

OVS收发包和流表卸载源码学习

流表模块数据结构与概要流程

数据结构

概要流程

收包模块

收包模块概要流程

用户态收包模块详细流程

flow到miniflow

flow和miniflow结构体介绍

flow压缩流程

三级流表

emc流表查找流程

datapath查找流程

openflow查找流程

流表模块卸载流程

OVS Conntrack和NAT

基础须知

key

action

命令行配置NAT

源码分析

数据结构

源码举例-命令行配置

源码举例-包触发

清除conntrack表项

发包模块

概要流程

日志模块

相关配置命令

后续计划

路由转发

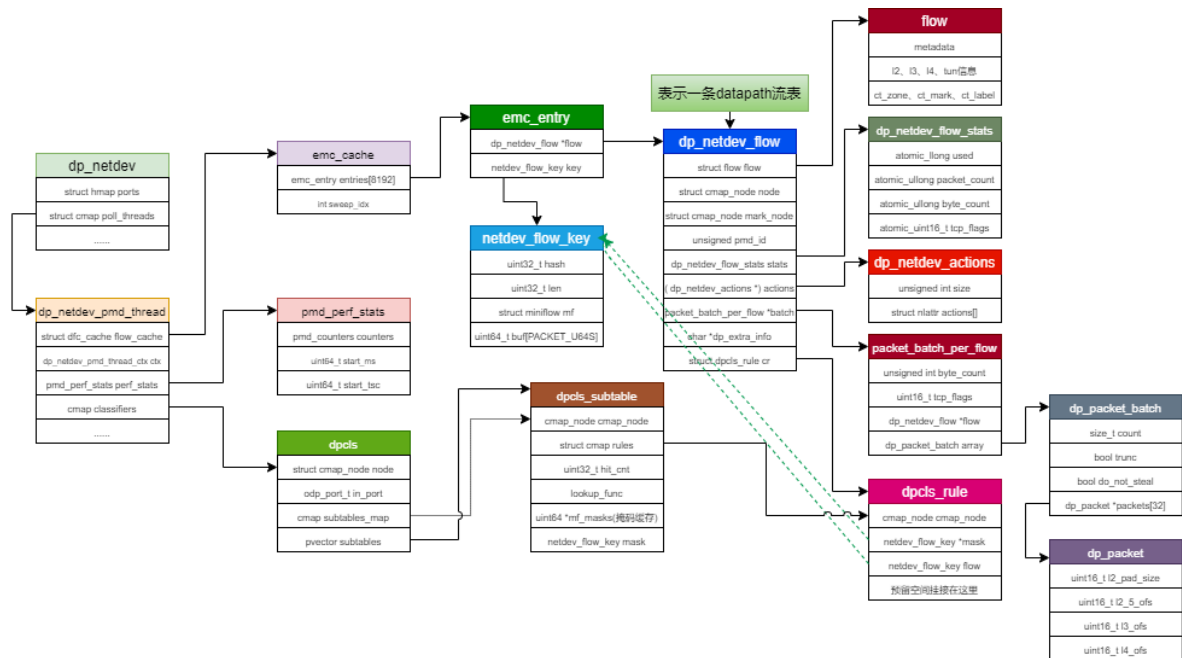
push eth

OVS收发包和流表卸载源码学习

流表模块数据结构与概要流程

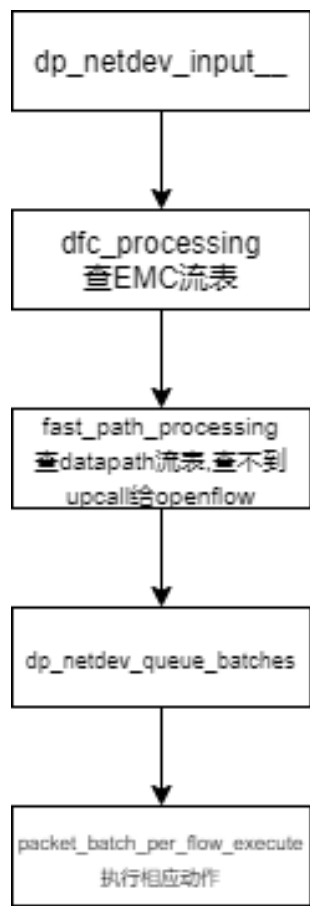
数据结构

看源码之前我们先学习下OVS的核心结构体之间的关系：



概要流程

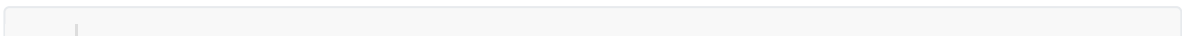
其主要流程如下图所示：



再明确一个结论：EMC是以PMD线程为边界的，每个PMD线程拥有自己的EMC；dpcls是以端口为边界的，每个端口拥有自己的dpcls；ofproto classifier是以桥为边界的，每个桥拥有自己的ofproto classifier。（后代码中均会有说明）

收包模块

收包模块概要流程





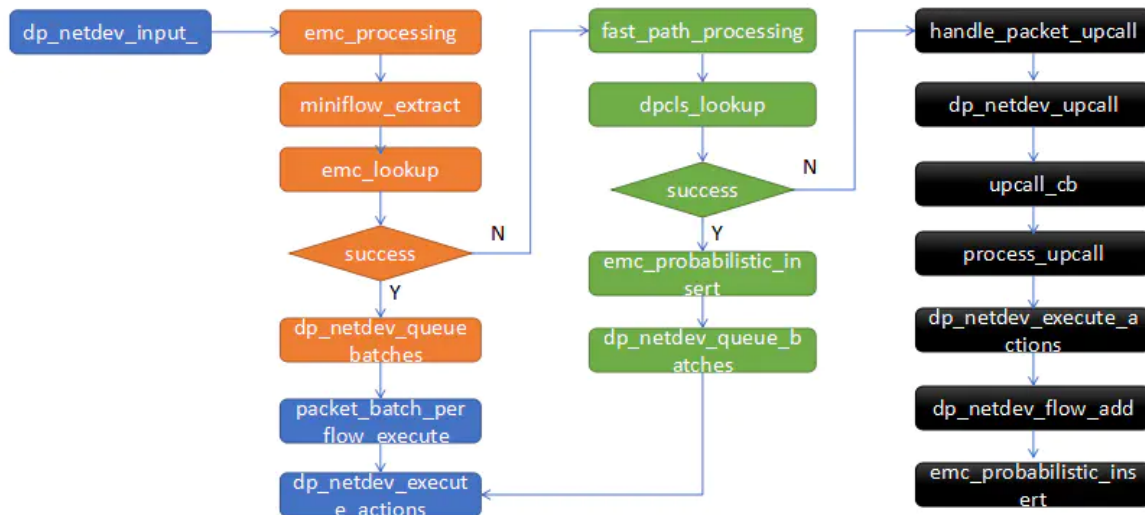
用户态收包模块详细流程

```

1 //主线程从非pmd类型端口收包
2 dpif_netdev_run
3 //pmd线程中从pmd类型端口收包
4 pmd_thread_main //会调用dfc_cache_init, 这也就是为什么emc是以pmd线程为边界的原因
5 dp_netdev_process_rxq_port -> dp_netdev_input -> dp_netdev_input__

```

收到包之后主要做三件事：`emc_processing` 查emc流表；`fast_path_processing` 查datapath流表和 `handle_packet_upcall` 查openflow流表已经后续的执行action。



flow到miniflow

flow和miniflow结构体介绍

介绍两个核心结构体：

```

1  //8字节对齐，总大小672个字节，负责存储流的所有信息。
2  struct flow {
3      /* Metadata */
4      struct flow_tnl tunnel; /* Encapsulating tunnel parameters. */
5      ovs_be64 metadata; /* OpenFlow Metadata. */
6      uint32_t regs[FLOW_N_REGS]; /* Registers. */
7      uint32_t skb_priority; /* Packet priority for QoS. */
8      uint32_t pkt_mark; /* Packet mark. */
9      uint32_t dp_hash; /* Datapath computed hash value. The exact
10                          * computation is opaque to the user space.
11                          */
12      union flow_in_port in_port; /* Input port.*/
13      uint32_t recirc_id; /* Must be exact match. */
14      uint8_t ct_state; /* Connection tracking state. */
15      uint8_t ct_nw_proto; /* CT orig tuple IP protocol. */
16      uint16_t ct_zone; /* Connection tracking zone. */
17      uint32_t ct_mark; /* Connection mark.*/
18      ovs_be32 packet_type; /* OpenFlow packet type. */
19      ovs_u128 ct_label; /* Connection label. */
20      uint32_t conj_id; /* Conjunction ID. */
21      ofp_port_t actset_output; /* Output port in action set. */
22
23      /* L2, Order the same as in the Ethernet header! (64-bit aligned) */
24      struct eth_addr dl_dst; /* Ethernet destination address. */
25      struct eth_addr dl_src; /* Ethernet source address. */
26      ovs_be16 dl_type; /* Ethernet frame type.
27                          * Note: This also holds the Ethertype for
28                          * packets of type PACKET_TYPE(1,
29                          * Ethertype) */
30      uint8_t pad1[2]; /* Pad to 64 bits. */
31      union flow_vlan_hdr vlans[FLOW_MAX_VLAN_HEADERS]; /* VLANs */
32      ovs_be32 mpls_label[ROUND_UP(FLOW_MAX_MPLS_LABELS, 2)]; /* MPLS label
33                          * stack
  
```

```

31                                                                 (with
padding). */
32     /* L3 (64-bit aligned) */
33     ovs_be32 nw_src;          /* IPv4 source address or ARP SPA. */
34     ovs_be32 nw_dst;          /* IPv4 destination address or ARP TPA. */
35     ovs_be32 ct_nw_src;       /* CT orig tuple IPv4 source address. */
36     ovs_be32 ct_nw_dst;       /* CT orig tuple IPv4 destination address.
*/
37     struct in6_addr ipv6_src; /* IPv6 source address. */
38     struct in6_addr ipv6_dst; /* IPv6 destination address. */
39     struct in6_addr ct_ipv6_src; /* CT orig tuple IPv6 source address. */
40     struct in6_addr ct_ipv6_dst; /* CT orig tuple IPv6 destination address.
*/
41     ovs_be32 ipv6_label;      /* IPv6 flow label. */
42     uint8_t nw_frag;          /* FLOW_FRAG_* flags. */
43     uint8_t nw_tos;           /* IP ToS (including DSCP and ECN). */
44     uint8_t nw_ttl;           /* IP TTL/Hop Limit. */
45     uint8_t nw_proto;         /* IP protocol or low 8 bits of ARP opcode.
*/
46     struct in6_addr nd_target; /* IPv6 neighbor discovery (ND) target. */
47     struct eth_addr arp_sha;    /* ARP/ND source hardware address. */
48     struct eth_addr arp_tha;    /* ARP/ND target hardware address. */
49     ovs_be16 tcp_flags;        /* TCP flags/ICMPv6 ND options type.
    * With L3 to avoid matching L4. */
50
51     ovs_be16 pad2;             /* Pad to 64 bits. */
52     struct ovs_key_nsh nsh;     /* Network Service Header keys */
53
54     /* L4 (64-bit aligned) */
55     ovs_be16 tp_src;           /* TCP/UDP/SCTP source port/ICMP type. */
56     ovs_be16 tp_dst;           /* TCP/UDP/SCTP destination port/ICMP code.
*/
57     ovs_be16 ct_tp_src;        /* CT original tuple source port/ICMP type.
*/
58     ovs_be16 ct_tp_dst;        /* CT original tuple dst port/ICMP code. */
59     ovs_be32 igmp_group_ip4;    /* IGMP group IPv4 address/ICMPv6 ND
reserved
60                                * field.
61                                * Keep last for BUILD_ASSERT_DECL below.
*/
62     ovs_be32 pad3;             /* Pad to 64 bits. */
63 };

```

下面介绍下miniflow，miniflow可以理解成压缩版的flow：

```

1  struct miniflow {
2      struct flowmap map;
3      /* flowmap用来记录和flow的对应关系，flowmap中的每一个bit对应struct flow中的一个
uint64_t字段，如果bit为1，则flow中对应的uint64_t字段为非0值，如果bit为0，则flow中对应的
uint64_t字段为0；另一部分是flowmap后面的内存，由调用者根据flowmap中bit为1的个数*8字
节申请内存，用来保存flow中非0的uint64_t。*/
4
5      /* 柔性数组后面接上：
    *      uint64_t values[n];
    * where 'n' is miniflow_n_values(miniflow). */
6
7  };
8

```

flow压缩流程

miniflow_extract用来从报文中提取flow信息，并保存到miniflow中。

```
1  /* Caller is responsible for initializing 'dst' with enough storage for
2   * FLOW_U64S * 8 bytes. */
3  void
4  miniflow_extract(struct dp_packet *packet, struct miniflow *dst)
5  {
6      const struct pkt_metadata *md = &packet->md;
7      const void *data = dp_packet_data(packet);
8      size_t size = dp_packet_size(packet);
9      ovs_be32 packet_type = packet->packet_type;
10     uint64_t *values = miniflow_values(dst);
11     struct mf_ctx mf = { FLOWMAP_EMPTY_INITIALIZER, values,
12                         values + FLOW_U64S };
13     const char *frame;
14     ovs_be16 dl_type = OVS_BE16_MAX;
15     uint8_t nw_frag, nw_tos, nw_ttl, nw_proto;
16     uint8_t *ct_nw_proto_p = NULL;
17     ovs_be16 ct_tp_src = 0, ct_tp_dst = 0;
18     ...
19     ...
20     //保存metadata信息到miniflow
21     if (md->skb_priority || md->pkt_mark) {
22         miniflow_push_uint32(mf, skb_priority, md->skb_priority);
23         miniflow_push_uint32(mf, pkt_mark, md->pkt_mark);
24     }
25     //保存md->dp_hash到miniflow
26     miniflow_push_uint32(mf, dp_hash, md->dp_hash);
27     //保存报文入端口到miniflow
28     miniflow_push_uint32(mf, in_port, odp_to_u32(md->in_port.odp_port));
29     ...
30     ...
31     //保存二层信息到miniflow
32     /* Link layer. */
33     ASSERT_SEQUENTIAL(dl_dst, dl_src);
34     miniflow_push_macs(mf, dl_dst, data);
35
36     /* VLAN */
37     union flow_vlan_hdr vlans[FLOW_MAX_VLAN_HEADERS];
38     size_t num_vlans = parse_vlan(&data, &size, vlans);
39
40     dl_type = parse_ethertype(&data, &size);
41     miniflow_push_be16(mf, dl_type, dl_type);
42     miniflow_pad_to_64(mf, dl_type);
43     if (num_vlans > 0) {
44         miniflow_push_words_32(mf, vlans, vlans, num_vlans);
45     }
46     ...
47     ...
48     //保存三层信息到miniflow
49     /* Push both source and destination address at once. */
50     miniflow_push_words(mf, nw_src, &nh->ip_src, 1);
51     ...
52     ...
53     //保存四层信息到miniflow
```

```

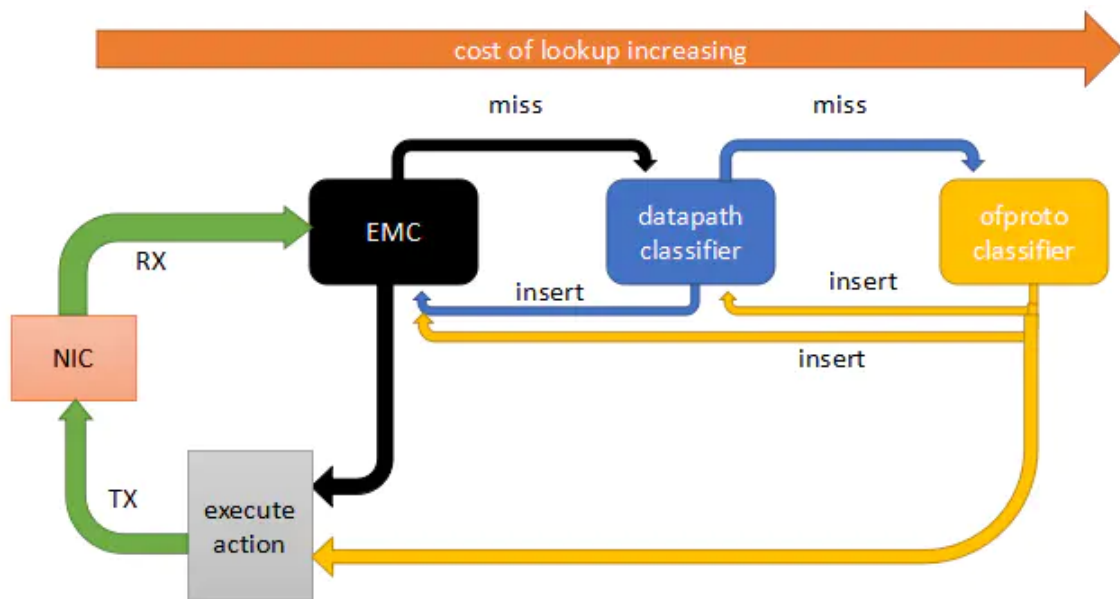
54     if (OVS_LIKELY(nw_proto == IPPROTO_TCP)) {
55         if (OVS_LIKELY(size >= TCP_HEADER_LEN)) {
56             const struct tcp_header *tcp = data;
57
58             miniflow_push_be32(mf, arp_tha.ea[2], 0);
59             miniflow_push_be32(mf, tcp_flags,
60                               TCP_FLAGS_BE32(tcp->tcp_ctl));
61             miniflow_push_be16(mf, tp_src, tcp->tcp_src);
62             miniflow_push_be16(mf, tp_dst, tcp->tcp_dst);
63             miniflow_push_be16(mf, ct_tp_src, ct_tp_src);
64             miniflow_push_be16(mf, ct_tp_dst, ct_tp_dst);
65         }
66     }

```

miniflow_expand用来将miniflow中的值恢复到flow结构体中,此处不做详细展开了。

三级流表

三级流表分别为EMC、datapath classifier、ofproto classifier。microflow在ovs+dpdk代码中, 又被称为EMC(exact match cache)。megaflow在ovs+dpdk代码中, 又被称为dpcls(datapath classifier)。从网卡接收到报文后, 首先查找EMC表项, 如果命中则直接执行action, 如果miss则查找dpcls。如果查找dpcls命中, 则将规则插入EMC, 并且执行action。还是miss的话, 就查找openflow流表。如果查找openflow命中, 则将规则插入dpcls和EMC, 并且执行action。还是miss的话, 就丢包或者发给controller。



三级流表中所使用到的hash、TSS算法我会抽时间单独写篇文章来介绍。本文主要介绍大致查找流程。

emc流表查找流程

```

1  static inline size_t
2  emc_processing(struct dp_netdev_pmd_thread *pmd,
3                 struct dp_packet_batch *packets_,
4                 struct netdev_flow_key *keys,
5                 struct packet_batch_per_flow batches[], size_t *n_batches,
6                 bool md_is_valid, odp_port_t port_no)
7  {
8      struct emc_cache *flow_cache = &pmd->flow_cache;
9      struct netdev_flow_key *key = &keys[0];
10     size_t n_missed = 0, n_dropped = 0;

```

```

11 struct dp_packet *packet;
12 const size_t size = dp_packet_batch_size(packets_);
13 uint32_t cur_min;
14 int i;
15 //获取插入emc的可能性, 如果为0, 说明emc流表被关闭
16 //此方案是为了解决EMC抖动问题, 会在流表算法文章中单独阐述这一内容。
17 //简而言之就是一个随机数获取函数 来实现概率性emc流表插入动作
18 atomic_read_relaxed(&pmd->dp->emc_insert_min, &cur_min);
19
20 //遍历报文
21 DP_PACKET_BATCH_REFILL_FOR_EACH (i, size, packet, packets_) {
22     struct dp_netdev_flow *flow;
23     //报文长度太小, 丢包
24     if (OVS_UNLIKELY(dp_packet_size(packet) < ETH_HEADER_LEN)) {
25         dp_packet_delete(packet);
26         n_dropped++;
27         continue;
28     }
29
30     if (i != size - 1) {
31         struct dp_packet **packets = packets_>packets;
32         /* 预取下一个包数据和metadata. */
33         OVS_PREFETCH(dp_packet_data(packets[i+1]));
34         pkt_metadata_prefetch_init(&packets[i+1]->md);
35     }
36     //初始化metadata
37     if (!md_is_valid) {
38         pkt_metadata_init(&packet->md, port_no);
39     }
40     //从报文提前flow信息, 保存到key->mf, 可以参考上面章节flow到miniflow
41     miniflow_extract(packet, &key->mf);
42     key->len = 0; /* Not computed yet. */
43     //获取hash值, 如果通过网卡的RSS收包则直接从mbuf获取hash值, 否则需要软件计算
44     hash key->hash = dpif_netdev_packet_get_rss_hash(packet, &key->mf);
45     //cur_min为0, 说明EMC是关闭的, 不用查询
46     flow = (cur_min == 0) ? NULL: emc_lookup(flow_cache, key);
47     if (OVS_LIKELY(flow)) {
48         //查到了EMC表项, 将报文插入flow的batch, 方便后面批量处理
49         dp_netdev_queue_batches(packet, flow, &key->mf, batches,
50                                 n_batches);
51     } else {
52         /* Exact match cache missed. Group missed packets together at
53          * the beginning of the 'packets' array. */
54         //没查到emc, 将报文重新插入packets_
55         dp_packet_batch_refill(packets_, packet, i);
56         /* 'key[n_missed]' contains the key of the current packet and
57          * must be returned to the caller. The next key should be
58          * to 'keys[n_missed + 1]'. */
59         key = &keys[++n_missed];
60     }
61 }
62 //增加命中emc的计数
63 dp_netdev_count_packet(pmd, DP_STAT_EXACT_HIT,
64                         size - n_dropped - n_missed);
65 //返回未命中emc并且没有被丢弃的报文个数

```



```

66     return dp_packet_batch_size(packets_);
67 }

```

datapath查找流程

下面再介绍下datapath的大致查找流程。

```

1  static inline void
2  fast_path_processing(struct dp_netdev_pmd_thread *pmd,
3                      struct dp_packet_batch *packets_,
4                      struct netdev_flow_key *keys,
5                      struct packet_batch_per_flow batches[], size_t
6                      *n_batches,
7                      odp_port_t in_port,
8                      long long now)
9  {
10     int cnt = packets_>count;
11     #if !defined(__CHECKER__) && !defined(_WIN32)
12     const size_t PKT_ARRAY_SIZE = cnt;
13     #else
14     /* Sparse or MSVC doesn't like variable length array. */
15     enum { PKT_ARRAY_SIZE = NETDEV_MAX_BURST };
16     #endif
17     struct dp_packet **packets = packets_>packets;
18     struct dpcls *cls;
19     struct dpcls_rule *rules[PKT_ARRAY_SIZE];
20     struct dp_netdev *dp = pmd->dp;
21     int miss_cnt = 0, lost_cnt = 0;
22     int lookup_cnt = 0, add_lookup_cnt;
23     bool any_miss;
24     size_t i;
25
26     for (i = 0; i < cnt; i++) {
27         /* Key length is needed in all the cases, hash computed on demand. */
28         keys[i].len =
29             netdev_flow_key_size(miniflow_n_values(&keys[i].mf));
30
31         /* Get the classifier for the in_port */
32         //根据入端口找到它的dpcls,这里也就是说为什么datapath是以端口为边界的
33         cls = dp_netdev_pmd_lookup_dpcls(pmd, in_port);
34         if (OVS_LIKELY(cls)) {
35             //调用dpcls_lookup进行匹配
36             any_miss = !dpcls_lookup(cls, keys, rules, cnt, &lookup_cnt);
37         } else {
38             any_miss = true;
39             memset(rules, 0, sizeof(rules));
40         }
41
42         //如果有miss的,则需要进行openflow流表查询
43         if (OVS_UNLIKELY(any_miss) && !fat_rwlock_tryrdlock(&dp-
44             >upcall_rwlock)) {
45             uint64_t actions_stub[512 / 8], slow_stub[512 / 8];
46             struct ofpbuf actions, put_actions;
47
48             ofpbuf_use_stub(&actions, actions_stub, sizeof actions_stub);
49             ofpbuf_use_stub(&put_actions, slow_stub, sizeof slow_stub);

```

```

47     for (i = 0; i < cnt; i++) {
48         struct dp_netdev_flow *netdev_flow;
49
50         if (OVS_LIKELY(rules[i])) {
51             continue;
52         }
53
54         /* It's possible that an earlier slow path execution installed
55         cheaper
56          * to catch it here than execute a miss. */
57         netdev_flow = dp_netdev_pmd_lookup_flow(pmd, &keys[i],
58                                                 &add_lookup_cnt);
59         if (netdev_flow) {
60             lookup_cnt += add_lookup_cnt;
61             rules[i] = &netdev_flow->cr;
62             continue;
63         }
64
65         miss_cnt++;
66         //开始查找openflow流表。
67         如果查找openflow流表成功并需要下发到dpcls时，需要判断是否超出最大流表限
68         制
69         handle_packet_upcall(pmd, packets[i], &keys[i], &actions,
70                             &put_actions, &lost_cnt, now);
71     }
72     ofpbuf_uninit(&actions);
73     ofpbuf_uninit(&put_actions);
74     fat_rwlock_unlock(&dp->upcall_rwlock);
75 } else if (OVS_UNLIKELY(any_miss)) {
76     for (i = 0; i < cnt; i++) {
77         if (OVS_UNLIKELY(!rules[i])) {
78             dp_packet_delete(packets[i]);
79             lost_cnt++;
80             miss_cnt++;
81         }
82     }
83 }
84
85 for (i = 0; i < cnt; i++) {
86     struct dp_packet *packet = packets[i];
87     struct dp_netdev_flow *flow;
88
89     if (OVS_UNLIKELY(!rules[i])) {
90         continue;
91     }
92
93     flow = dp_netdev_flow_cast(rules[i]);
94     //查找dpcls成功的，需要将相关rule下发到emc表项
95     emc_probabilistic_insert(pmd, &keys[i], flow);
96     dp_netdev_queue_batches(packet, flow, &keys[i].mf, batches,
97                             n_batches);
98     //统计信息
99     dp_netdev_count_packet(pmd, DP_STAT_MASKED_HIT, cnt - miss_cnt);
100    dp_netdev_count_packet(pmd, DP_STAT_LOOKUP_HIT, lookup_cnt);
101    dp_netdev_count_packet(pmd, DP_STAT_MISS, miss_cnt);

```

```

102     dp_netdev_count_packet(pmd, DP_STAT_LOST, lost_cnt);
103 }

```

openflow查找流程

如果EMC流表、datapath都查不到，则需要upcall调用查询openflow流表，内核情况下netlink实现，OVS+DPDK情况下，PMD线程中调用upcall_cb，通过回调函数调用的upcall_cb。

```

1  static int
2  upcall_cb(const struct dp_packet *packet, const struct flow *flow, ovs_u128
   *ufid,
3           unsigned pmd_id, enum dpif_upcall_type type,
4           const struct nlattr *userdata, struct ofpbuf *actions,
5           struct flow_wildcards *wc, struct ofpbuf *put_actions, void *aux)
6  {
7     static struct vlog_rate_limit rl = VLOG_RATE_LIMIT_INIT(1, 1);
8     struct udpif *udpif = aux;
9     struct upcall upcall;
10    bool megaflow;
11    int error;
12    //读取enable_megaflows，用来控制是否开启megaflow。
13    //可以通过命令开启 "ovs-appctl upcall/enable-megaflows"
14    atomic_read_relaxed(&enable_megaflows, &megaflow);
15
16    error = upcall_receive(&upcall, udpif->backer, packet, type, userdata,
17                          flow, 0, ufid, pmd_id);
18
19    if (error) {
20        return error;
21    }
22    //查找openflow流表，又是很长的函数调用
23    //后面单独讲解
24    error = process_upcall(udpif, &upcall, actions, wc);
25    if (error) {
26        goto out;
27    }
28
29    if (upcall.xout.slow && put_actions) {
30        ofpbuf_put(put_actions, upcall.put_actions.data,
31                  upcall.put_actions.size);
32    }
33
34    //如果关闭了megaflow，则将flow信息转换到wc，即megaflow也将变成精确匹配。
35    //如果开启megaflow，则wc查找openflow流表的通配符集合，小于等于flow信息。
36    if (OVS_UNLIKELY(!megaflow && wc)) {
37        flow_wildcards_init_for_packet(wc, flow);
38    }
39
40    //如果超过了最大流表限制，则返回ENOSPC
41    if (!should_install_flow(udpif, &upcall)) {
42        error = ENOSPC;
43        goto out;
44    }
45
46    if (upcall.ukey && !ukey_install(udpif, upcall.ukey)) {
47        VLOG_WARN_RL(&rl, "upcall_cb failure: ukey installation fails");
48        error = ENOSPC;
49    }
50
51    out:
52    return error;
53 }

```

```

49 out:
50     if (!error) {
51         upcall.ukey_persists = true;
52     }
53     upcall_uninit(&upcall);
54     return error;
55 }

```

这里单独介绍下process_upcall的主要流程:

```

1  //主要函数调用流程如下:
2  process_upcall -> upcall_xlate -> xlate_actions ->
   rule_dpif_lookup_from_table
3  //xbridge_lookup函数中就可以看出ofproto是基于bridge为边界的。
4  //xlate_actions函数就是先找到对应的bridge再去查流表
5
6  miss_config = OFPUTIL_TABLE_MISS_CONTINUE;
7  //从指定的table_id开始查找
8  for (next_id = *table_id;
9       next_id < ofproto->up.n_tables;
10      next_id++, next_id += (next_id == TBL_INTERNAL))
11  {
12      *table_id = next_id;
13      //到table的分类器中查找流表, 并设置wc
14      rule = rule_dpif_lookup_in_table(ofproto, version, next_id, flow,
   wc);
15      //如果查找到了, 则返回
16      if (rule) {
17          goto out; /* Match. */
18      }
19      //如果没有查找, 则根据配置选择继续查找下一个table还是结束查找
20      if (honor_table_miss) {
21          miss_config = ofproto_table_get_miss_config(&ofproto->up,
22                                                       *table_id);
23          if (miss_config == OFPUTIL_TABLE_MISS_CONTINUE) {
24              continue;
25          }
26      }
27      break;
28  }

```

流表模块卸载流程

```

1  流表卸载到驱动流程:
2
3      dpif_operate
4      (dpif-netdev) / \ (dpif-netlink)
5      / \
6      dpif_netdev_operate dpif_netlink_operate
7      | 用户态           | 内核态
8      dpif_netdev_flow_put try_send_to_netdev
9
10     |
11     |
12     queue_netdev_flow_put
13     | (回调)
14     netdev_offload_dpdk_flow_put
15
16     |
17     parse_flow_put
18     |
19     netdev_flow_put

```

```

13 |                                     | netdev_tc_flow_put
14 | netdev_offload_dpdk_flow_create
15 |                                     |
16 | netdev_offload_dpdk_mark_rss
17 |                                     |
18 | netdev_offload_dpdk_flow_create
19 |                                     |
20 | netdev_dpdk_rte_flow_create
21 |                                     |
22 | rte_flow_create
23 |                                     | 各厂商驱动回调
24 | sfc_flow_create
25 |
26 | 内核态:
27 | 当报文不匹配的时候, 会将报文上报, 会调用udpif_upcall_handler
28 | udpif_upcall_handler-->recv_upcalls-->handle_upcalls-->dpif_operate--
29 | >dpif_netlink_operate-->try_send_to_netdev-->parse_flow_put--
30 | >netdev_flow_put-->netdev_tc_flow_put
31 |
32 | 用户态:
33 | fast_path_processing->handle_packet_upcall->dp_netdev_flow_add-
34 | >queue_netdev_flow_put->dp_netdev_flow_offload_main-
35 | >dp_netdev_flow_offload_put--callback-->netdev_offload_dpdk_flow_put
36 |
37 | 内核态卸载是由handler线程处理, 而用户态OVS+DPDK是由dp_netdev_flow_offload线程处理,
38 | 此时的handler线程处于阻塞状态。

```

OVS Conntrack和NAT

基础须知

对于kernel datapath来说, 使用kernel的conntrack来实现, 对于userspace datapath来说, ovs本身来实现, 可参考 lib/conntrack.c 文件。

ct_state的可能值如下:

```

1 | new 通过ct action指定报文经过conntrack模块处理, 不一定有commit。通常是数据流的第一个数据包
2 | est 表示conntrack模块看到了报文双向数据流, 一定是在commit 的conntrack后
3 | rel 表示和已经存在的conntrack相关, 比如icmp不可达消息或者ftp的数据流
4 | rpl 表示反方向的报文
5 | inv 无效的, 表示conntrack模块没有正确识别到报文, 比如L3/L4 protocol handler没有加载, 或者L3/L4 protocol handler认为报文错误
6 | trk 表示报文经过了conntrack模块处理, 如果这个flag不设置, 其他flag都不能被设置
7 | snat 表示报文经过了snat, 源ip或者port
8 | dnat 表示报文经过了dnat, 目的ip或者port

```

这些flag得结合"+"或者 "-"来使用, "+"表示必须匹配, "-"表示必须不匹配。可以同时指定多个flag, 比如 ct_state=+trk+new。

key

匹配域和flow 中下的以下字段对应, 用来匹配流表:

```

1 | struct flow {
2 |     ...

```

```

3      uint8_t ct_state;           /* Connection tracking state. */
4      uint8_t ct_nw_proto;       /* CT orig tuple IP protocol. */
5      uint16_t ct_zone;          /* Connection tracking zone. */
6      uint32_t ct_mark;          /* Connection mark. */
7      ovs_be32 ct_nw_src;        /* CT orig tuple IPv4 source address. */
8      ovs_be32 ct_nw_dst;        /* CT orig tuple IPv4 destination address.
   */
9      struct in6_addr ct_ipv6_src; /* CT orig tuple IPv6 source address. */
10     struct in6_addr ct_ipv6_dst; /* CT orig tuple IPv6 destination address.
   */
11     ovs_be16 ct_tp_src;         /* CT original tuple source port/ICMP type.
   */
12     ovs_be16 ct_tp_dst;         /* CT original tuple dst port/ICMP code. */
13     ...
14 }

```

action

ovs通过ct action实现conntrack,ct会将报文送到conntrack模块进行处理,发送格式 `ct([argument][,argument...])`

支持参数如下:

```

1  commit //只有执行了commit, 才会在conntrack模块创建conntrack表项
2  force  //强制删除已存在的conntrack表项
3  table  //跳转到指定的table执行
4  zone   //设置zone, 隔离conntrack
5  exec   //执行其他action, 目前只支持设置ct_mark和ct_label, 比如exec(set_field: 1->ct_mark)
6         //alg=<ftp/tftp> 指定alg类型, 目前只支持ftp和tftp
7  nat    //指定ip和port

```

命令行配置NAT

```

1  #添加nat表项
2  ovs-ofctl add-flow br0 "table=0, priority=50, ct_state=-trk, tcp,
   in_port=veth_l0, actions=ct(commit,nat(src=10.1.1.240-10.2.2.2:2222-3333))"
3
4  //在一个ct里指定多次nat, 只有最后一个nat生效, 可参考do_xlate_actions中, ctx->
   >ct_nat_action = ofpact_get_NAT(a)只有一个ctx->ct_nat_action
5  ovs-ofctl add-flow br0 "table=0, priority=50, ct_state=-trk, tcp,
   actions=ct(commit,nat(src=10.1.1.240-10.2.2.2:2222-3333),
   nat(dst=10.1.1.240-10.2.2.2:2222-3333)), veth_r0"
6
7  //可以通过指定多个ct, 实现fullnat, 即同时转换源目的ip。
8  //但是这两个ct必须指定不同的zone, 否则只有第一个ct生效。因为在 handle_nat 中, 判断只有
   zone不一样才会进行后续的nat操作
9  //错误方式, 指定了src和dst nat, 但是zone相同, 只有前面的snat生效
10 ovs-ofctl add-flow br0 "table=0, priority=50, ct_state=-trk, tcp,
   actions=ct(commit,nat(src=10.1.1.240-10.2.2.2:2222-3333)),
   ct(commit,nat(dst=10.1.1.240-10.2.2.2:2222-3333)), veth_r0"
11
12 //正确方式, 使用不同zone, 指定fullnat
13 ovs-ofctl add-flow br0 "table=0, priority=50, ct_state=-trk, tcp,
   actions=ct(commit,zone=100, nat(src=10.1.1.240-10.2.2.2:2222-3333)),
   ct(commit, zone=200, nat(dst=10.1.1.240-10.2.2.2:2222-3333)), veth_r0"

```

源码分析

数据结构

主要存储CT信息的两个结构体如下：

```
1  /* OFPACT_CT.
2   *
3   * Used for NXAST_CT. */
4  struct ofpact_contrack {
5      OFPACT_PADDED_MEMBERS(
6          struct ofpact ofpact; /*eg:{ofpact = {type = OFPACT_CT, raw = 255
7          '\377', len为带上柔性数组后面内容的总长度} */
8          uint16_t flags; //NX_CT_F_COMMIT和NX_CT_F_FORCE
9          uint16_t zone_imm; //zone
10         struct mf_subfield zone_src;
11         uint16_t alg; //算法类型
12         uint8_t recirc_table; //跳转到指定table
13     );
14     struct ofpact actions[0];
15     /* struct ofpact_contrack用来保存ct后面的参数，并使用另一个结构体struct
16     ofpact_nat专门保存ct的nat信息。 */
17 };
18
19 struct ofpact_nat {
20     OFPACT_PADDED_MEMBERS(
21         struct ofpact ofpact; //类型为 OFPACT_NAT
22         uint8_t range_af; /* AF_UNSPEC, AF_INET, or AF_INET6 */
23         uint16_t flags; /* NX_NAT_F_* */
24         struct {
25             struct {
26                 uint16_t min;
27                 uint16_t max;
28             } proto;
29             union {
30                 struct {
31                     ovs_be32 min;
32                     ovs_be32 max;
33                 } ipv4;
34                 struct {
35                     struct in6_addr min;
36                     struct in6_addr max;
37                 } ipv6;
38             } addr;
39         } range;
40     );
41 };
```

源码举例-命令行配置

下面结合具体函数来进行代码流程分析,业务场景为命令行配置：

```
1  static char*
2  parse_CT(char *arg, const struct ofpact_parse_params *pp)
3  {
4      const size_t ct_offset = ofpacts_pull(pp->ofpacts);
```

```

5     struct ofpact_contrack *oc;
6     char *error = NULL;
7     char *key, *value;
8
9     //ofpact_put_CT可以参考
10    https://cloud.tencent.com/developer/article/1082684, 代码抽象的一种用法
11    //设置 OFPACT_CT
12    oc = ofpact_put_CT(pp->ofpacts);
13    oc->flags = 0;
14    oc->recirc_table = NX_CT_RECIRC_NONE;
15    while (ofputil_parse_key_value(&arg, &key, &value)) {
16        if (!strcmp(key, "commit")) {
17            oc->flags |= NX_CT_F_COMMIT;
18        } else if (!strcmp(key, "force")) {
19            oc->flags |= NX_CT_F_FORCE;
20        } else if (!strcmp(key, "table")) {
21            if (!ofputil_table_from_string(value, pp->table_map,
22                                          &oc->recirc_table)) {
23                error = xasprintf("unknown table %s", value);
24            } else if (oc->recirc_table == NX_CT_RECIRC_NONE) {
25                error = xasprintf("invalid table %#"PRIx8, oc->recirc_table);
26            }
27        } else if (!strcmp(key, "zone")) {
28            error = str_to_u16(value, "zone", &oc->zone_imm);
29
30            if (error) {
31                free(error);
32                error = mf_parse_subfield(&oc->zone_src, value);
33                if (error) {
34                    return error;
35                }
36            } else if (!strcmp(key, "alg")) {
37                error = str_to_connhelper(value, &oc->alg);
38            } else if (!strcmp(key, "nat")) {
39                const size_t nat_offset = ofpacts_pull(pp->ofpacts);
40                //解析NAT信息
41                error = parse_NAT(value, pp);
42                /* Update CT action pointer and length. */
43                pp->ofpacts->header = ofpbuf_push_uninit(pp->ofpacts, nat_offset);
44                oc = pp->ofpacts->header;
45            } else if (!strcmp(key, "exec")) {
46                /* Hide existing actions from ofpacts_parse_copy(), so the
47                 * nesting can be handled transparently. */
48                enum ofputil_protocol usable_protocols2;
49                const size_t exec_offset = ofpacts_pull(pp->ofpacts);
50
51                /* Initializes 'usable_protocols2', fold it back to
52                 * '*usable_protocols' afterwards, so that we do not lose
53                 * restrictions already in there. */
54                struct ofpact_parse_params pp2 = *pp;
55                pp2.usable_protocols = &usable_protocols2;
56                /* 解析 exec 参数, 比如 set_field ct(commit,exec(set_field:1->ct_mark)) (->后面的ct_mark为key, 前面的1为value) */
57                error = ofpacts_parse_copy(value, &pp2, false, OFPACT_CT);
58                *pp->usable_protocols &= usable_protocols2;

```



```

59     pp->ofpacts->header = ofpbuf_push_uninit(pp->ofpacts,
exec_offset);
60     oc = pp->ofpacts->header;
61 } else {
62     error = xasprintf("invalid argument to \"ct\" action: `%s'",
key);
63 }
64 if (error) {
65     break;
66 }
67 }
68 if (!error && oc->flags & NX_CT_F_FORCE && !(oc->flags &
NX_CT_F_COMMIT)) {
69     error = xasprintf("\"force\" flag requires \"commit\" flag.");
70 }
71
72 if (ofpbuf_oversized(pp->ofpacts)) {
73     free(error);
74     return xasprintf("input too big");
75 }
76 //更新 struct ofpact_contrack->ofpact.len, 包含nat的长度
77 ofpact_finish_CT(pp->ofpacts, &oc);
78 ofpbuf_push_uninit(pp->ofpacts, ct_offset);
79 return error;
80 }
81

```

该流程结束后大概内存格式如下: `struct ofpact_contrack(OPACT_CT) + struct ofpact_nat(OPACT_NAT) + struct ofpact_set_field(OPACT_SET_FIELD)`

源码举例-包触发

看到这里我们继续来分析下包触发CT流表的流程。ovs接收到报文后, 查找fastpath失败, 继续slowpath查找, 如果匹配到的流表的action为ct, 处理流程如下:

```

1  do_xlate_actions(const struct ofpact *ofpacts, size_t ofpacts_len,
2                  struct xlate_ctx *ctx, bool is_last_action,
3                  bool group_bucket_action)
4  {
5      struct flow_wildcards *wc = ctx->wc;
6      struct flow *flow = &ctx->xin->flow;
7      const struct ofpact *a;
8
9      OPACT_FOR_EACH (a, ofpacts, ofpacts_len) {
10         struct ofpact_controller *controller;
11         const struct ofpact_metadata *metadata;
12         const struct ofpact_set_field *set_field;
13         const struct mf_field *mf;
14         bool last = is_last_action && ofpact_last(a, ofpacts, ofpacts_len)
15                 && ctx->action_set.size;
16
17         //如果动作类型为CT, 调用对应处理函数
18         switch (a->type) {
19             case OPACT_CT:
20                 //ofpact_get_CT获取struct ofpact_contrack及其后面嵌套的action
21                 compose_contrack_action(ctx, ofpact_get_CT(a), last);
22             }
23     }
24 }

```

```

23     }
24 }

```

接下来是将struct ofpact_conntrack中的action信息转换到datapath能识别的aciton结构体。

```

1  static void
2  compose_conntrack_action(struct xlate_ctx *ctx, struct ofpact_conntrack
   *ofc)
3      //内部再次调用do_xlate_actions, 解析nat和ct_mark,ct_label信息
4      do_xlate_actions(ofc->actions, ofpact_ct_get_action_len(ofc), ctx);
5      //获取nat信息, 保存到 ctx->ct_nat_action, 如果指定了多次nat, 只有最后一次
   会生效
6      case OFPACT_NAT:
7          /* This will be processed by compose_conntrack_action(). */
8          ctx->ct_nat_action = ofpact_get_NAT(a);
9          break;
10
11     //解析 ct_mark 或者 ct_label 并保存到 flow->ct_mark和 flow->ct_label
12     case OFPACT_SET_FIELD:
13         set_field = ofpact_get_SET_FIELD(a);
14         mf = set_field->field;
15
16         /* Set the field only if the packet actually has it. */
17         if (mf_are_prereqs_ok(mf, flow, wc)) {
18             mf_mask_field_masked(mf, ofpact_set_field_mask(set_field),
   wc);
19             mf_set_flow_value_masked(mf, set_field->value,
20                                     ofpact_set_field_mask(set_field),
21                                     flow);
22         if (ofc->zone_src.field) {
23             zone = mf_get_subfield(&ofc->zone_src, &ctx->xin->flow);
24         } else {
25             zone = ofc->zone_imm;
26         }
27
28         //添加第一个 datapath action OVS_ACTION_ATTR_CT
29         //OVS_ACTION_ATTR_CT 开始
30         ct_offset = nl_msg_start_nested(ctx->odp_actions, OVS_ACTION_ATTR_CT);
31         if (ofc->flags & NX_CT_F_COMMIT) {
32             nl_msg_put_flag(ctx->odp_actions, ofc->flags & NX_CT_F_FORCE ?
   OVS_CT_ATTR_FORCE_COMMIT : OVS_CT_ATTR_COMMIT);
33             if (ctx->xbridge->support.ct_eventmask) {
34                 nl_msg_put_u32(ctx->odp_actions, OVS_CT_ATTR_EVENTMASK,
   OVS_CT_EVENTMASK_DEFAULT);
35             }
36         }
37
38         nl_msg_put_u16(ctx->odp_actions, OVS_CT_ATTR_ZONE, zone);
39         put_ct_mark(&ctx->xin->flow, ctx->odp_actions, ctx->wc);
40         if (wc->masks.ct_mark) {
41             struct {
42                 uint32_t key;
43                 uint32_t mask;
44             } *odp_ct_mark;
45
46             odp_ct_mark = nl_msg_put_unspec_uninit(odp_actions,
   OVS_CT_ATTR_MARK, sizeof(*odp_ct_mark));

```

```

47         odp_ct_mark->key = flow->ct_mark & wc->masks.ct_mark;
48         odp_ct_mark->mask = wc->masks.ct_mark;
49     }
50     put_ct_label(&ctx->xin->flow, ctx->odp_actions, ctx->wc);
51
52     put_ct_helper(ctx, ctx->odp_actions, ofc);
53     if (ofc->alg) {
54         switch(ofc->alg) {
55             case IPPORT_FTP:
56                 nl_msg_put_string(odp_actions, OVS_CT_ATTR_HELPER, "ftp");
57                 break;
58             case IPPORT_TFTP:
59                 nl_msg_put_string(odp_actions, OVS_CT_ATTR_HELPER,
60 "tftp");
61                 break;
62             default:
63                 xlate_report_error(ctx, "cannot serialize ct_helper %d",
64 ofc->alg);
65                 break;
66         }
67     }
68
69     put_ct_nat(ctx);
70     struct ofpact_nat *ofn = ctx->ct_nat_action;
71     nat_offset = nl_msg_start_nested(ctx->odp_actions,
72 OVS_CT_ATTR_NAT);
73     if (ofn->flags & NX_NAT_F_SRC || ofn->flags & NX_NAT_F_DST) {
74         nl_msg_put_flag(ctx->odp_actions, ofn->flags & NX_NAT_F_SRC
75 ? OVS_NAT_ATTR_SRC : OVS_NAT_ATTR_DST);
76     }
77     if (ofn->flags & NX_NAT_F_PERSISTENT) {
78         nl_msg_put_flag(ctx->odp_actions,
79 OVS_NAT_ATTR_PERSISTENT);
80     }
81     if (ofn->flags & NX_NAT_F_PROTO_HASH) {
82         nl_msg_put_flag(ctx->odp_actions,
83 OVS_NAT_ATTR_PROTO_HASH);
84     }
85     else if (ofn->flags & NX_NAT_F_PROTO_RANDOM) {
86         nl_msg_put_flag(ctx->odp_actions,
87 OVS_NAT_ATTR_PROTO_RANDOM);
88     }
89     ...
90 }
91 nl_msg_end_nested(ctx->odp_actions, nat_offset);
92
93 ctx->ct_nat_action = NULL;
94 //OVS_ACTION_ATTR_CT 结束
95 nl_msg_end_nested(ctx->odp_actions, ct_offset);
96 }
97
98 //如果配置 ct(table=x) 则需要添加第二个 datapath action
99 OVS_ACTION_ATTR_RECIRC
100 //recirc_table 值为table id, 表示需要转到其他table继续执行, 比如
101 actions=ct(table=0)
102 //值为 NX_CT_RECIRC_NONE, 说明不需要
103 //节点最后存储在ctx->xin->recirc_queue中
104 if (ofc->recirc_table == NX_CT_RECIRC_NONE) {

```

```

97     ctx->contracked = true;
98     compose_recirculate_and_fork(ctx, ofc->recirc_table);
99     uint32_t recirc_id;
100    ctx->freezing = true;
101    recirc_id = finish_freezing__(ctx, table);
102    struct frozen_state state = {
103        //保存需要跳转到的 table id, 即 recirc_table
104        .table_id = table,
105        .ofproto_uuid = ctx->xbridge->ofproto->uuid,
106        .stack = ctx->stack.data,
107        .stack_size = ctx->stack.size,
108        .mirrors = ctx->mirrors,
109        .contracked = ctx->contracked,
110        .xport_uuid = ctx->xin->xport_uuid,
111        .ofpacts = ctx->frozen_actions.data,
112        .ofpacts_len = ctx->frozen_actions.size,
113        .action_set = ctx->action_set.data,
114        .action_set_len = ctx->action_set.size,
115    };
116    frozen_metadata_from_flow(&state.metadata, &ctx->xin-
117>flow);
118    //获取 recirc_id, 保存到 odp_actions, 作为datapath的其中一个
119    action
120    id = recirc_alloc_id_ctx(&state);
121    uint32_t hash = frozen_state_hash(state);
122    struct recirc_id_node *node = recirc_ref_equal(state,
123hash);
124    node = recirc_alloc_id__(state, hash);
125    struct recirc_id_node *node = xzalloc(sizeof
126*node);
127    node->hash = hash;
128    ovs_refcount_init(&node->refcount);
129    frozen_state_clone(CONST_CAST(struct frozen_state
130*, &node->state), state);
131    cmap_insert(&id_map, &node->id_node, node->id);
132    cmap_insert(&metadata_map, &node->metadata_node,
133node->hash);
134    return node;
135    node->id;
136    nl_msg_put_u32(ctx->odp_actions, OVS_ACTION_ATTR_RECIRC,
137id);
138    ctx->contracked = false;
139    }

```

综上所述openflow中的action 可能包含两种action: OVS_ACTION_ATTR_CT和 OVS_ACTION_ATTR_RECIRC, 前者又包含了commit(OVS_CT_ATTR_FORCE_COMMIT), ct_mark(OVS_CT_ATTR_MARK), ct_label和nat(OVS_CT_ATTR_NAT)等信息, 后者仅仅包含了 recirc_id, 用来重新注入datapath后查看到table id, 即用来跳转到指定table执行。将action添加到 datapath后, 对当前报文执行action流程如下:

```

1  static inline void
2  packet_batch_per_flow_execute(struct packet_batch_per_flow *batch,
3                               struct dp_netdev_pmd_thread *pmd)
4  {
5      struct dp_netdev_actions *actions;
6      struct dp_netdev_flow *flow = batch->flow;

```

```

7         //获取action
8         actions = dp_netdev_flow_get_actions(flow);
9
10        //执行对应action
11        dp_netdev_execute_actions(pmd, &batch->array, true, &flow->flow,
actions->actions, actions->size, now);
12    }
13
14    //odp_execute_actions函数是具体执行流程函数
15    //遍历执行匹配流表的所有 actions
16    NL_ATTR_FOR_EACH_UNSAFE (a, left, actions, actions_len)
17        int type = nl_attr_type(a);
18        //调用dp_execute_cb, 如果是最后一个batch或者最后一个action 则执行完直
接返回
19        //具体操作见下述流程
20        dp_execute_action(dp, batch, a, may_steal);
21        //执行 ct action
22        case OVS_ACTION_ATTR_CT: {
23            NL_ATTR_FOR_EACH_UNSAFE (b, left, nl_attr_get(a),
nl_attr_get_size(a)) {
24                enum ovs_ct_attr sub_type = nl_attr_type(b);
25                switch(sub_type) {
26                    case OVS_CT_ATTR_FORCE_COMMIT:
27                        force = true;
28                        /* fall through. */
29                    case OVS_CT_ATTR_COMMIT:
30                        commit = true;
31                        break;
32                    case OVS_CT_ATTR_ZONE:
33                        zone = nl_attr_get_u16(b);
34                        break;
35                    case OVS_CT_ATTR_HELPER:
36                        helper = nl_attr_get_string(b);
37                        break;
38                    case OVS_CT_ATTR_MARK:
39                        setmark = nl_attr_get(b);
40                        break;
41                    case OVS_CT_ATTR_LABELS:
42                        setlabel = nl_attr_get(b);
43                        break;
44                    case OVS_CT_ATTR_EVENTMASK:
45                        /* silently ignored, as userspace datapath
does not generate
46                        * netlink events. */
47                        break;
48                    case OVS_CT_ATTR_NAT: {
49                        const struct nlattr *b_nest;
50                        unsigned int left_nest;
51                        bool ip_min_specified = false;
52                        bool proto_num_min_specified = false;
53                        bool ip_max_specified = false;
54                        bool proto_num_max_specified = false;
55                        memset(&nat_action_info, 0, sizeof
nat_action_info);
56                        nat_action_info_ref = &nat_action_info;
57
58                        NL_NESTED_FOR_EACH_UNSAFE (b_nest, left_nest,
b) {

```

```

59     enum ovs_nat_attr sub_type_nest =
nl_attr_type(b_nest);

60
61     switch (sub_type_nest) {
62     case OVS_NAT_ATTR_SRC:
63     case OVS_NAT_ATTR_DST:
64         nat_config = true;
65         nat_action_info.nat_action |=
66             ((sub_type_nest ==
OVS_NAT_ATTR_SRC)
67                 ? NAT_ACTION_SRC :
NAT_ACTION_DST);
68         break;
69     case OVS_NAT_ATTR_IP_MIN:
70         memcpy(&nat_action_info.min_addr,
71             nl_attr_get(b_nest),
72             nl_attr_get_size(b_nest));
73         ip_min_specified = true;
74         break;
75     case OVS_NAT_ATTR_IP_MAX:
76         memcpy(&nat_action_info.max_addr,
77             nl_attr_get(b_nest),
78             nl_attr_get_size(b_nest));
79         ip_max_specified = true;
80         break;
81     case OVS_NAT_ATTR_PROTO_MIN:
82         nat_action_info.min_port =
83             nl_attr_get_u16(b_nest);
84         proto_num_min_specified = true;
85         break;
86     case OVS_NAT_ATTR_PROTO_MAX:
87         nat_action_info.max_port =
88             nl_attr_get_u16(b_nest);
89         proto_num_max_specified = true;
90         break;
91         //persistent, hash和random在 userspace
datapath中没用到
92     case OVS_NAT_ATTR_PERSISTENT:
93     case OVS_NAT_ATTR_PROTO_HASH:
94     case OVS_NAT_ATTR_PROTO_RANDOM:
95         break;
96     case OVS_NAT_ATTR_UNSPEC:
97     case __OVS_NAT_ATTR_MAX:
98         OVS_NOT_REACHED();
99     }
100 }
101
102 if (ip_min_specified && !ip_max_specified) {
103     nat_action_info.max_addr =
nat_action_info.min_addr;
104 }
105 if (proto_num_min_specified &&
!proto_num_max_specified) {
106     nat_action_info.max_port =
nat_action_info.min_port;
107 }
108 if (proto_num_min_specified ||
proto_num_max_specified) {

```

```

109         if (nat_action_info.nat_action &
NAT_ACTION_SRC) {
110             nat_action_info.nat_action |=
NAT_ACTION_SRC_PORT;
111         } else if (nat_action_info.nat_action &
NAT_ACTION_DST) {
112             nat_action_info.nat_action |=
NAT_ACTION_DST_PORT;
113         }
114     }
115     break;
116 }
117 }
118 }
119
120     /* We won't be able to function properly in this case,
hence
121         * complain loudly. */
122     if (nat_config && !commit) {
123         static struct vlog_rate_limit rl =
VLOG_RATE_LIMIT_INIT(5, 5);
124         VLOG_WARN_RL(&rl, "NAT specified without
commit.");
125     }
126     //struct dp_netdev 是全局的，所以 dp->conntrack 也是全局
的，多个pmd共享dp->conntrack
127     struct dp_netdev *dp = pmd->dp;
128     conntrack_execute(&dp->conntrack, packets_, aux->flow-
>d1_type, force, commit, zone, setmark, setlabel, helper,
nat_action_info_ref);
129     //跳转到其他table执行
130     case OVS_ACTION_ATTR_RECIRC:
131         if (*depth < MAX_RECIRC_DEPTH) {
132             .....
133         }

```

清除conntrack表项

创建datapath时，会启动专门的线程clean_thread_main清除超期的conntrack表项：ct-

>clean_thread = ovs_thread_create("ct_clean", clean_thread_main, ct);

主要处理流程如下：

```

1  int
2  conntrack_execute(struct conntrack *ct, struct dp_packet_batch *pkt_batch,
3                   ovs_be16 d1_type, bool force, bool commit, uint16_t
zone,
4                   const uint32_t *setmark,
5                   const struct ovs_key_ct_labels *setlabel,
6                   const char *helper,
7                   const struct nat_action_info_t *nat_action_info)
8  for (size_t i = 0; i < cnt; i++) {
9      //从 pkts 中提取报文信息到 ct->key，并判断报文是否合法
10     if (!conn_key_extract(ct, pkts[i], d1_type, &ctx, zone))
11         ctx->key.zone = zone;
12         ctx->key.d1_type = d1_type;

```

```

13     extract_l3_ipv4(&ctx->key, 13, tail - (char *) 13, NULL,
!hwo1_good_l3_csum);
14     key->src.addr.ipv4 = ip->ip_src;
15     key->dst.addr.ipv4 = ip->ip_dst;
16     key->nw_proto = ip->ip_proto;
17     extract_l4(&ctx->key, 14, tail - 14, &ctx->icmp_related, 13,
!hwo1_good_l4_csum);
18     //计算 hash 值
19     ctx->hash = conn_key_hash(&ctx->key, ct->hash_basis);
20     {
21         //如果报文不合法, 则设置 CS_INVALID 后, 继续处理下一个报文
22         pkts[i]->md.ct_state = CS_INVALID;
23         write_ct_md(pkts[i], zone, NULL, NULL, NULL);
24         continue;
25     }
26     //开始处理合法报文
27     process_one(ct, pkts[i], &ctx, zone, force, commit, now, setmark,
setlabel, nat_action_info, helper);
28     struct conn *conn;
29     //根据 hash 值, 得出一个 hash 桶
30     unsigned bucket = hash_to_bucket(ctx->hash);
31     #define CONNTRACK_BUCKETS_SHIFT 8
32     #define CONNTRACK_BUCKETS (1 << CONNTRACK_BUCKETS_SHIFT)
33     //hash 桶大小 256
34     return (hash >> (32 - CONNTRACK_BUCKETS_SHIFT)) %
CONNTRACK_BUCKETS;
35
36     //根据 ctx->key 查找 conn, 如果是reply方向数据流, 则设置reply标志
37     conn_key_lookup(&ct->buckets[bucket], ctx, now);
38     uint32_t hash = ctx->hash;
39     struct conn *conn;
40     HMAP_FOR_EACH_WITH_HASH (conn, node, hash, &ctb-
>connections) {
41         if (!conn_key_cmp(&conn->key, &ctx->key)
&& !conn_expired(conn, now)) {
42             ctx->conn = conn;
43             ctx->reply = false;
44             break;
45         }
46         if (!conn_key_cmp(&conn->rev_key, &ctx->key)
&& !conn_expired(conn, now)) {
47             ctx->conn = conn;
48             ctx->reply = true;
49             break;
50         }
51     }
52     }
53     conn = ctx->conn;
54
55     /* Delete found entry if in wrong direction. 'force' implies
commit. */
56     if (conn && force && ctx->reply) {
57         conn_clean(ct, conn, &ct->buckets[bucket]);
58         conn = NULL;
59     }
60
61     bool create_new_conn = false;
62     struct conn conn_for_un_nat_copy;
63     conn_for_un_nat_copy.conn_type = CT_CONN_TYPE_DEFAULT;
64

```



```

65         bool ftp_ctl = is_ftp_ctl(pkt);
66
67         if (OVS_LIKELY(conn)) {
68             if (ftp_ctl) {
69                 /* Keep sequence tracking in sync with the source of
the
70                     * sequence skew. */
71                 if (ctx->reply != conn->seq_skew_dir) {
72                     handle_ftp_ctl(ct, ctx, pkt, conn, now,
CT_FTP_CTL_OTHER,
73                         !!nat_action_info);
74                     create_new_conn = conn_update_state(ct, pkt, ctx,
&conn, now,
75                         bucket);
76                 } else {
77                     create_new_conn = conn_update_state(ct, pkt, ctx,
&conn, now,
78                         bucket);
79                     handle_ftp_ctl(ct, ctx, pkt, conn, now,
CT_FTP_CTL_OTHER,
80                         !!nat_action_info);
81                 }
82             } else {
83                 create_new_conn = conn_update_state(ct, pkt, ctx,
&conn, now,
84                     bucket);
85             }
86             if (nat_action_info && !create_new_conn) {
87                 handle_nat(pkt, conn, zone, ctx->reply, ctx->icmp_related);
88             }
89
90             } else if (check_orig_tuple(ct, pkt, ctx, now, &bucket, &conn,
nat_action_info)) {
91                 create_new_conn = conn_update_state(ct, pkt, ctx, &conn,
now,
92                     bucket);
93             } else {
94                 if (ctx->icmp_related) {
95                     /* An icmp related conn should always be found; no new
connection is created based on an icmp related
96                     packet. */
97                     pkt->md.ct_state = CS_INVALID;
98                 } else {
99                     create_new_conn = true;
100                 }
101             }
102         }
103
104         if (OVS_UNLIKELY(create_new_conn)) {
105             conn = conn_not_found(ct, pkt, ctx, commit, now,
nat_action_info, &conn_for_un_nat_copy, helper, alg_exp);
106             unsigned bucket = hash_to_bucket(ctx->hash);
107             struct conn *nc = NULL;
108
109             //四层协议判断报文是否有效
110             if (!valid_new(pkt, &ctx->key))
111                 return 14_protos[key->nw_proto]->valid_new(pkt);
112             {

```

```

113         pkt->md.ct_state = CS_INVALID;
114         return nc;
115     }
116
117     //设置 CS_NEW
118     pkt->md.ct_state = CS_NEW;
119
120     //只有设置了 commit, 才会将conn添加到hash表
121     if (commit) {
122         //判断是否超过 conn 表项最大限制
123         unsigned int n_conn_limit;
124         atomic_read_relaxed(&ct->n_conn_limit,
125 &n_conn_limit);
126         if (atomic_count_get(&ct->n_conn) >= n_conn_limit)
127         {
128             COVERAGE_INC(conntrack_full);
129             return nc;
130         }
131
132         //创建新表项
133         nc = new_conn(&ct->buckets[bucket], pkt, &ctx-
134 >key, now);
135
136         struct conn *newconn;
137         //tcp_new_conn
138         newconn = 14_protos[key->nw_proto]-
139 >new_conn(ctb, pkt, now);
140
141         newconn->key = *key;
142         return newconn;
143
144         ctx->conn = nc;
145         nc->rev_key = nc->key;
146         //翻转key
147         conn_key_reverse(&nc->rev_key);
148
149         if (nat_action_info) {
150             nc->nat_info = xmemdup(nat_action_info, sizeof
151 *nc->nat_info);
152
153             if (alg_exp) {
154                 *conn_for_un_nat_copy = *nc;
155                 ct_rwlock_wrlock(&ct->resources_lock);
156                 //根据nat配置, 选择合适的ip和port
157                 bool nat_res = nat_select_range_tuple(ct,
158 nc, conn_for_un_nat_copy);
159
160                 bool new_insert =
161 nat_conn_keys_insert(&ct->nat_conn_keys, nat_conn, ct->hash_basis);
162                 //将 nat 的conn插入 nat_conn_keys
163                 hmap_insert(nat_conn_keys,
164 &nat_conn_key->node, nat_conn_key_hash);
165                 if (!nat_res) {
166                     goto nat_res_exhaustion;
167                 }
168                 /* Update nc with nat adjustments made to
169                  * conn_for_un_nat_copy by
170 nat_select_range_tuple(). */
171                 *nc = *conn_for_un_nat_copy;
172                 ct_rwlock_unlock(&ct->resources_lock);
173             }
174         }
175     }

```

```

162 //设置 conn_type 为 CT_CONN_TYPE_UN_NAT, 表示此表
    项需要nat
163 conn_for_un_nat_copy->conn_type =
    CT_CONN_TYPE_UN_NAT;
164 conn_for_un_nat_copy->nat_info = NULL;
165 conn_for_un_nat_copy->alg = NULL;
166 //将报文做nat转换
167 nat_packet(pkt, nc, ctx->icmp_related);
168 if (conn->nat_info->nat_action &
    NAT_ACTION_SRC) {
169     pkt->md.ct_state |= CS_SRC_NAT;
170     if (conn->key.dl_type ==
        htons(ETH_TYPE_IP)) {
171         struct ip_header *nh =
            dp_packet_l3(pkt);
172         packet_set_ipv4_addr(pkt, &nh-
            >ip_src, conn->rev_key.dst.addr.ipv4_aligned);
173     }
174     if (!related) {
175         pat_packet(pkt, conn);
176     }
177 } else if (conn->nat_info->nat_action &
    NAT_ACTION_DST) {
178     pkt->md.ct_state |= CS_DST_NAT;
179 }
180 }
181 //将新建表项插入hash表
182 hmap_insert(&ct->buckets[bucket].connections, &nc-
    >node, ctx->hash);
183 //增加表项个数
184 atomic_count_inc(&ct->n_conn);
185 }
186 return nc;
187 }
188
189 write_ct_md(pkt, zone, conn, &ctx->key, alg_exp);
190 pkt->md.ct_state |= CS_TRACKED;
191 pkt->md.ct_zone = zone;
192 pkt->md.ct_mark = conn ? conn->mark : 0;
193 pkt->md.ct_label = conn ? conn->label : OVS_U128_ZERO;
194
195 pkt->md.ct_orig_tuple_ipv6 = false;
196 if (key) {
197     if (key->dl_type == htons(ETH_TYPE_IP)) {
198         //ct_orig_tuple 保存原始报文(第一次进ct模块时)的五元组信
            息
199         pkt->md.ct_orig_tuple.ipv4 = (struct
            ovs_key_ct_tuple_ipv4) {
200             key->src.addr.ipv4_aligned,
201             key->dst.addr.ipv4_aligned,
202             key->nw_proto != IPPROTO_ICMP
203             ? key->src.port : htons(key->src.icmp_type),
204             key->nw_proto != IPPROTO_ICMP
205             ? key->dst.port : htons(key->src.icmp_code),
206             key->nw_proto,
207         };
208     }
209 }

```

```

210
211         if (conn && setmark) {
212             set_mark(pkt, conn, setmark[0], setmark[1]);
213         }
214
215         if (conn && setlabel) {
216             set_label(pkt, conn, &setlabel[0], &setlabel[1]);
217         }
218     }
219 }

```

发包模块

概要流程

```

1 kernel:
2 netdev_frame_hook-->netdev_port_receive-->ovs_vport_receive--
>ovs_dp_process_packet-->ovs_execute_actions-->do_execute_actions--
>do_output
3 这里会根据dev的mtu和报文的mru比较需要不需要分片，通常情况下报文的mru是0，不需要分片，但是
  因为重组函数handle_fragments重组时设置了报文的mru，所以分片时也根据此项来判定。调用
  ovs_fragment分片完成之后会调用ovs_vport_send-->dev_queue_xmit--
  >dev_hard_start_xmit-->ops->ndo_start_xmit-->bond_start_xmit进行报文发送。
4
5 另一种的报文分片是VM中设置了TSO、GSO、UFO等features，导致报文比较大，这个时候需要根据标
  志判定是否需要分片。
6
7 netdev_frame_hook-->netdev_port_receive-->ovs_vport_receive--
>ovs_dp_process_packet-->ovs_execute_actions-->do_execute_actions--
>do_output-->ovs_vport_send-->dev_queue_xmit-->dev_hard_start_xmit这里会根据
  gso进行分片，分片完成后继续调用ops->ndo_start_xmit-->bond_start_xmit进行报文发送。
  其实此处也会根据其他的一些特性来判定，一般比如带了TSO之类的需要网卡去分片。
8
9 OVS-DPDK:
10 pmd_thread_main-->dp_netdev_process_rxq_port-->dp_netdev_input--
  >dp_netdev_input__-->packet_batch_per_flow_execute--
  >dp_netdev_execute_actions-->dp_execute_cb-->netdev_send--
  >netdev_dpdk_eth_send

```

日志模块

相关配置命令

三种输出方式: syslog, console, file

查看log使能情况命令: `ovs-appctl -t ovsdb-server vlog/list`

设置所有模块log等级命令: `ovs-appctl -t ovsdb-server vlog/set dbg`

设置指定模块log等级命令, 如jsonrpc: `ovs-appctl -t ovsdb-server vlog/set jsonrpc,dbg`

后续计划

写后续计划的目的是为了让大家督促我继续总结学习文档，毕竟话撻这了，不做到也挺打脸的。后面计划是更新一篇DPDK相关内容和一篇流表算法内容。

push eth

```

1      fast_path_processing
2
3      dp_netdev_execute_actions
4
5      odp_execute_actions
6
7      |
8      dpif_execute
9
10
11
12     dpif_get_actions
13         /* Lookup actions in userspace cache. */
14         //upcall->ufid 流唯一的id
15         //upcall->pmd_id datapath pmd id
16         struct udpif_key *ukey = ukey_lookup(udpif, upcall->ufid,
17                                             upcall->pmd_id);
18         if (ukey) {
19             ukey_get_actions(ukey, actions, &actions_len);
20         }
21
22     udpif_revalidator
23         |
24     revalidate
25         |
26     ukey_acquire
27         |
28     ukey_create_from_dpif_flow
29         |
30     ukey_create__
31
32
33     xlate_actions
34
35

```

```

1  虚拟化笔记，内核态：
2  https://zhuanlan.zhihu.com/p/337143779
3
4  main
5      ↳bridge_run
6          ↳bridge_reconfigure
7              ↳ofproto_create
8                  ↳construct
9                      ↳open_dpif_backer
10                         ↳udpif_create
11                            |   ↳dpif_register_upcall_cb//kernel为NULL，用于DPDK
12                            ↳udpif_set_threads
13                            |   ↳dpif_handlers_set
14                            ↳udpif_start_threads

```

```

15         ↳udpip_upcall_handler
16         ↳dpif_recv
17         |   ↳dpif_netlink_recv
18         |       ↳dpif_netlink_recv__
19         |           ↳parse_odp_packet
20         ↳recv_upcalls
21
22 #linux kernel datapath
23 ovs_vport_receive
24     ↳ovs_flow_key_extract
25     |   ↳key_extract
26     ↳ovs_dp_process_packet
27         ↳if ovs_flow_tbl_lookup_stats
28         |   ↳ovs_dp_upcall
29         |       ↳queue_userspace_packet
30         |           ↳ovs_nla_put_key
31         |               ↳__ovs_nla_put_key
32     ↳else ovs_execute_actions
33         ↳case OVS_ACTION_ATTR_OUTPUT
34         ↳do_output
35         ↳ovs_vport_send
36
37
38 recv_upcalls
39     ↳odp_flow_key_to_flow
40     ↳upcall_receive
41     |   ↳xlate_lookup
42     |       ↳xlate_lookup_ofproto_
43     ↳process_upcall
44     |   ↳upcall_xlate
45     |       ↳xlate_in_init
46     |       ↳xlate_actions
47     |           ↳rule_dpif_lookup_from_table
48     |           |   ↳rule_dpif_lookup_in_table
49     |           ↳rule_get_actions
50     |           ↳do_xlate_actions
51     |           |   ↳freeze_unroll_actions
52     |           |       ↳freeze_put_unroll_xlate
53     |           ↳xlate_group_action
54     |           |   ↳pick_select_group//报文第一次找不到bucket
55     |           |       ↳pick_dp_hash_select_group
56     |           |           ↳ctx_trigger_recirculate_with_hash
57     |           |   ↳xlate_group_bucket//第二次找到bucket才执行这里
58     |           |       ↳ofpacts_execute_action_set
59     |           |           ↳do_xlate_actions//相当于递归了
60     |           ↳finish_freezing//报文第一次来执行这里
61     |           |   ↳xlate_commit_actions
62     |           |       ↳commit_odp_actions
63     |           ↳finish_freezing__
64     |           |   ↳recirc_alloc_id_ctx
65     |           |       ↳把OVS_ACTION_ATTR_HASH和OVS_ACTION_ATTR_RECIRC下给datapath
66     |           ↳ctx_cancel_freeze
67     ↳ukey_create_from_upcall
68     ↳handle_upcalls
69     ↳dpif_operate
70     |   ↳ukey_install
71     ↳dpif_operate//给datapath安装流表 and 把skb发回内核
72     |   ↳if dpif_netlink_operate

```

```

73         ↳dpif_execute_with_help
74         ↳odp_execute_actions
75         ↳计算hash
76         ↳dpif_execute_helper_cb
77         ↳dpif_operate
78
79
80 内核datapath流表老化:
81  udpif_start_threads
82     ↳udpif_revalidator
83         ↳revalidate
84         |   ↳dpif_flow_dump_next
85         |   ↳udpif_get_n_flows
86         |   ↳delete_op_init__
87         |   ↳push_dp_ops
88     ↳revalidator_sweep

```

```

1  用户态:
2  https://blog.csdn.net/qq\_20817327/article/details/106761936
3
4  dp_execute_cb
5  如果是OVS_ACTION_ATTR_OUTPUT，调用dp_netdev_lookup_port查找端口，
6  然后调用netdev_send进行报文发送。
7
8  如果是OVS_ACTION_ATTR_TUNNEL_PUSH，调用push_tnl_action进行tunnel封装，然后调用
9  dp_netdev_recirculate->dp_netdev_input__重新查表操作。
10
11  如果是OVS_ACTION_ATTR_TUNNEL_POP，调用netdev_pop_header解封装，
12  然后调用dp_netdev_recirculate->dp_netdev_input__重新查表操作。
13
14
15

```

```

1  struct nm_rule_action {
2      uint64_t flag; /* action flag, eg:set ipv4 src/redirect */
3      uint32_t drop_flag; /* drop or forward */
4      uint32_t soft_id; /* an unique value for this rule todo*/
5      uint32_t mark_id;
6      uint32_t counter_id;
7      uint8_t outer_type; //outer_tnl_type, eg:tunnel
8      uint8_t queue; /* assigned rx queue */
9      uint8_t queues_cnt;
10     uint8_t action_cnt; /* different action type total cnt */
11
12     uint16_t metadata; /* secondary flow table need */
13
14     /* set ops */
15     struct nm_fdir_l4 l4_outer;
16     struct nm_fdir_l2 l2_data_outer;
17     struct nm_fdir_l3 ip_outer;
18
19     /* push ops */
20     struct nm_flow_item_vlan vlan;
21     struct nm_flow_item_vxlan vxlan;
22     uint8_t set_flag;
23

```

```

24     uint8_t buf[NM_FLOW_ACTION_SET_OUTER_HEADER_MAX_LEN];
25     uint32_t l2_cnt:2;
26     uint32_t l3_cnt:2;
27     uint32_t l4_cnt:2;
28     uint32_t tn1_cnt:2;
29     uint32_t vlan_cnt:2;
30     uint32_t rsv:22;
31 };
32
33
34 /* SPDX-License-Identifier: BSD-3-Clause
35  * Copyright(c) 2021-2024 Nebula Limited.
36  */
37
38 #ifndef _NM_FLOW_H_
39 #define _NM_FLOW_H_
40
41 #define DESC (x)                1
42 #define FALSE                    0
43 #define TRUE                     1
44 #define ERR                      4
45 #define WARN                     3
46 #define NOTICE                 2
47 #define INFO                    1
48
49 #define NM_RTE_ETHER_TYPE_IPV4    4
50 #define NM_RTE_ETHER_TYPE_IPV6    6
51
52 /* 14 outer tn1 hash tab entries: 1024 */
53 #define NM_MAX_L4_OUTER_TNL_FILTER_NUM    (1<<10)
54 #define NM_MAX_L2_ETHTYPE_TNL_FILTER_NUM    (1<<10)
55 #define NM_MAX_L3_L2_TNL_FILTER_NUM    (1<<10)
56 #define NM_MAX_ACT_CNT_MAX_NUM    (1<<17) /* 131072 */
57
58
59 #define NM_FLOW_TABLE_LENGTH        3
60 #define NM_FLOW_TABLE_VXLAN_LENGTH    3
61
62 #define NM_FLOW_TAB_PATTERN_TIMES    1
63 #define NM_FLOW_TABLE_IPV4_DEFAULT_MASK    0xFFFFFFFF
64 #define NM_FLOW_TABLE_L4_PORT_DEFAULT_MASK    0xFFFF
65 #define NM_MAX_FILTER_ID            0x0FFF
66 #define NM_FLOW_TABLE_FULL_MASK_AS_U8    0xFF
67 #define NM_FLOW_ACTION_SET_OUTER_HEADER_MAX_LEN    128
68 #define NM_FLOW_ACTION_VLAN_CNT        2
69 #define NM_ETH_MAC_LEN                6
70 #define NM_ACTION_DEL_MIN_REF_CNT    1
71
72 /* input set */
73 #define NM_INSET_NONE                0ULL
74
75 #define nm_malloc(h, s)    rte_zmalloc(NULL, s, 0)
76 #define nm_calloc(h, c, s)    rte_zmalloc(NULL, (c) * (s), 0)
77 #define nm_free(h, m)        rte_free(m)
78
79 #define NM_IPV6_ADDR_LEN_AS_U32    4
80 #define NM_IPV6_ADDR_LEN_AS_U8    16
81

```



```

82  /* Field range used to indicate pattern */
83  #define NM_PROT_MAC_INNER          (1ULL << 1)
84  #define NM_PROT_MAC_OUTER          (1ULL << 2)
85  #define NM_PROT_VLAN_INNER         (1ULL << 3)
86  #define NM_PROT_VLAN_OUTER         (1ULL << 4)
87  #define NM_PROT_IPV4_INNER         (1ULL << 5)
88  #define NM_PROT_IPV4_OUTER         (1ULL << 6)
89  #define NM_PROT_IPV6_INNER         (1ULL << 7)
90  #define NM_PROT_IPV6_OUTER         (1ULL << 8)
91  #define NM_PROT_TCP_INNER          (1ULL << 9)
92  #define NM_PROT_TCP_OUTER          (1ULL << 10)
93  #define NM_PROT_UDP_INNER          (1ULL << 11)
94  #define NM_PROT_UDP_OUTER          (1ULL << 12)
95  #define NM_PROT_SCTP_INNER         (1ULL << 13)
96  #define NM_PROT_SCTP_OUTER         (1ULL << 14)
97  #define NM_PROT_VXLAN              (1ULL << 15)
98  #define NM_PROT_NVGRE              (1ULL << 16)
99  #define NM_PROT_GTPU               (1ULL << 17)
100 #define NM_PROT_PPPOE_S             (1ULL << 18)
101 #define NM_PROT_ESP                 (1ULL << 19)
102 #define NM_PROT_AH                  (1ULL << 20)
103 #define NM_PROT_L2TPV3OIP          (1ULL << 21)
104 #define NM_PROT_PFCP                (1ULL << 22)
105
106 //action flag
107 #define NM_FLOW_ACTION_DROP          (1ULL << 1)
108 #define NM_FLOW_ACTION_REDIRECT      (1ULL << 2)
109 #define NM_FLOW_ACTION_MIRRED        (1ULL << 3)
110 #define NM_FLOW_ACTION_VLAN_PUSH     (1ULL << 4)
111 #define NM_FLOW_ACTION_VLAN_POP      (1ULL << 5)
112 #define NM_FLOW_ACTION_TUNNEL_ENCAP  (1ULL << 5)
113 #define NM_FLOW_ACTION_TUNNEL_DECAP  (1ULL << 6)
114 #define NM_FLOW_ACTION_COUNTER        (1ULL << 7)
115 #define NM_FLOW_ACTION_SET_IPV4_SRC_IP (1ULL << 8)
116 #define NM_FLOW_ACTION_SET_IPV4_DST_IP (1ULL << 9)
117 #define NM_FLOW_ACTION_SET_IPV6_SRC_IP (1ULL << 10)
118 #define NM_FLOW_ACTION_SET_IPV6_DST_IP (1ULL << 11)
119 #define NM_FLOW_ACTION_SET_SRC_MAC    (1ULL << 12)
120 #define NM_FLOW_ACTION_SET_DST_MAC    (1ULL << 13)
121 #define NM_FLOW_ACTION_SET_SRC_PORT   (1ULL << 14)
122 #define NM_FLOW_ACTION_SET_DST_PORT   (1ULL << 15)
123 #define NM_FLOW_ACTION_SET_TTL        (1ULL << 16)
124 #define NM_FLOW_ACTION_SET_DSCP       (1ULL << 17)
125 #define NM_FLOW_ACTION_RSS            (1ULL << 18)
126 #define NM_FLOW_ACTION_QUEUE          (1ULL << 19)
127 #define NM_FLOW_ACTION_MARK           (1ULL << 20)
128 #define NM_FLOW_ACTION_PUSH_VLAN      (1ULL << 21)
129 #define NM_FLOW_ACTION_POP_VLAN       (1ULL << 22)
130 #define NM_FLOW_ACTION_METADATA_FLAG  (1ULL << 23)
131
132
133 // action:vlan encap, set_flag
134 #define NM_FLOW_ACTION_ETH_FLAG       (1ULL << 1)
135 #define NM_FLOW_ACTION_IPV4_FLAG      (1ULL << 2)
136 #define NM_FLOW_ACTION_IPV6_FLAG      (1ULL << 3)
137 #define NM_FLOW_ACTION_TCP_FLAG       (1ULL << 4)
138 #define NM_FLOW_ACTION_UDP_FLAG       (1ULL << 5)
139 #define NM_FLOW_ACTION_VLAN_FLAG      (1ULL << 7)

```

```

140 #define NM_FLOW_ACTION_VXLAN_FLAG (1ULL << 8)
141
142
143 enum nm_flow_service_type {
144     NM_FLOW_FILTER_NONE = 0,
145     NM_L4_UPPER_TUNNEL_OUTER,
146     NM_L4_UPPER_OUTER_AND_INNER_ETHTYPE,
147     NM_DOWNER_NON_IP,
148     NM_DOWNER_IP,
149     NM_DOWNER_IP_AND_SECGRP,
150     NM_SECURITY_GROUP,
151     NM_UPPER_L2_MULTICAST_VLAN,
152     NM_UPPER_L2_MULTICAST,
153     NM_UPPER_L3_MULTICAST,
154     NM_DOWNER_MULTICAST,
155 };
156
157 enum nm_fltr_ptype {
158     /* NONE - used for undef/error */
159     NM_FLTR_PTYPE_NONF_NONE = 0,
160     NM_FLTR_PTYPE_NONF_IPV4_UDP,
161     NM_FLTR_PTYPE_NONF_IPV4_TCP,
162     NM_FLTR_PTYPE_NONF_IPV4_SCTP,
163     NM_FLTR_PTYPE_NONF_IPV4_OTHER,
164     NM_FLTR_PTYPE_NONF_IPV4_GTPU_IPV4_UDP,
165     NM_FLTR_PTYPE_NONF_IPV4_GTPU_IPV4_TCP,
166     NM_FLTR_PTYPE_NONF_IPV4_GTPU_IPV4_ICMP,
167     NM_FLTR_PTYPE_NONF_IPV4_GTPU_IPV4_OTHER,
168     NM_FLTR_PTYPE_NONF_IPV6_GTPU_IPV6_OTHER,
169     NM_FLTR_PTYPE_NONF_IPV4_GTPU_EH_IPV4_OTHER,
170     NM_FLTR_PTYPE_NONF_IPV6_GTPU_EH_IPV6_OTHER,
171     NM_FLTR_PTYPE_NONF_IPV4_L2TPV3,
172     NM_FLTR_PTYPE_NONF_IPV6_L2TPV3,
173     NM_FLTR_PTYPE_NONF_IPV4_ESP,
174     NM_FLTR_PTYPE_NONF_IPV6_ESP,
175     NM_FLTR_PTYPE_NONF_IPV4_AH,
176     NM_FLTR_PTYPE_NONF_IPV6_AH,
177     NM_FLTR_PTYPE_NONF_IPV4_NAT_T_ESP,
178     NM_FLTR_PTYPE_NONF_IPV6_NAT_T_ESP,
179     NM_FLTR_PTYPE_NONF_IPV4_PFCP_NODE,
180     NM_FLTR_PTYPE_NONF_IPV4_PFCP_SESSION,
181     NM_FLTR_PTYPE_NONF_IPV6_PFCP_NODE,
182     NM_FLTR_PTYPE_NONF_IPV6_PFCP_SESSION,
183     NM_FLTR_PTYPE_NON_IP_L2,
184     NM_FLTR_PTYPE_FRAG_IPV4,
185     NM_FLTR_PTYPE_NONF_IPV6_UDP,
186     NM_FLTR_PTYPE_NONF_IPV6_TCP,
187     NM_FLTR_PTYPE_NONF_IPV6_SCTP,
188     NM_FLTR_PTYPE_NONF_IPV6_OTHER,
189     NM_FLTR_PTYPE_MAX,
190 };
191
192 struct rte_flow {
193     uint8_t filter_type;
194     void *rule;
195     uint32_t counter_id;
196 };
197

```

```

198 struct nm_fdir_l2 {
199     uint8_t dst_mac[NM_ETH_MAC_LEN]; /* dest MAC address */
200     uint8_t src_mac[NM_ETH_MAC_LEN]; /* src MAC address */
201     uint16_t ether_type; /* for NON_IP_L2 */
202 };
203
204 struct nm_fdir_session {
205     uint32_t session_id;
206 };
207
208 struct nm_fdir_l4 {
209     uint16_t dst_port;
210     uint16_t src_port;
211 };
212
213 struct nm_fdir_l3{
214     uint32_t dst_ip[NM_IPV6_ADDR_LEN_AS_U32];
215     uint32_t src_ip[NM_IPV6_ADDR_LEN_AS_U32];
216     uint8_t ip_ver;
217     uint8_t tos;
218     uint8_t ttl;
219     uint32_t l4_header; /* next header */
220     uint32_t sec_parm_idx; /* security parameter index */
221     uint8_t tc;
222     uint8_t proto;
223     uint8_t hlim;
224     uint16_t hdr_checksum;
225 };
226
227 /* upper outer tunnel tabl */
228
229 struct nm_l4_outer_tunnel_conf {
230     union {
231         uint32_t tunnel_id;
232         uint8_t vni[NM_FLOW_TABLE_LENGTH];
233     };
234
235     uint16_t dst_port;
236     uint8_t dip[NM_IPV6_ADDR_LEN_AS_U8]; /**< IP address of destination
237 host(s). */
238 };
239
240 /* upper tn1 inner tab */
241 struct nm_l2_inner_tunnel_conf {
242     uint32_t port; /* packet forward port */
243     uint16_t metadata;
244     struct nm_fdir_l2 l2;
245 };
246
247 /* downer non ip */
248 struct nm_l2_tunnel_conf {
249     uint32_t port;
250     struct nm_fdir_l2 l2;
251 };
252
253 /* down ip */
254 struct nm_l3_tunnel_conf {
255     uint32_t port;

```

```

255     struct nm_fdir_l2 l2;
256     uint8_t  proto;
257     uint8_t  sip[16];
258     uint8_t  dip[16];
259 };
260
261 /* security grooup */
262 struct nm_l4_sec_group_conf {
263     uint16_t metadata;
264     uint16_t src_port;
265     uint16_t dst_port;
266 };
267
268 struct nm_l4_outer_tunnel_conf_ele {
269     TAILQ_ENTRY(nm_l4_outer_tunnel_conf_ele) entries;
270     struct nm_l4_outer_tunnel_conf filter_info;
271     void *child;
272     void*  act;
273 };
274
275 struct nm_l2_inner_tunnel_conf_ele {
276     TAILQ_ENTRY(nm_l2_inner_tunnel_conf_ele) entries;
277     struct nm_l2_inner_tunnel_conf filter_info;
278     void *child;
279     void *act;
280 };
281
282 struct nm_ethertype_filter_ele {
283     TAILQ_ENTRY(nm_ethertype_filter_ele) entries;
284     struct nm_l2_tunnel_conf filter_info;
285     void *child;
286     void *act;
287 };
288
289 struct nm_l3_tunnel_conf_ele {
290     TAILQ_ENTRY(nm_l3_tunnel_conf_ele) entries;
291     struct nm_l3_tunnel_conf filter_info;
292     void *child;
293     void *act;
294 };
295
296 struct nm_l4_secgrp_conf_ele {
297     TAILQ_ENTRY(nm_l4_secgrp_conf_ele) entries;
298     struct nm_l4_sec_group_conf filter_info;
299     void *parent;
300     void *act;
301 };
302
303
304 /* rte_flow store */
305 struct nm_flow_mem {
306     TAILQ_ENTRY(nm_flow_mem) entries;
307     struct rte_flow *flow;
308 };
309
310 TAILQ_HEAD(nm_l4_outer_tun_filter_list, nm_l4_outer_tunnel_conf_ele);
311 TAILQ_HEAD(nm_l2_inner_tun_filter_list, nm_l2_inner_tunnel_conf_ele);
312 TAILQ_HEAD(nm_ethertype_filter_list, nm_ethertype_filter_ele);

```

```

313 TAILQ_HEAD(nm_l3_tunnel_filter_list, nm_l3_tunnel_conf_ele);
314 TAILQ_HEAD(nm_l4_secgrp_filter_list, nm_l4_secgrp_conf_ele);
315 TAILQ_HEAD(nm_flow_mem_list, nm_flow_mem);
316
317 /* hash-list struct */
318 struct nm_l4_outer_tnl_filter {
319     TAILQ_ENTRY(nm_l4_outer_tnl_filter)    entries;
320     struct nm_l4_outer_tunnel_conf        *key;
321 };
322
323 TAILQ_HEAD(nm_l4_outer_tnl_hash_list, nm_l4_outer_tnl_filter);
324
325 struct nm_l4_outer_tnl_hash_info {
326     struct nm_l4_outer_tnl_hash_list  l4_outer_tnl_hash_list;
327     struct nm_l4_outer_tnl_filter    **hash_map;
328     struct rte_hash                  *hash_handle;
329 };
330
331 struct nm_l2_ethtype_tnl_filter {
332     TAILQ_ENTRY(nm_l2_ethtype_tnl_filter)    entries;
333     struct nm_l2_tunnel_conf                *key;
334 };
335
336 TAILQ_HEAD(nm_l2_ethtype_tnl_hash_list, nm_l2_ethtype_tnl_filter);
337
338 struct nm_l2_ethtype_tnl_hash_info {
339     struct nm_l2_ethtype_tnl_hash_list  l2_ethtype_tnl_hash_list;
340     struct nm_l2_ethtype_tnl_filter    **hash_map;
341     struct rte_hash                  *hash_handle;
342 };
343
344 struct nm_l3_tnl_filter {
345     TAILQ_ENTRY( nm_l3_tnl_filter)    entries;
346     struct nm_l3_tunnel_conf        *key;
347 };
348
349 TAILQ_HEAD(nm_l3_hash_list, nm_l3_tnl_filter);
350
351 struct nm_l3_hash_info {
352     struct nm_l3_hash_list          l3_hash_list;
353     struct nm_l3_tnl_filter        **hash_map;
354     struct rte_hash                *hash_handle;
355 };
356
357 struct nm_flow_item_vxlan {
358     uint8_t flags;
359     uint8_t vni[NM_FLOW_TABLE_VXLAN_LENGTH];
360 };
361
362 struct nm_flow_item_vlan{
363     uint16_t vlan_tci; /*< Priority (3) + CFI (1) + Identifier Code (12)
364     */
365     uint16_t eth_proto;
366 };
367
368 struct nm_rule_action {
369     uint64_t flag; /* action flag, eg:set ipv4 src/redirect */

```

```

370     uint32_t drop_flag; /* drop or forward */
371     uint32_t soft_id; /* an unique value for this rule todo*/
372     uint32_t mark_id;
373     uint32_t counter_id;
374     uint8_t outer_type; //outer_tnl_type,eg:tunnel
375     uint8_t queue; /* assigned rx queue */
376     uint8_t queues_cnt;
377     uint8_t action_cnt; /* different action type total cnt */
378
379     uint16_t metadata; /* secondary flow table need */
380
381     /* set ops */
382     struct nm_fdir_l4 l4_outer;
383     struct nm_fdir_l2 l2_data_outer;
384     struct nm_fdir_l3 ip_outer;
385
386     /* push ops */
387     struct nm_flow_item_vlan vlan;
388     struct nm_flow_item_vxlan vxlan;
389     uint8_t set_flag;
390
391     uint8_t buf[NM_FLOW_ACTION_SET_OUTER_HEADER_MAX_LEN];
392     uint32_t l2_cnt:2;
393     uint32_t l3_cnt:2;
394     uint32_t l4_cnt:2;
395     uint32_t tnl_cnt:2;
396     uint32_t vlan_cnt:2;
397     uint32_t rsv:22;
398 };
399
400 struct nm_fdir_fltr {
401     struct nm_fdir_l3 ip;
402     struct nm_fdir_l3 ip_mask;
403
404     struct nm_fdir_l3 ip_outer;
405     struct nm_fdir_l3 ip_mask_outer;
406
407     struct nm_fdir_l4 l4;
408     struct nm_fdir_l4 l4_mask;
409
410     struct nm_fdir_l4 l4_outer;
411     struct nm_fdir_l4 l4_mask_outer;
412
413     struct nm_fdir_l2 l2_data_outer;
414     struct nm_fdir_l2 l2_mask_outer;
415
416     struct nm_fdir_l2 l2_data;
417     struct nm_fdir_l2 l2_mask;
418
419     uint16_t vlan_type; /* VLAN ethertype */
420     uint16_t vlan_tag; /* VLAN tag info */
421
422     struct nm_fdir_session l2tpv3_data;
423     struct nm_fdir_session l2tpv3_mask;
424
425     union{
426         uint32_t tunnel_id;
427         struct nm_flow_item_vxlan vxlan;

```

```

428     }tnl,tnl_mask;
429
430     uint32_t l2_cnt:1;
431     uint32_t l3_cnt:1;
432     uint32_t l4_cnt:1;
433     uint32_t tnl_cnt:1;
434     uint32_t rsv:28;
435
436     uint32_t port;
437 };
438
439 struct nm_acl_conf {
440     struct nm_fdir_fltr input;
441     uint64_t input_set;    //destroy rule by input-set
442 };
443
444 struct nm_action_filter {
445     TAILQ_ENTRY( nm_action_filter)    entries;
446     struct nm_rule_action              key;
447     uint32_t ref_cnt;    //ref times by different key type
448 };
449
450 TAILQ_HEAD(nm_action_hash_list, nm_action_filter);
451
452 struct nm_action_hash_info {
453     struct nm_action_hash_list    action_hash_list;
454     struct nm_action_filter        **hash_map;
455     struct rte_hash                *hash_handle;
456 };
457
458 struct nm_filter_info {
459     struct nm_flow_mem_list nm_flow_list;
460     struct nm_l4_outer_tun_filter_list filter_l4_outer_list;
461     struct nm_l2_inner_tun_filter_list filter_l2_inner_list;
462     struct nm_ethertype_filter_list filter_ethertype_list;
463     struct nm_l3_tunnel_filter_list filter_l3_tunnel_list;
464     struct nm_l4_secgrp_filter_list filter_l4_tunnel_list;
465
466     struct nm_l4_outer_tnl_hash_info l4_outer_hash;
467     struct nm_l2_ethertype_tnl_hash_info l2_eth_hash;
468     struct nm_l3_hash_info l3_hash;
469     struct nm_action_hash_info act;
470 };
471
472 #define NM_DEV_PRIVATE_TO_FILTER_INFO(adapter) \
473     (&(struct nm_adapter *)adapter->filter_info)
474
475 struct nm_flow_counter {
476     uint32_t shared:1;    /* Share counter ID with other flow rules. */
477     uint32_t ref_cnt:31; /* Reference counter. */
478     uint16_t id;    /* Counter ID. */
479     uint64_t hits; /* Number of packets matched by the rule. */
480 };
481
482 void nm_flow_init(struct rte_eth_dev *dev);
483 void nm_flow_fini(struct rte_eth_dev *dev);
484
485 #endif

```

