# XDP for the Rest of Us

Andy Gospodarek -- Principal Engineer, Broadcom
Jesper Dangaard Brouer --  Principal Engineer, Red Hat
Netdev 2.1, 6-8 April 2017

# Motivation for this talk

- XDP is still fairly new (v4.8)
- Discussion often centers around datacenter use-cases, often with expensive processor and NIC hardware
- Benefits of early packet drop (DDoS prevention) are significant on lower-powered systems
- Use some alternative hardware and encourage other driver maintainers to consider adding XDP support

# Example Application:
# XDP DDoS Blacklist

https://github.com/netoptimizer/prototype-kernel

# Example Application: XDP DDoS Blacklist

- Non-interactive eBPF program that drops traffic based on:
  - Source IPv4 addresses
  - Destination UDP and/or TCP ports
- Has a command-line tool that queries BPF maps for:
  - Configuration
  - Statistics (real-time and historical)
- Motivated to bring XDP processing to more networks using both x86 and ARM64 CPUs

# XDP DDoS Blacklist Configuration

```
01 # ./xdp_ddos01_blacklist_cmdline --list
02 {
03   "216.239.59.128" : 12024134,
04   "64.68.90.1" : 67809613,
05   "22" : {
06      "TCP" : 0
07   },
08   "80" : {
09      "TCP" : 10,
10      "UDP" : 4216804996
11   }
12 }
```

# XDP DDoS Blacklist Real-time Stats

```
01 # ./xdp_ddos01_blacklist_cmdline --stats
02
03 XDP_action     pps          pps-human-readable period/sec
04 XDP_ABORTED    0            0                  1.000571
05 XDP_DROP       5975205      5,975,205          1.000568
06 XDP_PASS       1102227      1,102,227          1.000568
07 XDP_TX         0            0                  1.000568
```

# Getting Started with eBPF and XDP

# Quick Review: eBPF Fundamentals

- Berkley Packet Filter (BPF) is a special-purpose virtual machine for filtering packets (circa 1992)
- eBPF is a 'universal in-kernel virtual machine' that is not necessarily network-specific (circa 2014)
- Programs are small and have programming restrictions
- Primary interaction with outside world is generic key/value store called an eBPF map
- Over a dozen different types of maps and growing

# Quick Review: eXpress Data Path (XDP)

- Programmable, high performance networking data-path
- The "packet-page" idea from previous Netdev/Netconf
- Operates on raw packet data (before SKB is allocated) directly in driver
- eBPF program result:
  - DROP (never to be seen again)
  - TX (with or without modification/encapsulation)
  - PASS (on to kernel stack for processing!)

# Getting Started with XDP (kernel)

- eBPF added in 3.15, XDP core in 4.8
- Native JIT compiling (higher performance execution): x86_64 (3.16), ARM (3.18), s390 (4.1), PowerPC64 (4.8)
- Mellanox: mlx4 (4.8) and mlx5 (4.9)
- QLogic/Cavium: qede (4.10)
- Virtio_net: (4.10)
- Netronome: nfp (4.10)
- Broadcom: bnxt_en (4.11)
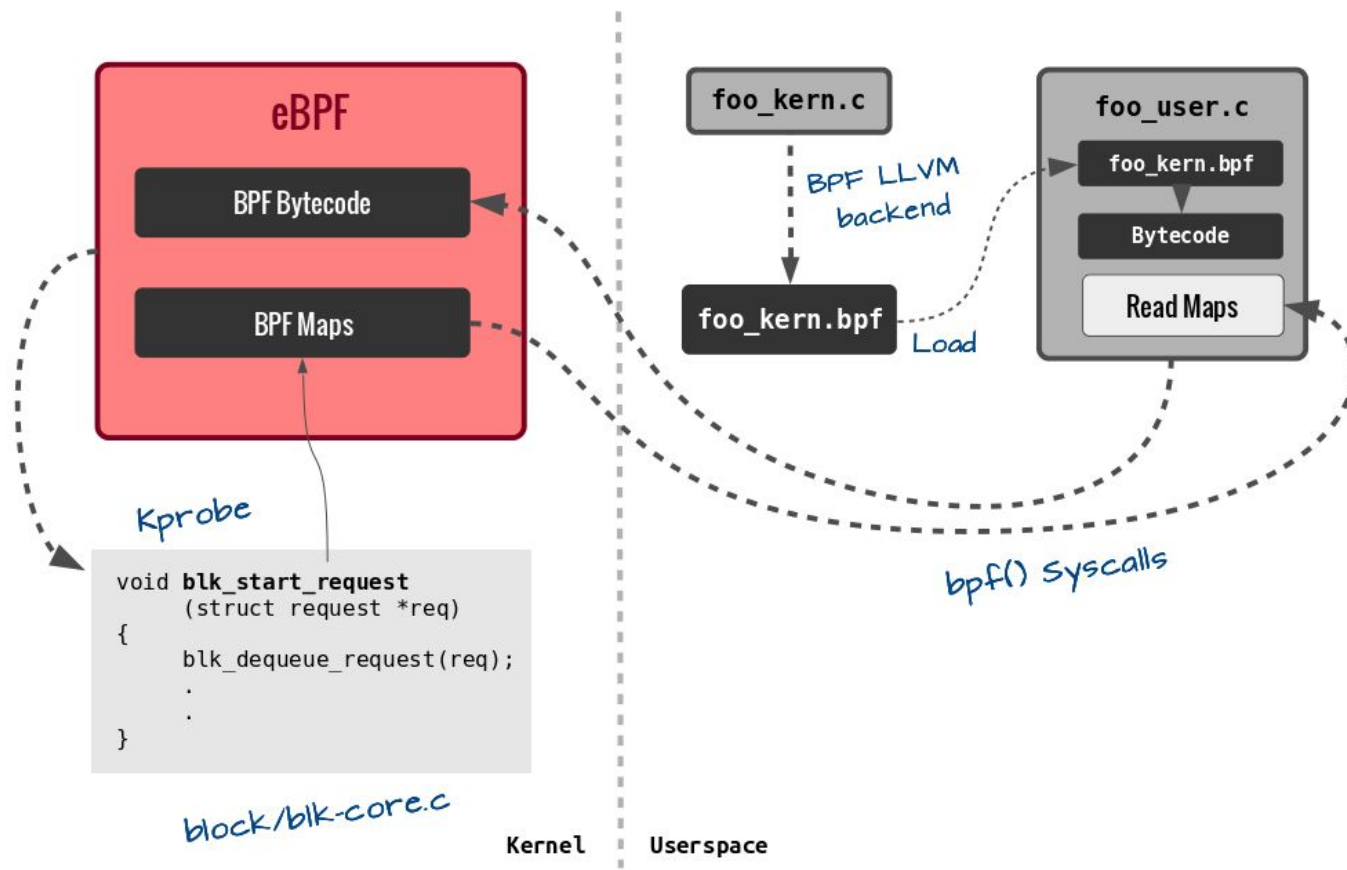
# Getting Started with XDP (compilers and distros)

- Roll your own compiler/tools with eBPF support:
    - LLVM >= 3.7.1
    - clang >= 3.4.0
- Fedora25 works with latest kernel and default LLVM/clang
- Ubuntu may require that you load a new kernel and build LLVM/clang yourself for latest features (depending on when you read these slides)

# Getting Started with XDP (free samples)

- Three important programs in samples/bpf in kernel tree:
  - xdp1 - drops all frames received on an interface
  - xdp2 - swaps src and dest MAC addresses and retransmits frames on same interface
  - xdp_tx_iptunnel - matches IP and dest port and encapsulates matching frames inside a new IP tunnel and sends them out the same interface

# Getting Started with XDP (free samples continued)

- Each program divides source into 2 files:
- eBPF code and map definitions in foo_kern.c
- Output ELF file with map layout and eBPF object code compiled to foo_kern.o by llvm with help from clang
- Userspace code used to load eBPF program and interact with maps compiled to foo_user.o with gcc

From: https://github.com/iovisor/bpf-docs/blob/master/bpf-internals-2.md

Enough already!  Let's see how this application actually works.

# File split similar to samples/bpf in kernel source

- **xdp_ddos01_blacklist_kern.c** - map definitions and restricted-C that is compiled into eBPF
- **xdp_ddos01_blacklist_user.c** - userspace code that loads xdp_ddos01_blacklist_kern.o, sets up sysfs links, and exits
- **xdp_ddos01_blacklist_common.h** - a few definitions and blacklist modify operations
- **xdp_ddos01_blacklist_cmdline.c** - code that is compiled into the userspace program used to configure blacklists

The following code has been modified from its original version. It has been formatted to fit this screen, to run in the presentation time allotted, edited for debug statements, code comments, and [possibly] bugs.

# xdp_ddos01_blacklist_kern.c

# Create per CPU Hash Map

```
01 struct bpf_map_def SEC("maps") blacklist = {
02        .type        = BPF_MAP_TYPE_PERCPU_HASH,
03        .key_size    = sizeof(u32),
04        .value_size  = sizeof(u64), /* Drop counter */
05        .max_entries = 100000,
06        .map_flags   = BPF_F_NO_PREALLOC,
07 };
```

# Create Function that Kernel Will Call

```c
01 SEC("xdp_prog")
02 int xdp_program(struct xdp_md *ctx)
03 {
04         void *data_end = (void *)(long)ctx->data_end;
05         void *data     = (void *)(long)ctx->data;
06         struct ethhdr *eth = data;
07         u16 eth_proto, l3_offset = 0;
08         u32 action;
09
10         if (!(parse_eth(eth, data_end, &eth_proto, &l3_offset)))
11                 bpf_debug("Cannot parse L2: L3off:%llu proto ...
12                          l3_offset, eth_proto);
13                 return XDP_PASS;
14         }
15         action = handle_eth_protocol(ctx, eth_proto, l3_offset);
16         stats_action_verdict(action);
17         return action;
18 }
```

# Start Parsing Packets

```c
01 static __always_inline
02 bool parse_eth(struct ethhdr *eth, void *data_end,
03                u16 *eth_proto, u64 *l3_offset)
04 {
05      u16 eth_type;
06      u64 offset;
07
08      offset = sizeof(*eth);
09      if ((void *)eth + offset > data_end)
10              return false;
11
12      eth_type = eth->h_proto;
13
```

# Continue Parsing Packets

```
01 static __always_inline
02 u32 handle_eth_protocol(struct xdp_md *ctx, u16 eth_proto, ...
03 {
04      switch (eth_proto) {
05      case ETH_P_IP:
06              return parse_ipv4(ctx, l3_offset);
07              break;
08      case ETH_P_IPV6: /* Not handler for IPv6 yet*/
09      case ETH_P_ARP:  /* Let OS handle ARP */
10              break;
11      default:
12              bpf_debug("Not handling eth_proto:0x%x\n", eth_proto);
13              return XDP_PASS;
14      }
15      return XDP_PASS;
16 }
```

# Parse IPv4 address and Lookup in Map

```
01 static __always_inline
02 u32 parse_ipv4(struct xdp_md *ctx, u64 l3_offset)
03 {
04        void *data_end = (void *)(long)ctx->data_end;
05        void *data     = (void *)(long)ctx->data;
06        struct iphdr *iph = data + l3_offset;
07        u64 *value;
08        u32 ip_src; /* type need to match map */
09
10        if (iph + 1 > data_end)
11                return XDP_ABORTED;
12        ip_src = iph->saddr;
13
14        value = bpf_map_lookup_elem(&blacklist, &ip_src);
15        if (value) {
17                *value += 1; /* Keep a counter for drop matches */
18                return XDP_DROP;
19        }
```

# Review xdp_ddos01_blacklist_kern.c functionality

- Define map for IPv4 blacklist
- Provide function written in restricted-C that will be called by kernel
- Track whether packet was acceptable or dropped which entry caused drop
- Report decision back to caller
- But how is this code loaded and attached to a netdev?

# xdp_ddos01_blacklist_user.c

# The Magic of Interacting with eBPF Maps

- BPF library in tools/lib/bpf of kernel tree is extremely helpful for userspace applications, but feels bit magical until you dig into it:
  - `prog_fd[]` entries are populated with each call to `load_bpf_file()`
  - `map_fd[]` entries align with maps declared in program loaded via `load_bpf_file()` (e.g. `foo_kern.o`)

# The importance of prog_fd[]

- File descriptors for a BPF program are loaded with
  `load_bpf_file(foo_kern.o)`, stored in `prog_fd[]`, and attached
  to netdev with call to `set_link_xdp_fd()`

```
01        if (!prog_fd[0]) {
02                printf("load_bpf_file: %s\n", strerror(errno));
03                return 1;
04        }
05
06        if (set_link_xdp_fd(ifindex, prog_fd[0]) < 0) {
07                printf("link set xdp fd failed\n");
08                return EXIT_FAIL_XDP;
09        }
```

# The importance of map_fd[]

- File descriptors for each BPF map are loaded with
  `load_bpf_file(foo_kern.o)`, and stored in `map_fd[]`, array.
  Those fds are used as first argument to operate on the
  map with library calls (from `tools/lib/bpf/bpf.h`)

```
01 int bpf_map_update_elem(int fd, const void *key, const void *value,
                                   __u64 flags);
02
03 int bpf_map_lookup_elem(int fd, const void *key, void *value);
04 int bpf_map_delete_elem(int fd, const void *key);
05 int bpf_map_get_next_key(int fd, const void *key, void *next_key);
06 int bpf_obj_pin(int fd, const char *pathname);
```

# Simple coding for an interactive program (`xdp1_user.c`)

```
01          if (load_bpf_file(filename)) {
02                  printf("%s", bpf_log_buf);
03                  return 1;
04          }
05
06          if (!prog_fd[0]) {
07                  printf("load_bpf_file: %s\n", strerror(errno));
08                  return 1;
09          }
10
11          signal(SIGINT, int_exit);
12
13          if (set_link_xdp_fd(ifindex, prog_fd[0]) < 0) {
14                  printf("link set xdp fd failed\n");
15                  return 1;
16          }
17
18          poll_stats(2);
```

# Slightly Different Implementation in this Case

- Skipped calling `load_bpf_file()` directly

- `load_bpf_file_fd()` was split into `load_bpf_elf_sections()` and `load_bpf_relocate_maps_and_attach()`

- Between two new functions called `bpf_obj_pin()` to create mapping between each sysfs file and fd for BPF map

- Allows access to eBPF maps by other processes via `bpf_obj_get()`

# Upon xdp_ddos01_blacklist_user exit

- All eBPF maps outlined in `xdp_ddos01_blacklist_kern.c` are pinned to files in `/sys/fs/bpf`
- Valid file descriptors for operating on those eBPF maps can be obtained by calling `bpf_obj_get()`
- eBPF program is attached to netdev and running!

# xdp_ddos01_blacklist_cmdline.c (and xdp_ddos01_blacklist_common.h)

# Command Line Tool Operations

- Print blacklist entries and historical stats

- Add/Delete IPv4 address from blacklist

- Add/Delete UDP and TCPs port from blacklist

- Print real-time XDP verdict stats

# Read Blacklist Entries

```
01          fd_blacklist = open_bpf_map(file_blacklist);
02          blacklist_list_all_ipv4(fd_blacklist);
03          close(fd_blacklist);
```

# Get file descriptor for map file

```
01 int open_bpf_map(const char *file)
02 {
03      int fd;
04
05      fd = bpf_obj_get(file);
06      if (fd < 0) {
07              printf("ERR: Failed to open bpf map file:%s err(%d):%s\n",
08                      file, errno, strerror(errno));
09              exit(EXIT_FAIL_MAP_FILE);
10      }
11      return fd;
12 }
```

# Find Each Key in the Hash

```c
01 static void blacklist_list_all_ipv4(int fd)
02 {
03        __u32 key = 0, next_key;
04        __u64 value;
05
06      while (bpf_map_get_next_key(fd, &key, &next_key) == 0) {
07              printf("%s", key ? "," : "" );
08              key = next_key;
09              value = get_key32_value64_percpu(fd, key);
10              blacklist_print_ipv4(key, value);
11        }
12      printf("%s", key ? "," : "");
13 }
```

# Lookup 32-bit key and Return 64-bit per CPU Sum

```
01 static __u64 get_key32_value64_percpu(int fd, __u32 key)
02 {
03        unsigned int nr_cpus = bpf_num_possible_cpus();
04        __u64 values[nr_cpus];
05        __u64 sum = 0;
06        int i;
07
08        if ((bpf_map_lookup_elem(fd, &key, values)) != 0) {
09                fprintf(stderr,
10                        "ERR: bpf_map_lookup_elem failed key:0x%X\n", key);
11                return 0;
12        }
13        for (i = 0; i < nr_cpus; i++) {
14                sum += values[i];
15        }
16        return sum;
17 }
```

# Print IPv4 address drop count

```
01 static void blacklist_print_ipv4(__u32 ip, __u64 count)
02 {
03         char ip_txt[INET_ADDRSTRLEN] = {0};
04
05         /* Convert IPv4 addresses from binary to text form */
06         if (!inet_ntop(AF_INET, &ip, ip_txt, sizeof(ip_txt))) {
07                 fprintf(stderr,
08                         "ERR: Cannot convert u32 IP:0x%X to IP-txt\n", ip);
09                 exit(EXIT_FAIL_IP);
10         }
11         printf("\n \"%s\" : %llu", ip_txt, count);
12 }
```

# Modify Blacklist

```
01          fd_blacklist = open_bpf_map(file_blacklist);
02          res = blacklist_modify(fd_blacklist, ip_string, action);
03          close(fd_blacklist);
```

# Add New IPv4 Address to Blacklist

```c
01 static int blacklist_modify(int fd, char *ip_string, uint action)
02 {
03         unsigned int nr_cpus = bpf_num_possible_cpus();
04         __u64 values[nr_cpus];
05         __u32 key;
06         int res;
07
08         memset(values, 0, sizeof(__u64) * nr_cpus);
09
10         /* Convert IP-string into 32-bit network byte-order value */
11 [...]
12         if (action == ACTION_ADD) {
13                 res = bpf_map_update_elem(fd, &key, values, BPF_NOEXIST);
14         } else if (action == ACTION_DEL) {
15                 res = bpf_map_delete_elem(fd, &key);
16         } else {
17                 return EXIT_FAIL_OPTION;
18         }
```

# XDP DDoS Blacklist Usage

https://github.com/netoptimizer/prototype-kernel

# Attaching eBPF Program to an Interface

```
01 # ./xdp_ddos01_blacklist --dev enp1s0f1d1
02 Documentation:
03   XDP: DDoS protection via IPv4 blacklist
04
05 This program loads the XDP eBPF program into the kernel.
06 Use the cmdline tool for add/removing source IPs to the blacklist
07 and read statistics.
08
09  - Attached to device:enp1s0f1d1 (ifindex:3)
10  - Blacklist      map file: /sys/fs/bpf/ddos_blacklist
11  - Verdict stats map file: /sys/fs/bpf/ddos_blacklist_stat_verdict
12  - Blacklist Port map file: /sys/fs/bpf/ddos_port_blacklist
13  - Verdict port stats map file: /sys/fs/bpf/ddos_port_blacklist_count_tcp
14  - Verdict port stats map file: /sys/fs/bpf/ddos_port_blacklist_count_udp
15 load_bpf_file: Success
16 # ip link show enp1s0f1d1
17 3: enp1s0f1d1: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 xdp qdisc mq state ...
18      link/ether 00:0a:f7:94:ef:2d brd ff:ff:ff:ff:ff:ff
```

# Blacklist configuration

```
01 # ./xdp_ddos01_blacklist_cmdline --add --ip 10.10.10.10
02 blacklist_modify() IP:10.10.10.10 key:0xA0A0A0A
03 # ./xdp_ddos01_blacklist_cmdline --list
04 {
05   "10.10.10.10" : 0,
06 }
07 # ./xdp_ddos01_blacklist_cmdline --del --ip 10.10.10.10
08 blacklist_modify() IP:10.10.10.10 key:0xA0A0A0A
09 # ./xdp_ddos01_blacklist_cmdline --list
10 {
11 }
```

# Automating Blacklist Configuration

- Fail2ban relies on external programs to block real hosts who attempt to access a system in a variety of ways

- Extremely easy to write a new module to call `xdp_ddos_blacklist` to initialize XDP and `xdp_ddos01_blacklist_cmdline` to add a host to blacklist

- Change needed committed here: https://github.com/gospo/fail2ban

# XDP DDos Blacklist Performance

# The Real Reason to use XDP

- CPU: Intel i7-6700K CPU @ 4.00GHz

- NIC: 50GbE Mellanox-CX4 (driver: mlx5)

- Single Stream UDP traffic from remote host

- Delivery to closed UDP port (after commit 9f2f27a9):

  - UdpNoPorts 3,143,931 pps (no iptables/netfilter)

- `iptables -t raw -I PREROUTING -p udp --dport 9 -j DROP`

  - Drop: 4,522,956 pps

- XDP blacklist:

  - Drop: 9,697,564 pps

# Can the DDoS Blacklist program protect us?

Simulate attack using pktgen_sample05_flow_per_thread.sh with 8 threads with 7 flows using blacklisted IPv4 addresses.

```
01 $ ./xdp_ddos01_blacklist_cmdline --stats
02
03 XDP_action    pps         pps-human-readable period/sec
04 XDP_ABORTED   0           0                  1.000089
05 XDP_DROP      30463237    30,463,237         1.000089
06 XDP_PASS      3438094     3,438,094          1.000089
07 XDP_TX        0           0                  1.000089
```

# Another Reason to use XDP

- CPU: Broadcom BCM5871x (Quad-core A57) @ 1.8GHz
- NIC 10GbE Embedded NIC (driver: bnxt_en)
- Single Stream UDP traffic from remote host
- Delivery to closed UDP port (after commit 9f2f27a9):
    - UdpNoPorts 666,823 pps (no iptables/netfilter)
- XDP blacklist:
    - Drop: 3,721,040 pps

# Tips and Tricks
# (Problems I encountered)

# Turn on in-kernel JIT

- If `perf top` or `perf report` indicate `__bpf_prog_run()` as the top consumer of cycles, check that `bpf_jit_enable=1` is set.

```
$ sysctl net/core/bpf_jit_enable=1
net.core.bpf_jit_enable = 1
```

- Reload all BPF/XDP programs that were running to have this sysctl setting take effect
- On both x86 and ARM64, close to double number of pps handled when in-kernel JIT is enabled

# ulimit Settings

- The eBPF maps uses locked memory, which is default very low. Your program likely need to increase resource limit RLIMIT_MEMLOCK see system call setrlimit(2) or run `ulimit -a` to see how to increase shell defaults.
- The bpf_create_map call will return errno EPERM (Operation not permitted) when the RLIMIT_MEMLOCK memory size limit is exceeded.

# Dumping maps with readelf or llvm-objdump

```
01 $ llvm-objdump -h xdp_ddos01_blacklist_kern.o
02
03 xdp_ddos01_blacklist_kern.o:  file format ELF64-BPF
04
05 Sections:
06 Idx Name             Size      Address            Type
07   0                  00000000 0000000000000000
08   1 .strtab          0000011d 0000000000000000
09   2 .text            000001a8 0000000000000000 TEXT DATA
10   3 .rel.text        00000030 0000000000000000
11   4 xdp_prog         00000360 0000000000000000 TEXT DATA
12   5 .relxdp_prog     00000050 0000000000000000
13   6 maps             00000064 0000000000000000 DATA
14   7 license          00000004 0000000000000000 DATA
15   8 .eh_frame        00000040 0000000000000000 DATA
16   9 .rel.eh_frame    00000020 0000000000000000
17  10 .symtab          00000258 0000000000000000
```

# printk debugging

```
01 #ifdef  DEBUG
02 /* Only use this for debug output. Notice output from bpf_trace_printk()
03  * end-up in /sys/kernel/debug/tracing/trace_pipe
04  */
05 #define bpf_debug(fmt, ...)                                          \
06                 ({                                                   \
07                         char ____fmt[] = fmt;                        \
08                         bpf_trace_printk(____fmt, sizeof(____fmt),   \
09                                         ##__VA_ARGS__);              \
10                 })
11 #else
12 #define bpf_debug(fmt, ...) { } while (0)
13 #endif
```

# Possibly Useful Hacks^W Ideas

# Whitelisting

```
# basically...
$ sed -i s/XDP_DROP/XDP_SAVEDROP/g *.[ch]
$ sed -i s/XDP_PASS/XDP_DROP/g *.[ch]
$ sed -i s/XDP_SAVEDROP/XDP_PASS/g *.[ch]
```

# Blacklisting subnets

- LPM map type (BPF_MAP_TYPE_LPM_TRIE) already implemented by Daniel Mack!
- Add new BPF map (xdp_ddos01_blacklist_kern.c)
- Add code in to lookup LPM entries for newly created map in function parse_ipv4() (xdp_ddos01_blacklist_kern.c)
- Add new command line options and JSON output (xdp_ddos01_blacklist_cmdline.c)

# IPv6 blacklisting

- New PERCPU Hash long enough to hold IPv6 addresses (xdp_ddos01_blacklist_kern.c)
- Create new function to check against new array similar to parse_ipv4() and call from handle_eth_protocol() (xdp_ddos01_blacklist_kern.c)
- Add new command line options and JSON output (xdp_ddos01_blacklist_cmdline.c)

# What's next for XDP DDoS Blacklist

# Keeping Maps During eBPF Program Attach/Detach

- Possible as long as application schema does not change in an incompatible manner
- Work from Daniel Borkmann in iproute2 codebase detects differences between map types and could be used to keep existing maps that are the same type across reloads

# --import Option to xdp_ddos01_blacklist_cmdline

- Parallel to the --list option
- Import/merge based on JSON output from previous run or from remote system

# Add support for custom processing pipeline

- Add a **BPF_MAP_TYPE_PROG_ARRAY** map that can be traversed with **bpf_tail_call()**
- Allow users to more easily create custom processing pipeline (i.e. use only IPv4 address filtering rather than checking addresses and TCP/UDP ports)

# References

- https://github.com/iovisor/bpf-docs/blob/master/Express_Data_Path.pdf
- https://prototype-kernel.readthedocs.io/en/latest/
- https://github.com/iovisor/bcc/blob/master/docs/kernel-versions.md
- https://www.slideshare.net/AlexeiStarovoitov/bpf-inkernel-virtual-machine
- https://www.slideshare.net/brendangregg/linux-bpf-superpowers
- man bpf
- http://www.tcpdump.org/papers/bpf-usenix93.pdf
- https://github.com/iovisor/bpf-docs/blob/master/bpf-internals-2.md
- http://people.netfilter.org/hawk/presentations/OpenSourceDays2017/XDP_DDoS_protecting_osd2017.pdf