

THE BPF TARGET IN LLVM



ABOUT ME

- Michal Rostecki
- **vadorovsky** @ Github, Discord, Mastodon etc.
- Software Engineer @ [Light Protocol](#)
- [Aya](#) maintainer

BPF

- "Berkeley Packet Filter"
- ...but that's not what it is, really.
- It's a virtual machine for running minimal, sandboxed programs.
- Restricted environment (instruction limit, verifier).
- Rules: execute fast, don't sleep*, don't allocate too much memory.

USE CASES

- **Linux kernel** - BPF originates from there
 - Network filtering
 - Tracing and profiling
 - Security policies
- **Solana**
 - Smart contracts

PROJECTS

- **Cilium** - Container Network Interface for Kubernetes.
- **Falco** - Security alert engine.
- Pretty much every Solana project. Orca, Mango, Audius, Metaplex, Light Protocol etc.
- **Good stuff from Deepfence coming soon!** Security, but with focus on enforcement.

BPF AND LLVM

- First and main compiler providing BPF support.
- BPF has a backend in LLVM.
- Supported by Clang (C) and Rust.

REGISTERS

- **R0** - return value from functions, exit values for programs.
- **R1 - R5** - function arguments.
- **R6 - R9** - callee saved registers that function calls will preserve.
- **R10** - read-only frame pointer to access stack.

R0 - R5: scratch registers, programs need to spill them if necessary across calls.

BASIC INSTRUCTION ENCODING

32 bits
(MSB)

16
bits

4 bits

4 bits

8 bits
(LSB)

Integer
Immediate
Value

Offset

Source
register

Destination
register

Opcode

WIDE INSTRUCTION ENCODING

ARITHMETIC INSTRUCTIONS

Code	Description
BPF_ADD	dst += src
BPF_SUB	dst -= src
BPF_MUL	dst *= src
BPF_DIV	dst = (src != 0) ? (dst / src) : 0
BPF_OR	dst
BPF_AND	dst &= src

BYTE SWAP INSTRUCTIONS

Code	Description
BPF_TO_LE	host byte order -> little endian
BPF_TO_BE	host byte order -> big endian

JUMP INSTRUCTIONS

Code	Description
BPF_CALL	function call
BPF_EXIT	program return

LOAD AND STORE INSTRUCTIONS

ATOMIC OPERATIONS

VERIFIER

- (So far) only in Linux kernel (Solana, rBPF and other user space implementations don't have it).
- Ensuring safe memory access - kinda, like, making C Rusty (◡‿◡).
- DAG check preventing unbound loops.
- Descending all possible instruction paths, observing the change of registers and stack.

VERIFIER ERRORS

- unreachable insn
- !read_ok
- invalid stack
- invalid indirect read from stack
- invalid mem access
- offset is outside of the packet

BPF TYPE FORMAT (BTF)

- Debug info format for BPF, way more lightweight than DWARF.
- Used for offsets across Linux kernel versions.
- Used for stack traces in BPF verifier.
- But there is no debugger (yet).

EXAMPLE: C CODE

```
struct foo {  
    __u32 a;  
    __u64 b;  
};
```

EXAMPLE: LLVM DEBUG INFO

```
!49 = distinct !DICompositeType(tag: DW_TAG_structure_type,  
  name: "foo", file: !3, line: 21, size: 128, elements: !50)  
!50 = !{!51, !54}  
!51 = !DIDerivedType(tag: DW_TAG_member, name: "a", scope: !49,  
  file: !3, line: 22, baseType: !52, size: 32)  
!52 = !DIDerivedType(tag: DW_TAG_typedef, name: "__u32",  
  file: !53, line: 27, baseType: !12)  
!54 = !DIDerivedType(tag: DW_TAG_member, name: "b", scope: !49,  
  file: !3, line: 23, baseType: !55, size: 64, offset: 64)  
!55 = !DIDerivedType(tag: DW_TAG_typedef, name: "__u64",  
  file: !53, line: 31, baseType: !56)
```

EXAMPLE: BTF

```
[8] STRUCT 'foo' size=16 vlen=2  
    'a' type_id=9 bits_offset=0  
    'b' type_id=10 bits_offset=64
```

EXAMPLE: RUST CODE

```
pub struct Foo {  
    a: u32,  
    b: u64,  
}
```

EXAMPLE: LLVM DEBUG INFO

```
!42 = !DIBasicType(name: "u32", size: 32,  
    encoding: DW_ATE_unsigned)  
!60 = !DICompositeType(tag: DW_TAG_structure_type, name: "Foo",  
    scope: !2, file: !5, size: 128, align: 64, elements: !61,  
    templateParams: !65,  
    identifier: "63dcf8d9f7a7a7ed6f05eaed70c4b12f")  
!61 = !{!62, !63}  
!62 = !DIDerivedType(tag: DW_TAG_member, name: "a", scope: !60,  
    file: !5, baseType: !42, size: 32, align: 32, offset: 64)  
!63 = !DIDerivedType(tag: DW_TAG_member, name: "b", scope: !60,  
    file: !5, baseType: !64, size: 64, align: 64)  
!64 = !DIBasicType(name: "u64", size: 64,  
    encoding: DW_ATE_unsigned)
```

EXAMPLE: BTF

```
[20] STRUCT 'Foo' size=16 vlen=2  
      'a' type_id=14 bits_offset=64  
      'b' type_id=21 bits_offset=0
```

LOCAL BTF

Each modern Linux kernel comes with BTF info:

```
# bpftool btf dump file /sys/kernel/btf/vmlinux
[...]
[335] STRUCT 'pid' size=112 vlen=8
      'count' type_id=330 bits_offset=0
      'level' type_id=59  bits_offset=32
      'lock'  type_id=324 bits_offset=64
      'tasks' type_id=336 bits_offset=128
      'inodes' type_id=134 bits_offset=384
      'wait_pidfd' type_id=328 bits_offset=448
      'rcu' type_id=139 bits_offset=640
      'numbers' type_id=337 bits_offset=768
```


ELF SECTIONS

.BTF SECTION

.BTF.EXT SECTION

BTF_IDS SECTION

BTF RELOCATIONS

- BPF programs are adjusted to read type fields at the offset specified in **local** BTF info.
- Regardless of the memory layout of the type.
- Types with BTF-based access are annotated with `llvm.preserve.*.access.index` intrinsics.

CHALLENGES WITH RUST

- BPF support introduced later than in Clang.
- BTF emission not supported, but close to be done!
- BTF relocations not supported.

WHAT'S THE PROBLEM?

- Kernel expects specific BTF layout.
- It's very C-specific.
 - BPF maps definitions have to be anonymous structs (which Rust doesn't support).
 - Complex Rust types (e.g. data carrying enums) are not supported.

SOLUTIONS

- Temporary: modify DI in bpf-linker.
- Long-term: # [btf_export] macro in Rust.

BPF-LINKER

Currently working PoC. Transforms DI to meet kernel expectations:

- Removes names from pointer types and BTF map structs.
- Tweaks the DI of Rust-specific types to be C-compatible.

BTF (FROM RUST) AFTER MODIFICATIONS

```
[10] STRUCT '(anon)' size=40 vlen=5
      'type' type_id=1 bits_offset=0
      'key' type_id=5 bits_offset=64
      'value' type_id=5 bits_offset=128
      'max_entries' type_id=6 bits_offset=192
      'map_flags' type_id=8 bits_offset=256
```

```
[6] PTR '(anon)' type_id=7
```

DEBUG INFO INCLUDED

```
; let parent_pid: i32 = unsafe {  
    ctx.read_at(PARENT_PID_OFFSET)? };  
9: 63 1a f4 ff 00 00 00 00 *(u32 *) (r10 - 0xc) = r1  
[...]  
; let child_pid: i32 = unsafe {  
    ctx.read_at(CHILD_PID_OFFSET)? };  
19: 63 1a f8 ff 00 00 00 00 *(u32 *) (r10 - 0x8) = r1  
20: bf a2 00 00 00 00 00 00 r2 = r10
```

LLVM CHANGES

- Already merged, but to be released in LLVM 17:
 - `LLVMGetDINodeTag` function to get the tag of DI Node.
 - `LLVMReplaceMDNodeOperandWith` function to modify DI.

#[BTF_EXPORT]

- Decoupled from `-C debuginfo`.
- Generates DI, which produces correct BTF for annotated types.
- Raises a compiler error when used on BTF-incompatible type.

IF YOU WANT TO TRY IT OUT

- github.com/vadorovsky/aya-btf-map - structs and macros for BTF maps.
- github.com/vadorovsky/aya-btf-maps-experiments - example project using it.
- Requires LLVM and bpf-linker patches.

THANK YOU

- aya-rs.dev
 - github.com/aya-rs/aya
 - [Discord](#)
- lightprotocol.com

