BPF Memory Model

Presented by Alexei Starovoitov on behalf of Paul McKenney



© 2023 Meta Platforms



Overview

2

- psABI and Memory Model
- BPF Memory Model Context
- BPF Instructions
- JITs Must Respect BPF Memory Model
- Validation: GCC Atomic Built-Ins
- Future Work

el ·ovt

Aemory Model Built-Ins

processor specific ABI (psABI)

All psABIs define:

- Function calling convention
- Register usage
- Stack usage and unwinding
- Type convention. ex: sizeof(void*)
- ELF object file format
- Relocations and linking
- Libraries
- Code model and address space <-

<- this is not a Memory Model !</p>

Goals of psABls

x86, RISC-V, and all others psABIs:

• A manual for the compilers: generate compatible binary code

Goals of psABls

x86, RISC-V, and all others psABIs:

• A manual for the compilers: generate compatible binary code

BPF psABI:

5

- A manual for the compilers: generate compatible binary code • A manual for JITs: how to map BPF ISA to native ISA

- It exists !
- JITs use it to translate BPF ISA to native ISA
- GCC and LLVM use it to compile C into BPF assembly
- Largely undocumented
- LLVM/GCC source code is a source of truth
- JITs source code is a source of truth

BPF psABI



BPF Memory-Model Context



- I want to concurrently access data:
- 1. Between BPF programs
- 2.Between BPF program and user space
- 3.Between BPF program and the kernel
- How would I do that?

BPF Memory-Model Context



8

- I want to concurrently access data:
- 1. Between BPF programs
- 2.Between BPF program and user space
- 3.Between BPF program and the kernel
- How would I do that?
- Just use Linux-kernel memory model (LKMM)

BPF Memory-Model Context



9

- I want to concurrently access data:
- 1. Between BPF programs
- 2.Between BPF program and user space
- 3.Between BPF program and the kernel
- How would I do that?
- Just use Linux-kernel memory model (LKMM)
- **Other language MMs are a strict subset of LKMM**



10 Flavors of Memory Models

Language Memory Model (C, C++, LKMM, ...)

Compiler

BPF Instruction Set Has No Memory Model

JIT

Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)

Flavors of Memory Models

11

Language Memory Model (C, C++, LKMM, ...)

Compiler

Instruction-Level BPF Memory Model

JIT

Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)







Aside on Linux-Kernel Memory Model

- The Linux kernel uses assembly, C, and Rust
- LKMM relies not just on the language memory model, but also on strict coding conventions:
 - memory-barriers.txt "CONTROL DEPENDENCIES"
 - rcu_dereference.rst
- Language MMs do not handle dependencies
 - And hence are plagued by OOTA issues
 - Therefore, a hardware-level model for BPF instruction set

https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2019/p1780r0.html https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2019/p1916r0.pdf

¹³ Example OOTA, x == y == 0 initially</sup>

- r1 = x.load(relaxed);
- y.store(r1, relaxed);

- r1 = y.load(relaxed);
- x.store(r1, relaxed);

¹⁴ Example OOTA, x == y == 0 initially

r1 = x.load(relaxed);
y.store(r1, relaxed);

According to the mathematical core of the C and C++ memory models, this code can result in x == y == 42!!!

- r1 = y.load(relaxed);
- x.store(r1, relaxed);



Example OOTA, x == y == 0 initially

r1 = x.load(relaxed);y.store(r1, relaxed);

According to the mathematical core of the C and C++ memory models, this code can result in x == y == 42!!!

- r1 = y.load(relaxed);
- x.store(r1, relaxed);

But only in theory, not in practice!



- It exists !
- Largely undocumented
- LLVM source code is a source of truth
- JITs and interpreter source code is a source of truth

Defining BPF Memory Model



BPF Instructions

17

- BPF Atomic Instructions
- BPF Conditional Jump Instructions
- BPF Load instructions

BPF Atomic Instructions

- BPF_XCHG, BPF_CMPXCHG
- BPF_ADD, BPF_OR, BPF_AND, BPF_XOR
- BPF_FETCH with one of the above

BPF Atomic Instructions 1/3

- BPF_XCHG and BPF_CMPXCHG instructions are fully ordered
- All CPUs and tasks agree that all instructions preceding or following a given BPF_XCHG or BPF_CMPXCHG instruction are ordered before or after, respectively, that same instruction
 - Consistent with Linux-kernel atomic_xchg() and atomic_cmpxchg(), respectively
 - Alternatively, consistent with the following:
 - smp_mb(); atomic_cmpxchg_relaxed(); smp_mb();

BPF Atomic Instructions 2/3

- BPF_ADD, BPF_OR, BPF_AND, BPF_XOR instructions are unordered
- CPUs and JITs can reorder these instructions freely
 - Consistent with Linux-kernel atomic_add(), atomic_or(), atomic_and(), and atomic_xor() APIs

BPF Atomic Instructions 3/3

- When accompanied by BPF_FETCH, BPF_ADD, BPF_OR, BPF_AND, BPF_XOR instructions are fully ordered
- respectively, that same instruction
 - Consistent with Linux-kernel atomic_fetch_add(),

 All CPUs and tasks agree that all instructions preceding or following a given instruction adorned with BPF_FETCH are ordered before or after,

atomic_fetch_or(), atomic_fetch_and(), and atomic_fetch_xor() APIs

- Modifiers to BPF_JMP32 and BPF_JMP instructions:
 - BPF_JEQ, BPF_JGT, BPF_JGE, BPF_JSET, BPF_JNE, BPF_JSGT, BPF_JSGE, BPF_JLT, BPF_JLE, BPF_JSLT, and BPF_JSLE
- Unconditional jump instructions (BPF_JA, BPF_CALL, BPF_EXIT) provide no memory-ordering semantics

- weak ordering:
 - BPF_JSGE, BPF_JLT, BPF_JLE, BPF_JSLT, and BPF_JSLE
- Too-smart JITs might need to be careful

These modifiers to BPF_JMP32 and BPF_JMP instructions provide

- BPF_JEQ, BPF_JGT, BPF_JGE, BPF_JSET, BPF_JNE, BPF_JSGT,

- This weak ordering applies when:
 - Either the src or dst registers depend on a prior load instruction (BPF_LD or BPF_LDX), *and*
 - There is a store instruction (BPF_ST or BPF_STX) before control flow converges, and
 - The restrictions outlined in the "CONTROL DEPENDENCIES" section of Documentation/memory-barriers.txt are faithfully followed
 - Compilers do not understand control dependencies, and happily break them.
 - Optimizations involving conditional-move instructions requires the "before control flow converges" restriction

Conditional Jump Example





26 **Control-Dependency Breakage**

```
r0 = READ ONCE(*x);
if (r0) {
    WRITE ONCE (*y, 1);
} else {
   WRITE_ONCE(*y, 1);
```

²⁷ Control-Dependency Breakage

```
r0 = READ_ONCE(*x);
if (r0) {
    WRITE_ONCE(*y, 1);
} else {
    WRITE_ONCE(*y, 1);
}
```

JIT or compiler Optimization

r0 = READ_ONCE(*x); WRITE_ONCE(*y, 1);

²⁸ Control-Dependency Breakage

```
r0 = READ_ONCE(*x);
if (r0) {
    WRITE_ONCE(*y, 1);
} else {
    WRITE_ONCE(*y, 1);
}
```

JIT or compiler Optimization

Broken Dependency!!!

r0 = READ_ONCE(*x);
WRITE_ONCE(*y, 1);

- Different hardware architectures order control dependencies in different ways:
 - Strongly ordered (x86, s390, …):
 - Prior load instructions are ordered before later store instructions, courtesy of TSO
 - Weakly ordered (ARMv8, PowerPC, ...):
 - Control dependencies are tracked by hardware

- What do you mean by "weak"???
 - CPU2
 - spanning CPU 1 and 2

CPU 0's control dependency is visible to CPU 1, and separately to

But CPU 0's control dependency is not necessarily visible to code

Example of Weakness in Play

CPU 0

WRITE_ONCE(*x, 1);

CPU1

r0 = READ_ONCE(*x); if (r0) { // Control dependency WRITE_ONCE(*y, 1); }

Both r0 instances and r1 can all be zero!!!

CPU 2

r0 = smp_load_acquire(y); $r1 = READ_ONCE(*x);$

32 Example For Converging Control Flow

cmov

uses

- r1 = READ ONCE(x);
- if (r1)
 - WRITE ONCE (y, 1);
- else
 - WRITE ONCE (y, 2);
- WRITE ONCE(z, 1); // Converged here

Which is why control dependencies extend only to control-flow convergence!!!



No ordering constraint!!!



BPF Load Instructions

- **BPF_LD** and **BPF_LDX** instructions
 - If the value returned by a given load instruction is used to compute the address of ____ a later load or store instruction, address-dependency ordering is guaranteed
 - If the value returned by a given load instruction is used to compute the value ____ stored by a given store instruction, data-dependency ordering is guaranteed
 - These are used by RCU readers, which must faithfully follow the restrictions outlined in Documentation/RCU/rcu_dereference.rst
 - Compilers do not understand address or data dependencies, and happily break them.
 - Address and data dependencies are weak, similar to control dependencies

BPF Load Instructions

- Different hardware architectures order address and data dependencies in different ways:
 - Strongly ordered (x86, s390, …):
 - Prior load instructions are ordered before later load and store instructions, courtesy of TSO
 - Weakly ordered (ARMv8, PowerPC, ...):
 - Address and data dependencies are tracked by hardware

Example of Weakness in Play

CPU 0

CPU 1

WRITE_ONCE(*x, 1);

r0 = READ_ONCE(*x);
// Data dependency
WRITE_ONCE(*y, r0);

Both r0 instances and r1 can all be zero!!!

CPU 2

r0 = smp_load_acquire(y); r1 = READ_ONCE(*x);

JITs must respect BPF Memory Model

- Viable strategies:
 - Preserve address, control, and data dependencies ____
 - Just generate instructions that match the BPF assembly code most closely •
 - Put atomic_signal_fence(memory_order_seq_cst) everywhere
 - Trace and explicitly preserve dependencies
 - Order prior loads before later stores: —
 - JIT every BPF_LD and BPF_LDX into a target-machine load-acquire instruction sequence
 - Place at least one target-machine load-to-store memory-barrier instruction between each BPF load/store instruction pair
 - atomic_signal_fence(memory_order_acquire) works on x86
 - Rely on source-level code having followed Linux-kernel coding standards —



Discovering BPF Memory Model via GCC Atomic Built-Ins

Note that GCC defined the built-ins, but this section uses only Clang/LLVM

GCC Atomic Memory Orders

- ____ATOMIC_RELAXED: Relaxed ordering
- ___ATOMIC_ACQUIRE: Acquire ordering
- ____ATOMIC_CONSUME: Treated as acquire
- ____ATOMIC_RELEASE: Release ordering
- _ATOMIC_ACQ_REL: Acquire/release ordering
- ____ATOMIC_SEQ_CST: Sequential consistency

No BPF C/C++ Weak Orderings

- ____ATOMIC_RELAXED: Relaxed ordering
- _ATOMIC_ACQUIRE, __ATOMIC_CONSUME, __ATOMIC_RELEASE, __ATOMIC_ACQ_REL, __ATOMIC_SEQ_CST: Full ordering
- Revisit when BPF does acquire and release

GCC Full Memory Barriers

- __atomic_thread_fence(__ATOMIC_SEQ_CST)
- BPF has none, but it can emulate them:
 - "BPF_ATOMIC | BPF_DW | BPF_STX" with an imm field of "BPF_ADD | BPF_FETCH" and a src register value of zero
 - Or: "lock *(u32 *)(r2 + 0) += r1"
 - Call it bpf_mb() for short

GCC Atomic Loads

- __atomic_load_n() & __atomic_load()
- Relaxed ordering:
 - BPF_LD or BPF_LDX
- Non-relaxed ordering:
 - BPF_LD or BPF_LDX followed by bpf_mb()

Clang/LLVM does not yet support this

GCC Atomic Stores

- __atomic_store_n() & __atomic_store()
- Relaxed ordering:
 - BPF_ST or BPF_STX
- Non-relaxed ordering:
 - bpf_mb() followed by BPF_ST or BPF_STX —

Clang/LLVM does not yet support this

GCC Atomic Exchange

- __atomic_exchange_n() & __atomic_exchange()
- No matter what ordering:

- "BPF_ATOMIC | BPF_DW | BPF_STX" with an immediate field of "BPF_XCHG | BPF_FETCH", which supplies full ordering

GCC Atomic Compare and Exchange

- __atomic_compare_exchange_n() & _atomic_compare_exchange()
- No matter what ordering:

- "BPF_ATOMIC | BPF_DW | BPF_STX" with an immediate field of "BPF_CMPXCHG | BPF_FETCH", which supplies full ordering

GCC Atomic Fetch-Op

- ___atomic_fetch_add(), __atomic_fetch_sub(), __atomic_fetch_and(), __atomic_fetch_xor() \bullet
 - "BPF_ATOMIC | BPF_DW | BPF_STX'' with an immediate field of "BPF_XXX | BPF_FETCH", which supplies full ordering, needed or not
 - Where "XXX" is ADD, SUB, AND, and XOR, respectively
- ___atomic_fetch_or(), __atomic_fetch_nand() lacksquare
 - Loop containing "BPF_ATOMIC | BPF_DW | BPF_STX'' with an immediate field of "BPF_CMPXCHG | BPF_FETCH", which supplies full ordering, needed or not
 - Use BPF_OR or a combination of BPF_AND with best bit-complement code, respectively

Clang/LLVM does not yet support __atomic_fetch_nand()

GCC Atomic Op-Fetch

- __atomic_add_fetch(), __atomic_sub_fetch(), __atomic_and_fetch(), __atomic_xor_fetch(), __atomic_or_fetch(), __atomic_nand_fetch()
 - Implement in the same way as for atomic_fetch_xxx()
 - Except that it is necessary to fix up return value to provide the after-operation value
 - Full ordering is supplied whether it is needed or not ____

Clang/LLVM does not yet support __atomic_nand_fetch()

GCC Miscellaneous Atomics

- __atomic_test_and_set()
 - Implement the same as ___atomic_exchange()
 - Except casting the return value to boolean if needed
- __atomic_clear()
 - Implement as an ___atomic_store() of zero _____

Clang/LLVM does not yet support these

GCC Fences

- __atomic_thread_fence()
 - Implement as bpf_mb()
- __atomic_signal_fence()
 - Implement as the Linux-kernel barrier() macro
 - Unless relaxed, in which case this is a no-op

Clang/LLVM does not yet support these

hel barrier() macro e this is a no-op

Complication: BPF Helper Ordering?

https://man7.org/linux/man-pages/man7/bpf-helpers.7.html

⁵⁰ Complication: BPF Helper Ordering?

Language Memory Model (C, C++, LKMM, ...)

Compiler

BPF Instruction Set (Proposing Memory Model)

Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)

JIT

https://man7.org/linux/man-pages/man7/bpf-helpers.7.html

C-Language BPF Helper (LKMM)

Compiler

Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)



⁵Complication: BPF Helper Ordering?

Language Memory Model (C, C++, LKMM, ...)

Compiler

BPF Instruction Set (Proposing Memory Model)

JIT

Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)

https://man7.org/linux/man-pages/man7/bpf-helpers.7.html

0 0

C-Language BPF Helper (LKMM)

Compiler

Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)



52 **Complication: BPF Helper Ordering?**

Language Memory Model (C, C++, LKMM, ...)

BPF Instruct (Proposing Mem by the definition of mpiler

Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)

https://man7.org/linux/man-pages/man7/bpf-helpers.7.html

ct 3 **C-Language BPF Helper** (LKMM) No ordering unless specified the helper in question 5 Hardware ISA Memory Model (x86, ARMv8, RISC-V, ...)



BPF ISA extensions

- Possible additions:
 - BPF_LDX_ACQ load acquire
 - **BPF_STX_REL** store release store ____
 - Full barrier
 - Possibly one variant with I/O semantics and another variant having only normal-memory semantics
 - Normal-memory semantics (smp_mb()) is more urgent

54 For More Information

- "Instruction-Level BPF Memory Model" •
 - ____
- "IETF eBPF Instruction Set Specification, v1.0" \bullet
 - https://www.ietf.org/archive/id/draft-thaler-bpf-isa-00.html —
- "Towards a BPF Memory Model" (2021 BPF & Networking Summit at Linux Plumbers Conference)
 - https://lpc.events/event/11/contributions/941/ ____
- "GCC Atomic Compiler Built-Ins"
 - https://gcc.gnu.org/onlinedocs/gcc/_005f_005fatomic-Builtins.html ____
- lacksquaresection, Documentation/RCU/rcu_dereference.rst
- "Is Parallel Programming Hard, And, If So, What Can You Do About It?" \bullet
 - Chapter 15 ("Advanced Synchronization: Memory Ordering") —
 - Appendic C ("Why Memory Barriers?") ____

https://docs.google.com/document/d/1TaSEfWfLnRUi5KqkavUQyL2tThJXYWHS15qcbxIsFb0/edit?usp=sharing

Linux kernel source tree: tools/memory-model, Documentation/memory-barriers.txt "CONTROL DEPENDENCIES"

https://mirrors.edge.kernel.org/pub/linux/kernel/people/paulmck/perfbook/perfbook.html