD
together we advance_

# Definitions

- Parallel: Two events, A, and B, are parallel iff there is no order relationship between them.
- Static instruction: A sequence of bytes that adheres an element of the ISA. Defined by the value of the bytes.
- Dynamic instruction: An instance of a static instruction, defined by the static instruction as well as the execution context.

```
1  b: b.eq r0, r1, end
2     fadd [r3], d0
3     add r3, #8
4     inc r0
5     jmp b
6  end:
```

```
1    b.eq r0, r1, end # r0 = 3, r1 = 5
2    fadd [r3], d0    # r3 = 0x10018, d0 = 2.0
3    add r3, #8       # r3 = 0x10018 | r3 <- 0x10020
4    inc r0           # r0 = 3       | r0 <- 4
5    jmp b
6    b.eq r0, r1, end # r0 = 4, r1 = 5
7    fadd [r3], d0    # r3 = 0x10020, d0 = 2.0
8    add r3, #8       # r3 = 0x10020 | r3 <- 0x10028
9    inc r0           # r0 = 4       | r0 <- 5
10   jmp b
11   ...
```

AMD
together we advance_

# Definitions

- Parallel: Two events, A, and B, are parallel iff there is no order relationship between them.
- Static instruction: A sequence of bytes that adheres an element of the ISA. Defined by the value of the bytes.
- Dynamic instruction: An instance of a static instruction, defined by the static instruction as well as the execution context.

- Speculative: means executing without waiting for **memory order** guarantee, as opposed to **control order** guarantee.
- $L(A)$: Load/Read `Mem[A]`.
- $S(B)$: Store/Write `Mem[B]`.
- $A \neq B \Rightarrow Mem[A] \cap Mem[B] = \Phi$.

AMD
together we advance_

# Definitions:Order relationships
**All about not being one and only one order**

- Program order: Is the order relationship between **dynamic instructions** implied by the trace (sequential execution) of the program, $<_p$.
- Coherence order: Is the order relationship between **dynamic instructions** implied by the cache coherence protocol.
- Memory order: Is the order in which (memory) operations arrive to the shared global memory, $<_m$.

The rest of the talk is about when, given two instructions $I_i$, $I_j$, $I_i <_p I_j \Rightarrow I_i <_m I_j$.

AMD
together we advance_

# Definitions:Order relationships
**All about not being one and only one order**

- Program order: Is the order relationship between **dynamic instructions** implied by the trace (sequential execution) of the program, $<_p$.
- Coherence order: Is the order relationshipt between **dynamic instructions** implied by the cache coherence protocol.
- Memory order: Is the order in which (memory) operations arrive to the shared global memory, $<_m$.

The rest of the talk is about when, given two instructions $I_i$, $I_j$, $I_i <_p I_j \Rightarrow I_i <_m I_j$.

- Fence: the ultimate tool. A Fence is a special instruction such that: for any two instructions $I_i$, $I_j$, for any two fences $F_u$, $F_v$
    - $I_i <_p F_u \Rightarrow I_i <_m F_u$
    - $F_u <_p I_j \Rightarrow F_u <_m I_j$
    - $F_u <_p F_v \Rightarrow F_u <_m F_v$

*i.e.*, fences enforce order.

- There can be fences that provide only a subset of the guarantees.

# Motivating example
## Why do we want a consistency model?

Assuming `P0` and `P1` are processors in a shared memory system.

```
Mem[X] := 0
Mem[Y] := 0
```

```
    [P0]
St [X], 0x1 # S(X)
St [Y], 0x1 # S(Y)
```

```
    [P1]
Ld r0, [Y] # L(Y)
Ld r1, [X] # L(X)
```

Which of the following are possible outcomes?

Knowing what is allowed and what is not is needed to validate HW and SW.

AMD
together we advance_

# Motivating example
## Why do we want a consistency model?

Assuming `P0` and `P1` are processors in a shared memory system.

```
Mem[X] := 0
Mem[Y] := 0
```

```
[P0]
St [X], 0x1 # S(X)
St [Y], 0x1 # S(Y)
```

```
[P1]
Ld r0, [Y] # L(Y)
Ld r1, [X] # L(X)
```

Which of the following are possible outcomes?

- `[P1] r0 = 1, r1 = 1`

Knowing what is allowed and what is not is needed to validate HW and SW.

AMD
together we advance_

# Motivating example
## Why do we want a consistency model?

Assuming `P0` and `P1` are processors in a shared memory system.

```
Mem[X] := 0
Mem[Y] := 0
```

```
     [P0]
St [X], 0x1 # S(X)
St [Y], 0x1 # S(Y)
```

```
     [P1]
Ld r0, [Y] # L(Y)
Ld r1, [X] # L(X)
```

Which of the following are possible outcomes?

- `[P1] r0 = 1, r1 = 1`
- `[P1] r0 = 0, r1 = 1`

Knowing what is allowed and what is not is needed to validate HW and SW.

AMD
together we advance_

# Motivating example
## Why do we want a consistency model?

Assuming `P0` and `P1` are processors in a shared memory system.

```
Mem[X] := 0
Mem[Y] := 0
```

```
[P0]
St [X], 0x1 # S(X)
St [Y], 0x1 # S(Y)
```

```
[P1]
Ld r0, [Y] # L(Y)
Ld r1, [X] # L(X)
```

Which of the following are possible outcomes?

- `[P1] r0 = 1, r1 = 1`
- `[P1] r0 = 0, r1 = 1`
- `[P1] r0 = 1, r1 = 0`

Knowing what is allowed and what is not is needed to validate HW and SW.

AMD
together we advance_

# Motivating example
## Why do we want a consistency model?

Assuming `P0` and `P1` are processors in a shared memory system.

```
Mem[X] := 0
Mem[Y] := 0
```

```
[P0]
St [X], 0x1 # S(X)
St [Y], 0x1 # S(Y)
```

```
[P1]
Ld r0, [Y] # L(Y)
Ld r1, [X] # L(X)
```

Which of the following are possible outcomes?

- `[P1] r0 = 1, r1 = 1`
- `[P1] r0 = 0, r1 = 1`
- `[P1] r0 = 1, r1 = 0`
- `[P1] r0 = 0, r1 = 0`

Knowing what is allowed and what is not is needed to validate HW and SW.

AMD
together we advance_

# What is memory Consistency?

Memory consistency is a contract between a programmer and a system that describes which outcomes are possible for a particular program.

- **Strict consistency**: writes to memory are instantaneously observable by all processors in the system. *Purely theoretical*.
- **Sequential consistency** (Lamport, 1979): *"The result of any execution is the same as if the (read and write) operations of all processes on the data store were executed in some sequential order,and the operations of each individual processor appear in this sequence in the order specified by its program"* (MIPS R10000).
- **Release consistency**
  - **TSO**, strong: Younger loads may be reordered w.r.t. older stores. *I.e.* writes can be buffered locally (x86).
  - **RMO**, weak: A read or write may be reordered w.r.t. any other read or write to a different *location* (Power, Arm).

AMD
together we advance_

# What does that mean?
## The informal view

```
Mem[X] := 0
Mem[Y] := 0
Mem[Z] := 0
```

```
        [P0]
St [X], 0x1
St [Y], 0x1
Ld r0, [Z]
```

```
        [P1]
St[Z], 0x1
Ldr0, [Y]
Ldr1, [X]
```

| P0.r0 | P1.r0 | P1.r1 | SC | TSO | RMO |
|-------|-------|-------|----|----|-----|
| 0 | 0 | 0 | N | Y | Y |
| 0 | 0 | 1 | N | Y | Y |
| 0 | 1 | 0 | N | N | Y |
| 0 | 1 | 1 | Y | Y | Y |
| 1 | 0 | 0 | Y | Y | Y |
| 1 | 0 | 1 | N | Y | Y |
| 1 | 1 | 0 | N | N | Y |
| 1 | 1 | 1 | Y | Y | Y |

# Push towards stricter consistency
**Or why did you say we care about this?**

```cpp
 1  class Peterson {
 2      bool _request[2]{false, false};
 3      int _yields{0};
 4  public:
 5      void lock(int me) {
 6          _request[me] = true;
 7          _yields = me;
 8          while (_request[1 - me] && _yields == me)
 9              continue;
10      }
11      void unlock(int me) {
12          _request[me] = false;
13      }
14  };
```

AMD
together we advance_

# Push towards stricter consistency
**Or why did you say we care about this?**

```cpp
1   class Peterson {
2       bool _request[2]{false, false};
3       int _yields{0};
4   public:
5       void lock(int me) {
6           _request[me] = true;
7           _yields = me;
8           while (_request[1 - me] && _yields == me)
9               continue;
10      }
11      void unlock(int me) {
12          _request[me] = false;
13      }
14  };
```

This algorithm is SC-sound but does not work in TSO or RMO

- Thinking SC is easier.

AMD
together we advance_

# The (a bit more) formal view
## What is the TSO guarantee?

No matter whether A and B are the same

- $L(A) <_p L(B) \Rightarrow L(A) <_m L(B)$: Loads cannot overtake loads.
- $L(A) <_p S(B) \Rightarrow L(A) <_m S(B)$: Stores cannot overtake loads.
- $S(A) <_p S(B) \Rightarrow S(A) <_m S(B)$: Stores cannot overtake stores.

AMD
together we advance_

# The (a bit more) formal view
## What is the TSO guarantee?

No matter whether A and B are the same

- $L(A) <_p L(B) \Rightarrow L(A) <_m L(B)$: Loads cannot overtake loads.
- $L(A) <_p S(B) \Rightarrow L(A) <_m S(B)$: Stores cannot overtake loads.
- $S(A) <_p S(B) \Rightarrow S(A) <_m S(B)$: Stores cannot overtake stores.

- $L(A)$ yields the value written by

$$\max_{<_m} \{S(A)|S(A) <_m L(A) \text{ or } S(A) <_p L(A)\}$$

Loads take the value of
- If the $<_p$-latest store before the Load is $L(A) <_m S(A)$, then S(A) (Store-to-load-forwarding).
- Otherwise he $<_m$-latest store before the Load.

AMD
together we advance_

# Back to the example

```
I1: W(_request[0], 1)
I2: W(_yields, 0)
I3: R(_request[1])
```

```
J1: W(_request[1], 1)
J2: W(_yields, 1)
J3: R(_request[0])
```

$I1 <_p I2 <_p I3$

$J1 <_p J2 <_p J3$

According to TSO

- $J3 <_m J1$, and
- $I3 <_m I1$ are valid outcomes

so $I3 <_m J3 <_m I1 <_m I2 <_m J1 <_m J2$ is a valid outcome.

AMD
together we advance_

# What about the good old flag based synchronisation?

```
1  class Channel {
2      bool _ready{0};
3      int _data;
4  public:
5      int consume() {
6          while (!_ready) {}
7          int d = _data;
8          _ready = false;
9      }
10     void produce(int v) {
11         while (_ready) {}
12         _data = v;
13         _ready = true;
14     }
15 };
```

```
C1: R0 <- R(_ready)
C2: BIfZero R0, C1
C3: R1 <- R(_data)
C4: W(_ready, 0)
```

```
P1: R0 <- R(_ready)
P2: BIfNotZero R0, P1
P3: W(_data, R1)
P4: W(_ready, 1)
```

AMD
together we advance_

# Flags in TSO
**Why does it work?**

```
C1: R0 <- R(_ready)
C2: BIfZero R0, C1
C3: R1 <- R(_data)
C4: W(_ready, 0)
```

```
P1: R0 <- R(_ready)
P2: BIfNotZero R0, P1
P3: W(_data, R1)
P4: W(_ready, 1)
```

- $C1_n <_p C3 \Rightarrow C1_n <_m C3$
- $C3 <_p C4 \Rightarrow C3 <_m C4$
- $P1_n <_p P3 \Rightarrow P1_n <_m P3$
- $P3 <_p P4 \Rightarrow P3 <_m P4$

Putting all together **+<2>** If $C1_i$ observes P4, *i.e.* $P4 <_m C1_i$, then it is guaranteed by TSO that $P3 <_m C3$, *i.e.* _data gets the expected v. **+<3>** If $P1_i$ observes C4, *i.e.* $C4 <_m P1_i$, then it is guaranteed by TSO that $C3 <_m P3$, *i.e.* _data is not overwritten before being read.

AMD
together we advance_

# High performance techniques

- We want to do as much Out-of-order as needed to maximise ILP/MLP extraction.
- Break name dependencies, honour real dependencies.
- Execute instructions as soon as operands are ready. Commit in order to preserve precise exceptions.

What does consistency have to do with all this?

- The moment there are more than one threads sharing the memory out-of-order execution needs to be done carefully.

Four techniques, as described in [Shen and Lipasti] and [Sorin *et al*]

- Load-load bypassing: Executing a younger load ahead of an older load waiting.
- Load-store bypassing: Executing a younger load ahead of an older store waiting.
- Store-Load Forwarding: Sending the data from an older store to a younger load.
- Store coalescing: Merging writes before sending them to memory.

AMD
together we advance_

# Some code to reason on
## Who does not love I-DAXPY?

```c
1  void idaxpy(double** x, double** y, double** z, double A) {
2      for (int i = 0; i < N; ++i) {
3          *(z[i]) = A * *(x[i]) + *(y[i]);
4      }
5  }
```

```
Loop:
I1   Rx  <= Ld  [Rpx]      ; x[i]
I2   F2  <= Ld  [Rx]       ; *x[i]
I3   F2  <= Mul F2, F0     ; * A
I4   Ry  <= Ld  [Rpy]      ; y[i]
I5   F4  <= Ld  [Ry]       ; *y[i]
I6   F4  <= Add F4, F0     ; + *y[i]
I7   Rz  <= Ld  [Rpz]      ; z[i]
I8   St  F4, [Rz]          ; *z[i] =
I9   Rpx <= Add Rpx, 8     ; ++i
I10  Rpy <= Add Rpy, 8     ; ++i
I11  Rpz <= Add Rpz, 8     ; ++i
I12  Cmp Rpx, End
I13  Bneq Loop
```

**AMD**
together we advance_

# Load-Load bypassing
**We want loads executed ASAP**

Scenario, n-th iteration

- $I2_n$: `F2 <= Ld [Rx]` misses in DL1 @Rx.
- We do not want to wait for the miss to do $I4_n$: `Ry <= Ld [Rpy]` and $I5_n$: `F4 <= Ld [Ry]`.
- So we execute $I4_n$ as soon as $I10_{n-1}$ is done, and $I5_n$ as soon as $I4_n$ is done.

Thoughts?

AMD🔺
together we advance_

# Load-Load bypassing
**We want loads executed ASAP**

Scenario, n-th iteration

- $I2_n$: F2 <= Ld [Rx] misses in DL1 @Rx.
- We do not want to wait for the miss to do $I4_n$: Ry <= Ld [Rpy] and $I5_n$: F4 <= Ld [Ry].
- So we execute $I4_n$ as soon as $I10_{n-1}$ is done, and $I5_n$ as soon as $I4_n$ is done.

Thoughts?

- Doing so would violate $I2_n <_p I4_n \Rightarrow I2_n <_m I4_n$.
    - If we are in relaxed consistency that is okay.
    - If we are in TSO or SC that is not okay!
- One way to fake it is as R10000 did: The eviction or invalidation of a cache block squashes any load to that block and all subsequent instructions.
    - If $S_1 : S(X) <_p S_2 : S(Y) <_m L_1 : L(Y) <_p L_2 : L(X)$, both TSO and SC have to guarantee that if $L_1$ reads the value written by $S_2$, $L_2$ needs to read the value from $S_1$ (or younger).
    - If $L_2$ gets a value older than $S_1$ it is guaranteed to observe an invalidation from $S_1$ before the data from $S_2$ arrives.

AMD🠪
together we advance_

# Load-Store bypassing
**We want loads executed ASAP**

Scenario, n-th iteration

- $I3_{n-1}$: `F2 <= Mul F2, F0` is a long latency operation/$I7_{n-1}$: `Rz <= Ld [Rpz]` misses in DL1 @Rz.
- Stores cannot write memory until they are the oldest instruction.
- We don't want $I1_n$ ($I2_n$, $I4_n$, $I5_n$) to have to wait until $I8_{n-1}$ is the oldest instruction/retires.
- So we execute $I1_n$ as soon as as $I9_{n-1}$ is done.

Thoughts?

AMD
together we advance_

# Load-Store bypassing
**We want loads executed ASAP**

Scenario, n-th iteration

- $I3_{n-1}$: F2 <= Mul F2, F0 is a long latency operation/$I7_{n-1}$: Rz <= Ld [Rpz] misses in DL1 @Rz.
- Stores cannot write memory until they are the oldest instruction.
- We don't want $I1_n$ ($I2_n$, $I4_n$, $I5_n$) to have to wait until $I8_{n-1}$ is the oldest instruction/retires.
- So we execute $I1_n$ as soon as as $I9_{n-1}$ is done.

Thoughts?

- If z[n-1] == x[n] we have to do Store-Load Forwarding to preserve program semantics.
- If z[n-1] != x[n], doing so would violate $I7_{n-1} <_p I1_n \Rightarrow I7_{n-1} <_m I1_n$.
  - If we are in relaxed consistency or TSO that is okay.
  - If we are in SC that is not okay!
- The same solution as in the previous case would handle it.

AMD
together we advance_

# Load-Store bypassing
## We want loads executed ASAP

Scenario, n-th iteration

- $I3_{n-1}$: `F2 <= Mul F2, F0` is a long latency operation/$I7_{n-1}$: `Rz <= Ld [Rpz]` misses in DL1 @Rz.
- Stores cannot write memory until they are the oldest instruction.
- We don't want $I1_n$ ($I2_n$, $I4_n$, $I5_n$) to have to wait until $I8_{n-1}$ is the oldest instruction/retires.
- So we execute $I1_n$ as soon as as $I9_{n-1}$ is done.

Thoughts?

- What if $I7_{n-1}$ misses in DL1 @Rz?
  - This is a program semantics issue =) nothing new as far as consistency is concerned.
  - Enforcing issue in order handles the situation.

AMD
together we advance_

# Store coalescing
**If stores overlap, let's send them together**

Scenario, n-th iteration

- `*z[n-2]` and `*z[n]` lie consecutively in the same cache line.
- `*z[n-1]` lies in a different cache line.
- We want to minimise the amount of cache writes, so we use a buffer to coalesce the data and perform two writes to DL1 instead of 3.

Thoughts?

AMD
together we advance_

# Store coalescing
**If stores overlap, let's send them together**

Scenario, n-th iteration

- `*z[n-2]` and `*z[n]` lie consecutively in the same cache line.
- `*z[n-1]` lies in a different cache line.
- We want to minimise the amount of cache writes, so we use a buffer to coalesce the data and perform two writes to DL1 instead of 3.

Thoughts?

- Doing so would violate $\text{I8}_{n-2} <_p \text{I8}_{n-1} <_p \text{I8}_n \Rightarrow \text{I8}_{n-2} <_m \text{I8}_{n-1} <_m \text{I8}_n$.
  - If we are in relaxed consistency that is okay.
  - If we are in TSO or SC that is not okay!
- There is literature handling it by delaying coherence responses to still enforce TSO.

AMD
together we advance_

# Push towards more relaxed consistency

- This should make it obvious why HW designers may want relaxed consistency.
- Why do we just not do all relaxed consistency then?
  - Code is not portable towards more restricted consistency model.
  - Ordering is a responsibility of the programmer (or compiler).
    - Harder to write (correct) software.
    - Achived through *fences*.
  - Fences are instructions such that impose order. If F is a *fence*, and I and J are any two instructions
    - $I <_p F \Rightarrow I <_m F$.
    - $F <_p J \Rightarrow F <_m J$.
- The flipside of fences is that
  - The most straigtforward way to enforce a fence is squashing all subsequent instructions upon retire.
  - Fences may have implications as they travel the memory hierarchy.

AMD⏉
together we advance_

# Conclusions

- You cannot ignore your memory model
  - If you are making SW, you need to know what outcomes are possible, so your memory sharers communicate the way you want them to.
  - If you are making HW, you do not need to implement it, faking it is just enough, and enables better performance.

AMD
together we advance_

## Copyright and Disclaimer