

## BPF – in-kernel virtual machine

# BPF

- BPF - Berkeley Packet Filter
- inspired by BSD
- introduced in linux in 1997 in version 2.1.75
- initially used as socket filter by packet capture tool tcpdump (via libpcap)

# Classic BPF

- two 32-bit registers: A, X
- implicit stack of 16 32-bit slots (LD\_MEM, ST\_MEM insns)
- full integer arithmetic
- explicit load/store from packet (LD\_ABS, LD\_IND insns)
- conditional branches (with two destinations: jump true/false)

# Ex: tcpdump syntax and classic BPF assembler

- tcpdump 'ip and tcp port 22'

```
(000) ldh      [12]           // fetch eth proto
(001) jeq     #0x800   jt 2   jf 12 // is it IPv4 ?
(002) ldb     [23]           // fetch ip proto
(003) jeq     #0x6      jt 4   jf 12 // is it TCP ?
(004) ldh     [20]           // fetch frag_off
(005) jset    #0x1fff   jt 12  jf 6  // is it a frag?
(006) ldx     4*([14]&0xf) // fetch ip header len
(007) ldh     [x + 14]       // fetch src port
(008) jeq     #0x16      jt 11  jf 9  // is it 22 ?
(009) ldh     [x + 16]       // fetch dest port
(010) jeq     #0x16      jt 11  jf 12 // is it 22 ?
(011) ret     #65535        // trim packet and pass
(012) ret     #0           // ignore packet
```

# Classic BPF for packets

- input: skb
- output: use case specific
- socket filters:
  - ret 0 – don't pass to user space
  - ret N – trim packet to N bytes and pass to user space
- cls\_bpf:
  - ret 0 – no match
  - ret -1 – default classid
  - ret > 0 – overwrite classid

# Classic BPF for seccomp

- introduced in 2012 to filter syscall arguments with bpf program (used by chrome for sandboxing)
- input: struct seccomp\_data  
LD\_ABS insn was overloaded
- output:
  - ret 0 – kill task
  - ret 30000 | data – force sigsys
  - ret 50000 | errno – return errno
  - ret 7ff00000 – pass to trace
  - ret 7fff0000 – allow

# Classic BPF safety

- verifier checks all instructions, forward jumps only, stack slot load/store, etc
- instruction set has some built-in safety (no exposed stack pointer, instead load instruction has 'mem' modifier)
- dynamic packet-boundary checks

# Classic BPF extensions

- over years multiple extensions were added in the form of 'load from negative hard coded offset'
- LD\_ABS -0x1000 – skb->protocol  
LD\_ABS -0x1000+4 – skb->pkt\_type  
LD\_ABS -0x1000+56 – get\_random()

# Extended BPF

- design goals:
  - parse, lookup, update, modify network packets
  - loadable as kernel modules on demand, on live traffic
  - safe on production system
  - performance equal to native x86 code
  - fast interpreter speed (good performance on all architectures)
  - calls into bpf and calls from bpf to kernel should be free (no FFI overhead)

# in kernel 3.15

tcpdump

dhclient

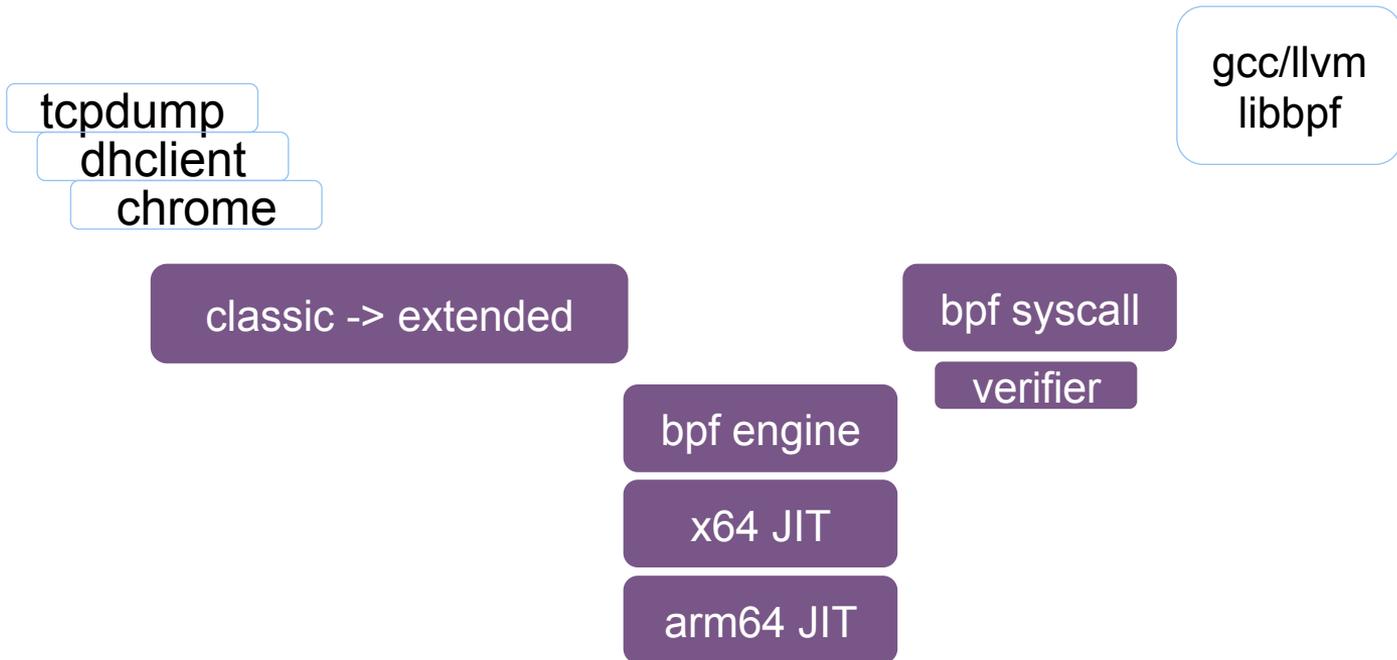
chrome

cls,  
xt,  
team,  
ppp,  
...

classic -> extended

bpf engine

# in kernel 3.18



# Early prototypes

- Failed approach #1 (design a VM from scratch)
  - performance was too slow, user tools need to be developed from scratch as well
- Failed approach #2 (have kernel disassemble and verify x86 instructions)
  - too many instruction combinations, disasm/verifier needs to be rewritten for every architecture

# Extended BPF

- take a mix of real CPU instructions: 70% x86 + 25% arm64 + 5% risc
- rename every x86 instruction 'mov rax, rbx' into 'mov r1, r2'
- analyze x86/arm64/risc calling conventions and define a common one for this 'renamed' instruction set
- make instruction encoding fixed size (for high interpreter speed)
- reuse classic BPF instruction encoding (for trivial classic->extended conversion)

# Extended vs classic BPF

- ten 64-bit registers vs two 32-bit registers
- arbitrary load/store vs stack load/store
- call instruction

# Performance

- user space compiler ‘thinks’ that it’s emitting simplified x86 code
- kernel verifies this ‘simplified x86’ code
- kernel JIT translates each ‘simplified x86’ insn into real x86
  - all registers map one-to-one
  - most of instructions map one-to-one
  - bpf ‘call’ instruction maps to x86 ‘call’

# Performance

- bpf 'call' and set of in-kernel helper functions define what bpf programs can do
- bpf code itself is a 'glue' between calls to in-kernel helper functions

# Extended BPF calling convention

- BPF calling convention was carefully selected to match a subset of amd64/arm64 ABIs to avoid extra copy in calls:
- R0 – return value
- R1..R5 – function arguments
- R6..R9 – callee saved
- R10 – frame pointer

# Mapping of BPF registers to x86

- R0 – rax return value from function
- R1 – rdi 1<sup>st</sup> argument
- R2 – rsi 2<sup>nd</sup> argument
- R3 – rdx 3<sup>rd</sup> argument
- R4 – rcx 4<sup>th</sup> argument
- R5 – r8 5<sup>th</sup> argument
- R6 – rbx callee saved
- R7 – r13 callee saved
- R8 – r14 callee saved
- R9 – r15 callee saved
- R10 – rbp frame pointer

# BPF compilers

- BPF backend for LLVM is in trunk and will be released as part of 3.7
- BPF backend for GCC is being worked on
- C front-end (clang) is used today to compile C code into BPF
- tracing and networking use cases may need custom languages
- BPF backend knows how to emit instructions (calls to helper functions look like normal calls)

# Extended BPF assembler

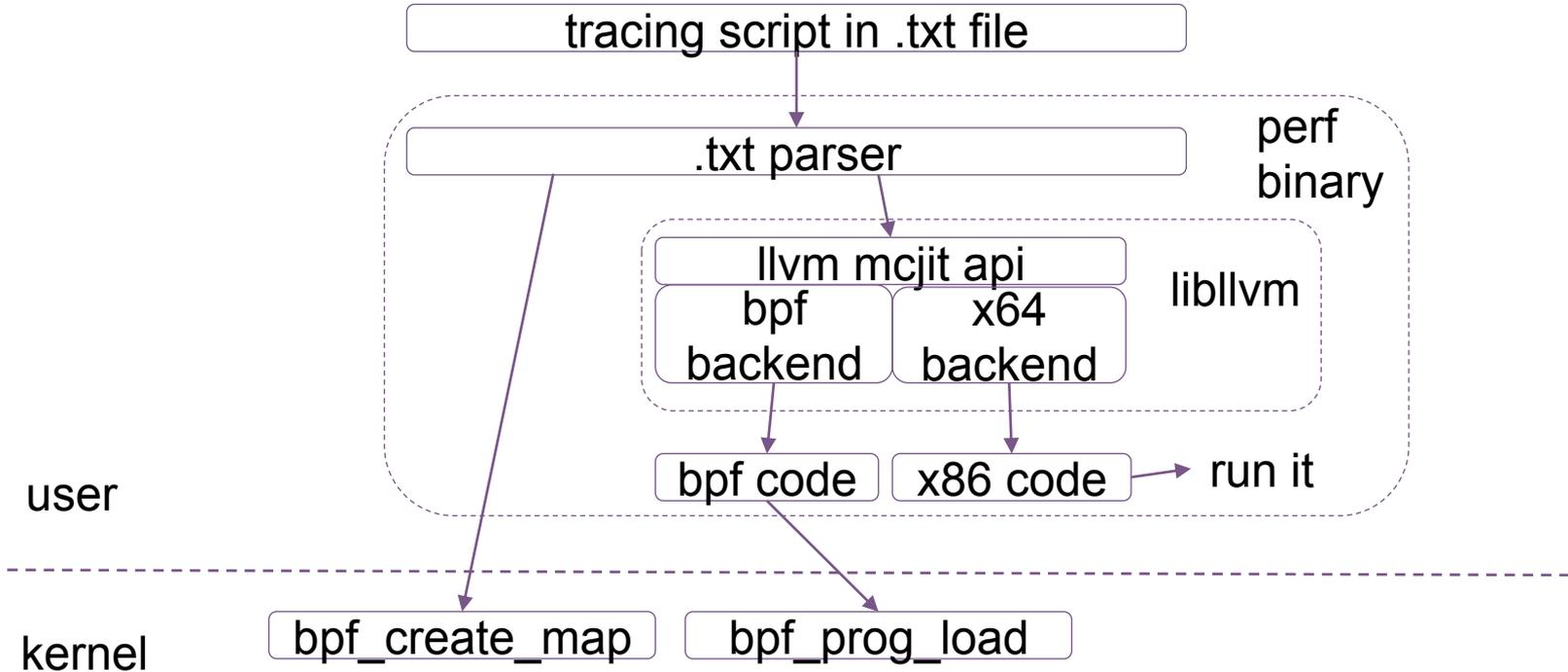
```
int bpf_prog(struct bpf_context *ctx)
{
    u64 loc = ctx->arg2;
    u64 init_val = 1;
    u64 *value;

    value = bpf_map_lookup_elem(&my_map, &loc);
    if (value)
        *value += 1;
    else
        bpf_map_update_elem(&my_map, &loc,
                            &init_val, BPF_ANY);
    return 0;
}
```

compiled by LLVM from C to bpf asm

```
0: r1 = *(u64 *)(r1 +8)
1: *(u64 *)(r10 -8) = r1
2: r1 = 1
3: *(u64 *)(r10 -16) = r1
4: r1 = map_fd
6: r2 = r10
7: r2 += -8
8: call 1
9: if r0 == 0x0 goto pc+4
10: r1 = *(u64 *)(r0 +0)
11: r1 += 1
12: *(u64 *)(r0 +0) = r1
13: goto pc+8
14: r1 = map_fd
16: r2 = r10
17: r2 += -8
18: r3 = r10
19: r3 += -16
20: r4 = 0
21: call 2
22: r0 = 0
23: exit
```

# compiler as a library



# BPF maps

- maps is a generic storage of different types for sharing data between kernel and userspace
- The maps are accessed from user space via BPF syscall, which has commands:
  - create a map with given type and attributes  
`map_fd = bpf(BPF_MAP_CREATE, union bpf_attr *attr, u32 size)`
  - lookup key/value, update, delete, iterate, delete a map
- userspace programs use this syscall to create/access maps that BPF programs are concurrently updating

# BPF verifier (CFG check)

- To minimize run-time overhead anything that can be checked statically is done by verifier
- all jumps of a program form a CFG which is checked for loops
  - DAG check = non-recursive depth-first-search
  - if back-edge exists -> there is a loop -> reject program
  - jumps back are allowed if they don't form loops
  - bpf compiler can move cold basic blocks out of critical path
  - likely/unlikely() hints give extra performance

# BPF verifier (instruction walking)

- once it's known that all paths through the program reach final 'exit' instruction, brute force analyzer of all instructions starts
- it descends all possible paths from the 1st insn till 'exit' insn
- it simulates execution of every insn and updates the state change of registers and stack

# BPF verifier

- at the start of the program:
  - type of R1 = PTR\_TO\_CTX  
type of R10 = FRAME\_PTR  
other registers and stack is unreadable
- when verifier sees:
  - 'R2 = R1' instruction it copies the type of R1 into R2
  - 'R3 = 123' instruction, the type of R3 becomes CONST\_IMM
  - 'exit' instruction, it checks that R0 is readable
  - 'if (R4 == 456) goto pc+5' instruction, it checks that R4 is readable and forks current state of registers and stack into 'true' and 'false' branches

# BPF verifier (state pruning)

- every branch adds another fork for verifier to explore, therefore branch pruning is important
- when verifiers sees an old state that has more strict register state and more strict stack state then the current branch doesn't need to be explored further, since verifier already concluded that more strict state leads to valid 'exit'
- two states are equivalent if register state is more conservative and explored stack state is more conservative than the current one

# BPF for tracing

- BPF is seen as alternative to systemtap/dtrace
- provides in-kernel aggregation, event filtering

# BPF for tracing (kernel part)

```
struct bpf_map_def SEC("maps") my_hist_map = {  
    .type = BPF_MAP_TYPE_ARRAY,  
    .key_size = sizeof(u32),  
    .value_size = sizeof(u64),  
    .max_entries = 64,  
};
```

} sent to kernel as bpf map via bpf() syscall

```
SEC("events/syscalls/sys_enter_write")  
int bpf_prog(struct bpf_context *ctx)  
{  
    u64 write_size = ctx->arg3;  
    u32 index = log2(write_size);  
    u64 *value;  
  
    value = bpf_map_lookup_elem(&my_hist_map, &index);  
    if (value)  
        __sync_fetch_and_add(value, 1);  
    return 0;  
}
```

} name of elf section - tracing event to attach via perf\_event ioctl

} compiled by llvm into .o and loaded via bpf() syscall

# BPF for tracing (user part)

```
u64 data[64] = {}; u32 key; u64 value;
```

```
for (key = 0; key < 64; key++) {  
    bpf_lookup_elem(fd, &key, &value);  
    data[key] = value;  
    if (value && key > max_ind)  
        max_ind = key;  
    if (value > max_value)  
        max_value = value;  
}  
printf("syscall write() stats\n");
```

} user space walks the map and fetches elements via bpf() syscall

```
syscall write() stats
```

byte_size	count	distribution
1 -> 1	: 9	*****
2 -> 3	: 0	
4 -> 7	: 0	
8 -> 15	: 2	*****
16 -> 31	: 0	
32 -> 63	: 10	*****
64 -> 127	: 12	*****
128 -> 255	: 1	**
256 -> 511	: 2	*****

# Cute Tracing Logos

(Brendan Gregg's slide)



ftrace



perf\_events



SystemTap



ktap



LTTng



dtrace4linux

# Extended BPF



# BPF for networking

- Extended BPF programs can be attached to sockets similar to classic BPF:
  - `setsockopt(sock, SOL_SOCKET, SO_ATTACH_BPF, prog_fd, sizeof(prog_fd))`
  - input: `skb`
  - output: 0 – ignore packet, N – trim packet to N bytes and pass to userspace

# BPF for networking

- three different attachment points in ovs+bpf
  - bpf as an action on flow-hit
  - bpf as packet parser before flow lookup
  - bpf as fallback on flow-miss
- two attachment points in tc+bpf
  - cls – packet parser and classifier
  - act – action

demo