

Chasing Page Pool Leaks

Dragoș Tătulea, NetDev 0x19 March 2025, Zagreb - Croatia

... with drgns



Agenda

- Scope and Motivation
- Introduction
- Toolbox
- Examples / Tutorial







Scope What is this talk about?

- I have a page pool leak, is this a bug or a false positive?
- Where is it coming from?



Motivation

- Encountered quite a few of these leaks.
- Discussed on netdev mailing list as well.
- Wanted to learn drgn.



> but I think those page pointed to by anything...

From: Jakub Kicinski <kuba@kernel.org>

To: Yonglong Liu <liuyonglong@huawei.com> Cc: Yunsheng Lin <linyunsheng@huawei.com>, <netdev@vger.kernel.org>, <davem@davemloft.net>, <edumazet@google.com>, <pabeni@redhat.com>, <ilias.apalodimas@linaro.org>, Jesper Dangaard Brouer <hawk@kernel.org>, Alexander Duyck <alexander.duyck@gmail.com> Subject: Re: [RFC net] net: make page pool stall netdev unregistration to avoid IOMMU crashes Date: Wed, 14 Aug 2024 07:56:03 -0700 [thread overview] Message-ID: <20240814075603.05f8b0f5@kernel.org> (raw) In-Reply-To: <cec044dc-504c-47e6-8ffa-58e8c9b42713@huawei.com>

On Wed, 14 Aug 2024 18:09:59 +0800 Yonglong Liu wrote: > On 2024/8/10 11:57, Jakub Kicinski wrote: > > On Fri, 9 Aug 2024 14:06:02 +0800 Yonglong Liu wrote: >>> [7724.272853] hns3 0000:7d:01.0: page_pool_release_retry(): eno1v0 >> stalled pool shutdown: id 553, 82 inflight 6706 sec (hold netdev: 1855491) > > Alright : (You gotta look around for those 82 pages somehow with drgn. > > bpftrace+kfunc the work that does the periodic print to get the address > > of the page pool struct and then look around for pages from that pp.. :(

> I spent some time to learn how to use the drgn, and found those page,

> is allocated by the hns3 driver, how to find out who own those page now?

Scan the entire system memory looking for the pointer to this page. Dump the memory around location which hold that pointer. If you're lucky the page will be held by an skb, and the memory around it will look like struct skb_shared_info. If you're less lucky the page is used by sk_buff for the head and address will not be exact. If you're less lucky still the page will be directly leaked by the driver, and not

I think the last case is most likely, FWIW.



Introduction

- What is a page pool?
- What is page inflight state?
- When is a leak?
- Is could be a:
 - Bug: miscounting of resources.
 - False positive: page still referenced for a valid reason.
- _refcount vs pp_ref_count





unsigned long flags;

/* Atomic flags, some possibly * updated asynchronously */

* Five words (20/40 bytes) are available in this union. * WARNING: bit 0 of the first word is used for PageTail(). That * means the other users of this union MUST NOT use the bit to * avoid collision and false-positive PageTail().

/* Usage count. *D0 NOT USE DIRECTLY*. See page_ref.h */ atomic_t _refcount;





Inflight Page Lifetime

- Full page.
- Page fragment:
 - Size known in advance





Inflight Page Lifetime with bias



- Size not known in advance
- 2 step release

• How about queue teardown?





Types of Leaks

- Driver side
 - Incorrect refcounting at runtime
 - Incorrect refcounting at teardown
- System side
 - False positive
 - Bug



runtime teardown



Toolbox





Let's begin





- A leak was found ...
- Where does it come from?



+ /w/linux/tools/net/ynl/cll.py [{'detach-time': 149, 'id': 3, 'ifindex': 3, 'inflight': 1, 'inflight-mem': 4096}]

+ /w/linux/tools/net/ynl/cli.py --no-schema --spec /w/linux/Documentation/netlink/specs/netdev.yaml --dump page-pool-get [{'detach-time': 149,

ge-pool-get

Variation 1

- Search sockets for SKBs with leaked page_pool pages.
- Use tcp_sock.py.



tcp_sock.py

tcp_hashinfo = prog.object("tcp_hashinfo") for i in range(tcp_hashinfo.ehash_mask + 1): head = tcp_hashinfo.ehash[i].chain if hlist_nulls_empty(head): continue for sk in sk_nulls_for_each(head): # Filter by interface: continue first_skb = sk.sk_receive_queue.next skb = first_skb while skb != None: try: # Check linear part of skb: page = virt_to_page(skb.data) # Check fragments: shinfo = skb_shinfo(skb) for i in range(0, shinfo.nr_frags): frag = shinfo.frags[i] except FaultError: continue # Move to next skb: skb = skb_next if skb == first_skb: break

```
if ifindex > 0 and sk.sk_rx_dst_ifindex.value_() != ifindex:
             if page.pp_magic == PP_SIGNATURE and page.pp.user.detach_time:
    print(f"Found leaked page {hex(page)} in linear part of skb: {hex(skb.address_of_())}")
                  page = Object(prog, "struct page", address=frag.netmem)
if page.pp_magic == PP_SIGNATURE and page.pp.user.detach_time:
    print(f"Found leaked page {hex(page.address_of_())} in skb frag {i} of skb: {hex(skb.address_of_())}")
```

Found leaked page 0xffffea000016ab40 in linear part of skb: 0xffff88800d7689e8



Variation 1

- Check "struct page"
- Peek in page.





```
>>> Object(prog, "struct page", address=0xffffea000016ab40)
(struct page){
        _flags = (unsigned long)36028797018963968,
        .lru = (struct list_head){
                 .next = (struct list_head *)0xdead0000000000040,
                 .prev = (struct list_head *)0xffff88800b1ac800,
         ζ,
        .__filler = (void *)0xdead000000000000040,
        .mlock_count = (unsigned int)186304512,
        .buddy_list = (struct list_head){
                 .next = (struct list_head *)0xdead0000000000040,
                 .prev = (struct list_head *)0xffff88800b1ac800,
        ζ,
        .pcp_list = (struct list_head){
                 .next = (struct list_head *)0xdead0000000000040,
                 .prev = (struct list_head *)0xffff88800b1ac800,
         ζ,
        .mapping = (struct address_space *)0x0,
        .index = (unsigned long)95080448,
        .share = (unsigned long)95080448,
        .private = (unsigned long)1,
        .pp_magic = (unsigned long)16045481047390945344,
.pp = (struct page_pool *)0xffff88800b1ac800,
        ._pp_mapping_pad = (unsigned long)0,
        .dma_addr = (unsigned long)95080448,
        .pp_ref_count = (atomic_long_t){
                 .counter = (s64)1, 🭊
         ζ,
        .compound_head = (unsigned long)16045481047390945344,
        .pgmap = (struct dev_pagemap *)0xdead0000000000040,
        .zone_device_data = (void *)0xffff88800b1ac800,
        .callback_head = (struct callback_head){
                 .next = (struct callback_head *)0xdead0000000000040,
                 .func = (void (*)(struct callback_head *))0xffff88800b1ac800,
        .page_type = (unsigned int)4294967295,
        ._mapcount = (atomic_t){
                 .counter = (int)-1,
       },
        ._refcount = (atomic_t){
                 .counter = (int)1,
        .memcg_data = (unsigned long)0,
```

<pre>> print_annotate</pre>	ed_memory(page_to_	_virt(page.address_of_()),	128)
DRESS	VALUE		
ff888005aad000:	0000000000000000000		
ff888005aad008:	0000000000000000000		
ff888005aad010:	00000000000000000000		
ff888005aad018:	00000000000000000000		
ff888005aad020:	0000000000000000000		
ff888005aad028:	0000000000000000000		
ff888005aad030:	0000000000000000000		
ff888005aad038:	0000000000000000000		
ff888005aad040:	79fe573412005452		
ff888005aad048:	0045000853158c69		
ff888005aad050:	06400040e6bb5200		
ff888005aad058:	01010a010101b37a		
ff888005aad060:	1840393038d40101		
ff888005aad068:	1880944e2a81aef4		
ff888005aad070:	010100002b78f601		
ff <u>888005aad078</u> :	7e05558744e50a08		

Variation 1

• Let's look at the skb.



```
>>> Object(prog, "struct sk_buff *", address=0xffff88800d7689e8)
*(struct sk_buff *)0xffff888004fdb000 = {
        .next = (struct sk_buff *)0xffff88800d7689e8,
        .prev = (struct sk_buff *)0xffff88800d7689e8,
        .dev = (struct net_device *)0x0,
        .dev_scratch = (unsigned long)0,
        .rbnode = (struct rb_node){
                .__rb_parent_color = (unsigned long)18446612682295904744,
                .rb_right = (struct rb_node *)0xffff888800d7689e8,
                .rb_left = (struct rb_node *)0x0,
       },
.list = (struct list_head){
                .next = (struct list_head *)0xffff88800d7689e8,
                .prev = (struct list_head *)0xffff88800d7689e8,
        },
        .ll_node = (struct llist_node){
                .next = (struct llist_node *)0xffff88800d7689e8,
        },
        .sk = (struct sock *)0xffff88800d768940,
        .tstamp = (ktime_t)0,
        .skb_mstamp_ns = (u64)0,
        .cb = (char [48])'' xae xf4 x18@ xcc xf4 x18@'',
        ._skb_refdst = (unsigned long)0,
        .destructor = (void (*)(struct sk_buff *))sock_rfree+0x0 = 0xffffffff81902740,
        .tcp_tsorted_anchor = (struct list_head){
                .next = (struct list_head *)0x0,
                .prev = (struct list_head *)sock_rfree+0x0 = 0xffffffff81902740,
        },
        ._sk_redir = (unsigned long)0,
        ._nfct = (unsigned long)0,
        .len = (unsigned int)30,
        .data_len = (unsigned int)0,
        mac_len = (\_u16)14,
        .hdr_len = (__u16)0,
        .queue_mapping = (___u16)1,
        \_\_cloned_offset = (\__u8 [0]){},
        .cloned = (\_u8)0,
        .nohdr = (__u8)0,
        fclone = (\_u8)0,
        peeked = (\_u8)0,
        head_frag = (\_u8)1,
        .pfmemalloc = (\_u8)0,
        pp_recycle = (\_u8)1,
        .active_extensions = (__u8)0,
        .___pkt_type_offset = (___u8 [0]){},
        .pkt_type = (\_u8)0,
```

```
.tail = (sk_buff_data_t)160,
.end = (sk_buff_data_t)192,
.head = (unsigned char *)0xffff888005aad000 = "", 
.truesize = (unsigned int)768,
.users = (refcount_t){
      .refs = (atomic_t){
            .counter = (int)1,
      },
},
.extensions = (struct skb_ext *)0x0,
```





Example 1 What was found

- Leaked page in SKB linear area.
- Direction: Socket -> SKB -> Page
- ... let's try the other way around: Page -> SKB -> Socket



• • •

head

data



Variation 2

• Scan all leaked page pool pages.



ls_pp_leaks.py

try:

```
for page in for_each_page():
          if page.pp_magic == PP_SIGNATURE and page.pp.user.detach_time:
             print(page)
             print_annotated_memory(page_to_virt(page), 128)
       except FaultError:
          continue
```

Leaked page: 0xff	ffea000016ab40
Page content:	
ADDRESS	VALUE
ffff888005aad000:	000000000000000000
ffff888005aad008:	000000000000000000
ffff888005aad010:	000000000000000000
ffff888005aad018:	0000000000000000000
ffff888005aad020:	0000000000000000000
ffff888005aad028:	000000000000000000
ffff888005aad030:	000000000000000000
ffff888005aad038:	0000000000000000000
ffff888005aad040:	79fe573412005452
ffff888005aad048:	0045000853158c69
ffff888005aad050:	06400040e6bb5200
ffff888005aad058:	01010a010101b37a
ffff888005aad060:	1840393038d40101



Variation 2

• This looks familiar...



```
>>> Object(prog, "struct page", address=0xffffea000016ab40)
(struct page){
        .flags = (unsigned long)36028797018963968,
        .lru = (struct list_head){
                .next = (struct list_head *)0xdead0000000000040,
                .prev = (struct list_head *)0xffff88800b1ac800,
        .___filler = (void *)0xdead0000000000040,
        .mlock_count = (unsigned int)186304512,
        .buddy_list = (struct list_head){
                .next = (struct list_head *)0xdead0000000000040,
                .prev = (struct list_head *)0xffff88800b1ac800,
        ነ,
        .pcp_list = (struct list_head){
                .next = (struct list_head *)0xdead0000000000040,
                .prev = (struct list_head *)0xffff88800b1ac800,
        .mapping = (struct address_space *)0x0,
        .index = (unsigned long)95080448,
        .share = (unsigned long)95080448,
        .private = (unsigned long)1,
        .pp_magic = (unsigned long)16045481047390945344,
        .pp = (struct page_pool *)0xffff88800b1ac800,
       ._pp_mapping_pad = (unsigned long)0,
        .dma_addr = (unsigned long)95080448,
        .pp_ref_count = (atomic_long_t){
                .counter = (s64)1, 
        ζ,
        .compound_head = (unsigned long)16045481047390945344,
        .pgmap = (struct dev_pagemap *)0xdead0000000000040,
        .zone_device_data = (void *)0xffff88800b1ac800,
        .callback_head = (struct callback_head){
                .next = (struct callback_head *)0xdead0000000000040,
                .func = (void (*)(struct callback_head *))0xffff88800b1ac800,
        },
        .page_type = (unsigned int)4294967295,
        ._mapcount = (atomic_t){
                .counter = (int)-1,
        <u>۲</u>
        ._refcount = (atomic_t){
                .counter = (int)1,
        了,
        .memcg_data = (unsigned long)0,
```



Variation 2

- Find SKBs referencing address in page range.
- Search for most significant 6 bytes.
- Check next 4 bits.
- Once found, check that address is within page range.



snippet of find.py

 $ptr_size = 8$ else: results = [] return results

```
def search_page_reference(page):
    val = page_to_virt(page).value_()
    skip_bytes = math.ceil(PAGE_SHIFT / 8)
    val_endian = val.to_bytes(ptr_size, byteorder)
    if byteorder == "little":
        big_needle = val_endian[skip_bytes:ptr_size - skip_bytes]
        big_needle = val_endian[0:ptr_size - skip_bytes]
    small_needle = val >> PAGE_SHIFT
    # Search for first 6 bytes:
    for addr in search_memory(prog, big_needle):
        if byteorder == "little":
            # Adjust address to skipped bytes:
            addr = addr - skip_bytes
        mem_bytes = prog read(addr, ptr_size)
        mem_val = int.from_bytes(mem_bytes, <u>byteorder</u>)
        if mem_val >> PAGE_SHIFT == small_needle:
            results_append((addr, mem_val))
```



Variation 2

• Found something. Probably much more.

• Interpret as SKB.

••• DVIDIA. 19

Fou	nd	re	fere	ence	at	0 x	fff	f88	80	04 [.]
ADD	RES	S				VAL	UE			
fff	f88	80	04fc	lb00	0:	fff	f88	800	d7	689
fff	f88	80	04fc	lb00	8:	fff	f88	800	d7	689
fff	f88	80	04fc	lb01	0:	000	000	000	00	000
fff	f88	80	04fc	lb01	8:	fff	f88	800	d7	689
fff	f88	80	04fc	lb02	0:	000	000	000	00	000
fff	f88	80	04fc	lb02	8:	401	8f4	lcc4	01	8f4
fff	f88	80	04fc	lb03	0:	000	000)180	00	000
fff	f88	80	04fc	lb03	8:	000	000	008	12	a40
fff	f88	80	04fc	lb04	0:	000	000	000	00	000
fff	f88	80	04fc	lb04	8:	000	000	000	00	00
fff	f88	80	04fc	lb05	0:	000	000	000	00	000
fff	f88	80	04fc	lb05	8:	000	000	000	00	000
fff	f88	80	04fc	lb06	0:	fff	fff	ff8	19	027
fff	f88	80	04fc	lb06	8:	000	000	000	00	000
fff	f88	80	04fc	lb07	0:	000	000	000	00	000
fff	f88	80	04fc	lb07	8:	00a	000	010	00	00
fff	f88	80	04fc	lb08	0:	000	000	001	.042	200
fff	f88	80	04fc	lb08	8:	000	000	000	00	0a
fff	f88	80	04fc	lb09	0:	ec8	0bf	⁼ c30	00	000
fff	f88	80	04fc	lb09	8:	000	002	020	00	000
fff	f88	80	04fc	lb0a	0:	000	000	000	00	000
fff	f88	80	04fc	lb0a	8:	000	000	000	00	000
fff	f88	80	04fc	lb0b	0:	004	000)4e0	06	200
fff	f88	80	04fc	lb0b	8:	000	000)c00	00	000
fff	f88	80	04fc	lb0c	0:	fff	f88	800	5a	ad
fff	f88	80	04fc	lb0c	8:	fff	f88	800	5a	ad

```
fdb0c0. slab object: skbuff_head_cache+0xc0
9e8 [slab object: TCP+0xa8]
9e8 [slab object: TCP+0xa8]
 000
940 [slab object: TCP+0x0]
 900
 lae
 000
 e94
 003
 000
740 [function symbol: sock_rfree+0x0]
 01e
 00e
 040
 efb
 003
 000
 000
 000
 008
d000 [page offset in page 0xffffea000016ab40 page_pool id: 3, dev eth1, detached: yes]
d082 [page offset in page 0xffffea000016ab40 page_pool id: 3, dev eth1, detached: yes]
```



Example 1 Variation 2

- Interpret as SKB.
- This also looks familiar!



.user

```
ject(prog, "struct sk_buff", address=0xffff8888004fdb000)
uff){
= (struct sk_buff *)0xffff88800d7689e8,
= (struct sk_buff *)0xffff88800d7689e8,
= (struct net_device *)0x0,
scratch = (unsigned long)0,
ode = (struct rb_node){
   .___rb_parent_color = (unsigned long)18446612682295904744,
   .rb_right = (struct rb_node *)0xffff888800d7689e8,
   .rb_left = (struct rb_node *)0x0,
: = (struct list_head){
   .next = (struct list_head *)0xffff88800d7689e8,
   .prev = (struct list_head *)0xffff88800d7689e8,
node = (struct llist_node){
   .next = (struct llist_node *)0xffff888800d7689e8,
 (struct sock *)0xffff88800d768940,
amp = (ktime_t)0,
_mstamp_ns = (u64)0,
 (char [48])"\xae\xf4\x18@\xcc\xf4\x18@",
_refdst = (unsigned long)0,
tructor = (void (*)(struct sk_buff *))sock_rfree+0x0 = 0xfffffffff81902740,
_tsorted_anchor = (struct list_head){
   .next = (struct list_head *)0x0,
   .prev = (struct list_head *)sock_rfree+0x0 = 0xfffffff81902740,
_redir = (unsigned long)0,
t = (unsigned long)0,
= (unsigned int)30,
len = (unsigned int)0,
_len = (__u16)14,
_len = (__u16)0,
ue_mapping = (\_u16)1,
loned_offset = (\_u8 [0]){},
ned = (\_u8)0,
lr = (\__u8)0,
one = (\__u8)0,
ked = (\__u8)0,
_frag = (\__u8)1,
emalloc = (\_u8)0,
recycle = (\_u8)1,
lve_extensions = (\__u8)0,
t_type_offset = (\_u8 [0]){},
_type = (\__u8)0,
```

```
(sk_buff_data_t)192,
= (unsigned char *)0xffff888005aad000 = "",
size = (unsigned int)768,
s = (refcount_t)
  .refs = (atomic_t){
       .counter = (int)1,
  },
nsions = (struct skb_ext *)0x0,
```





Conclusions

- Found leaking page.
- Was in linear part of SKB.
- pp_ref_count = 1
- Diagnosis: false positive



- A new day, a new leak.
- Repeat steps...



+ /w/linux/tools/net/ynl/cli.py --no-schema --spec /w/linux/Documentation/netlink/specs/netdev.yaml --dump page-pool-get [{'detach-time': 431, 'id': 3, 'ifindex': 3, 'ifindex': 3, 'inflight': 2, 'inflight-mem': 8192}]



Variation 1

- Search sockets for SKBs with leaked page_pool pages.
- This time pages are frags...



Found leaked Found leaked

From tcp_sock.py

k	page	0xffffea0000495940	in	skb	frag	0	of	skb:	0xffff888
ł	page	0xffffea0000495900	in	skb	frag	1	of	skb:	0xffff888



Variation 1

- Let's peek into the SKB.
- Linear part not a page_pool page.





```
>>> skb = Object(prog, "struct sk_buff *", address=0xffff888005dc1328)
*(struct sk_buff *)0xffff888005ba8800 = {
        .next = (struct sk_buff *)0xffff888005dc1328,
        .prev = (struct sk_buff *)0xffff888005dc1328,
        .dev = (struct net_device *)0x0,
        dev_scratch = (unsigned long)0,
        .rbnode = (struct rb_node){
                .___rb_parent_color = (unsigned long)18446612682168341288,
               .rb_right = (struct rb_node *)0xffff888005dc1328,
               .rb_left = (struct rb_node *)0x0,
       },
       .list = (struct list_head){
                .next = (struct list_head *)0xffff8888005dc1328,
               .prev = (struct list_head *)0xffff888005dc1328,
        },
        .ll_node = (struct llist_node){
                .next = (struct llist_node *)0xffff888005dc1328,
        },
        .sk = (struct sock *)0xffff888005dc1280,
        .tstamp = (ktime_t)0,
        .skb_mstamp_ns = (u64)0,
        .cb = (char [48])''x8ax1dxx1c1xx81',
        ._skb_refdst = (unsigned long)0,
        .destructor = (void (*)(struct sk_buff *))sock_rfree+0x0 = 0xffffffff81902740,
        .tcp_tsorted_anchor = (struct list_head){
               .next = (struct list_head *)0x0,
               },
        ._sk_redir = (unsigned long)0,
        ._nfct = (unsigned long)0,
        len = (unsigned int)5010,
        .data_len = (unsigned int)4820,
        .mac_len = (__u16)14,
        hdr_len = (\_u16)0,
        .queue_mapping = (___u16)1,
        \ldots_cloned_offset = (__u8 [0]){},
        .cloned = (\_u8)0,
        .nohdr = (__u8)0,
        fclone = (__u8)0,
        peeked = (\_u8)0,
        head_frag = (\__u8)1,
        pfmemalloc = (\_u8)0,
        .pp_recycle = (__u8)1,
        .active_extensions = (__u8)0,
       .___pkt_type_offset = (___u8 [0]){},
.pkt_type = (___u8)0,
```

```
.tail = (sk_buff_data_t)320,
      .end = (sk_buff_data_t)704,
      .head = (unsigned char *)0xffff888005b99000 = "",
      .data = (unsigned char *)0xffff888005b99082 = "more data......data.....
.truesize = (unsigned int)6144,
      .users = (refcount_t){
            .refs = (atomic_t){
                  .counter = (int)1,
            },
      .extensions = (struct skb_ext *)0x0,
```





Variation 1

- Let's look at SKB shinfo.
- 2 frags, each a leaking page.





S,

```
*(struct skb_shared_info *)0xffff888005b992c0 = {
        .flags = (\_u8)0,
        .meta_len = (\_u8)0,
        nr_frags = (\_u8)2,
        tx_flags = (\_u8)0,
        .gso_size = (unsigned short)0,
        .gso_segs = (unsigned short)0,
        .frag_list = (struct sk_buff *)0x0,
        .hwtstamps = (struct skb_shared_hwtstamps){
                .hwtstamp = (ktime_t)0,
                .netdev_data = (void *)0x0,
        ζ,
        .xsk_meta = (struct xsk_tx_metadata_compl){
                tx_timestamp = (\_u64 *)0x0,
        ζ,
        .gso_type = (unsigned int)0,
        .tskey = (u32)0,
        .dataref = (atomic_t){
                .counter = (int)1,
        },
        .xdp_frags_size = (u32)0,
        .xdp_frags_truesize = (u32)0,
        destructor_arg = (void *)0x0,
        .frags = (skb_frag_t [17]){
                        .netmem = (netmem_ref)18446719884458547520,
                        .len = (unsigned int)1792,
                        .offset = (unsigned int)2304,
                ر کر
ر
                        netmem = (netmem_ref) 18446719884458547456,
                        len = (unsigned int)3028,
                        .offset = (unsigned int)0,
                ζ,
```



What was found

• Leaked pages in SKB fragments.

• Direction: Socket -> SKB -> shinfo -> Page

• ... let's try the other way around: Page -> shinfo -> SKB -> Socket



SKB	
head	
data	
end	
	_



Variation 2

- Scan all pages for leaked page pool pages.
- Found 2. Showing 1.



Output of ls_pp_leaks.py

Leak	ed	pa	age		02	xf	f	ff	ea	00	00	49	590
Page	СС	ont	en	it:									
ADDR	ESS	5						V	AL	UE			
ffff	888	301	L25	64	.0(00	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	.0(08	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	0	10	•	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	0	18	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	02	20	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	02	28	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	0	30	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	0	38	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	.04	40	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	.04	48	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	0!	50	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	0!	58	:	2	e2	e2	e2	e2	e2e
ffff	888	301	L25	64	.0	60	•	0	00	00	00	02	e2e
Leak	ed	pa	age	2	02	xf	f	ff	ea	00	00	49	594
Page	CC	ont	en	it:									
ADDR	ESS	5						V	٩L	UE			
ffff	888	301	L25	65	0	00	•	7	Øf	ef	ff	ff	fff
ffff	888	301	L25	65	0	08	•	0	10	00	60	89	b01
ffff	888	301	L25	65	0	10	•	7	0f	e0	10	00	406
ffff	888	301	L25	65	0	18	•	0	a0	20	10	19	b01
ffff	888	301	L25	65	02	20	•	0	10	10	00	00	000
ffff	888	301	L25	65	02	28	•	0	00	00	00	00	000
ffff	888	301	L25	65	0	30	•	0	00	00	00	00	000
ffff	888	301	L25	65	0	38	•	0	00	00	00	00	000
ffff	888	301	L25	65	04	40	•	0	00	00	00	00	001
ffff	888	301	L25	65	04	48	•	0	00	00	00	00	005
ffff	888	301	L25	65	0!	50	•	0	00	00	0 2	40	000
ffff	888	301	L25	65	0!	58	:	0	00	00	00	00	003
F F F F	000	101		C F	0	20	-	0	20	00	00	00	MAE

e2e2e ffff .af02 00008 .af02 00000 00102 00000 1bdfc 570d0 3393c 570df

```
>>> page = Object(prog, "struct page", address=0xffffea0000495940)
>>> page
(struct page){
       flags = (unsigned long)36028797018963968,
       .lru = (struct list_head){
               .prev = (struct list_head *)0xffff88800e236800,
       .___filler = (void *)0xdead00000000000040,
       .mlock_count = (unsigned int)237201408,
       .buddy_list = (struct list_head){
               .next = (struct list_head *)0xdead0000000000040,
               .prev = (struct list_head *)0xffff88800e236800,
       },
       .pcp_list = (struct list_head){
               .prev = (struct list_head *)0xffff888800e236800,
        ζ,
       .mapping = (struct address_space *)0x0,
       index = (unsigned long)307646464,
       .share = (unsigned long)307646464,
       .private = (unsigned long)1,
       .pp_magic = (unsigned long)16045481047390945344,
       .pp = (struct page_pool *)0xffff88800e236800,
       ._pp_mapping_pad = (unsigned long)0,
       dma_addr = (unsigned long)307646464,
       .pp_ref_count = (atomic_long_t){
               .counter = (s64)1,
        },
       .compound_head = (unsigned long)16045481047390945344,
       .pgmap = (struct dev_pagemap *)0xdead000000000040,
       .zone_device_data = (void *)0xffff88800e236800,
       .callback_head = (struct callback_head){
               .next = (struct callback_head *)0xdead0000000000040,
               .func = (void (*)(struct callback_head *))0xffff88800e236800,
        ζ,
       page_type = (unsigned int)4294967295,
       ._mapcount = (atomic_t){
               .counter = (int)-1,
        ζ,
       ._refcount = (atomic_t){
               .counter = (int)1,
        ζ,
```

```
.memcg_data = (unsigned long)0,
```



Variation 2

- Find page references.
 - Actual pointer.
- Does it look like SKB or shinfo?
- shinfo -> SKB



F	0	u	n	d		r	`e	t	e	r	`e	er)(C
A	D	D	R	E	S	S)							
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)2	2
f	f	f	f	8	8	8	0	0	5	b	9)Ç)]	3
f	f	f	f	8	8	8	0	0	5	b	9)Ç)]	3
f	f	f	f	8	8	8	0	0	5	b	9)Ç)]	3
f	f	f	f	8	8	8	0	0	5	b	9)Ç)]	3
(d	r	g	ne	eı	۱۱)		d	t	a	t	u	
E	~		n	Ч	1	~	~ f	-	ĥ	~	n	~	~	
Δ	ט ח	u N	RI RI	u F	י ק	ן נ ק	- 1	e	: 1	e	11	C	e	
f	f	f	f	с. 88	88	ן 19	30)5	b	а	8	8	0	(
f	f	f	f	88	88	30)e)5	b	ă	8	8	0	8
f	f	f	f	88	88	8(90)5	b	а	8	8	1	(
f	f	f	f	88	88	8(90)5	b	а	8	8	1	8
f	f	f	f	88	88	8(90)5	b	а	8	8	2	(
f	f	f	f	88	88	8(90)5	b	а	8	8	2	8
Ť	Ť	Ť	T	88	88	36	96)5	b	а	8	8	3	ł

F	0	u	n	d		r	e	f	e	r	e	n	С	e	
A	D	D	R	E	S	S									
f	f	f	f	8	8	8	0	0	5	b	а	8	8	0	l
f	f	f	f	8	8	8	0	0	5	b	а	8	8	0	8
f	f	f	f	8	8	8	0	0	5	b	а	8	8	1	6
f	f	f	f	8	8	8	0	0	5	b	а	8	8	1	8
f	f	f	f	8	8	8	0	0	5	b	а	8	8	2	6
f	f	f	f	8	8	8	0	0	5	b	а	8	8	2	8
f	f	f	f	8	8	8	0	0	5	b	a	8	8	3	6
f	f	f	f	8	8	8	0	0	5	b	a	8	8	3	8
f	f	f	f	8	8	8	0	0	5	b	a	8	8	4	0
f	f	f	f	8	8	8	0	0	5	b	a	8	8	4	8
f	f	f	f	8	8	8	0	0	5	ĥ	ă	8	8	5	0
f	f	f	f	8	8	8	õ	õ	5	ĥ	ă	8	8	5	2
f	f	f	י f	8	8	8	õ	õ	5	h	a	ุ้่ง	о 8	6	0
י f	י f	י f	י f	8	8	о 8	a	a	5	h	a	о 8	ง 8	6	2
י f	י f	י f	י f	о 8	2 8	о 8	a	a	5	h	a	2 8	о 8	7	0
י f	י f	י f	י f	0 8	0 8	0 8	a	a	5	h	a a	2 8	2 8	' 7	2
י f	י f	י f	י f	0 8	0 8	0 8	a	a	5	h	а а	0 2	0 8	י 2	0
י f	י f	י f	י f	0 0	0 0	0 0	0 A	0 A	5	h	a a	0 2	0 0	2	ں ج
י ר	י ר	י ר	י f	0 0	0 0	0 Q	0 0	0 0	5	h	a a	0 0	0 Q	a	0
י ר	י ר	י ר	י ר	0 Q	0 Q	0 Q	0 0	0 0	55	h	a a	0 Q	0 Q	0	۲ ۲
ו ר	ו ר	ו ר	ו ר	0	0	0	0 0	0 0	ך ב	b h	a a	0	0	ש כ	0
ן ר	ו ר	ן ר	ן ר	0	0	0 0	ย ก	ย ก	С Б	ม ห	a ~	0	0 0	a 2	e c
ן ר	ן ר	ן ר	ן ר	0	0	0 0	0 0	0 0	С С	ม ห	a	0	0 0	d L	C D
ן ר	ן ר	ן ר	ן ר	0	0	0	0	0	С Г	ม เก	a	0	0	0 6	۲ C
T ع	٦ ع	٦ ع	T ع	8	8	80	0	0	С Г	D	a	8	80	D	C C
T	T	T	T _	8	8	8	0	0	С Г	D	a	8	8	C	l

Output of find.py

e a	t_0xffff888005b992	2f0: va	alue	0x1	fffea0000	49594	40.		·	
	VALUE				and the second se					
90:	000000000000000000000000000000000000000									
98:	000000000000000000000000000000000000000									
a0:	000000000000000000000000000000000000000									
a8:	000000000000000000000000000000000000000									
:0 0	000000000000000000000000000000000000000									
: 8c	000000000000000000000000000000000000000									
c0:	0000000000020000									
:8	00000000000000000000									
:0b	000000000000000000000000000000000000000									
: 8t	000000000000000000000000000000000000000									
e0:	000000000000000000000000000000000000000									
e8:	000000000000000000000000000000000000000									
f0:	ffffea0000495940	[page	ref	in	page_pool	_ id:	3,	dev	eth1,	detached
f8:	0000090000000700									
00:	ffffea0000495900	[page	ref	in	page_pool	_ id:	3,	dev	eth1,	detached
08:	00000000000000000000000000000000000000									
10:	000000000000000000000000000000000000000									
18:	00000000000000000000									

lea@virtme-ng /x> ./find.py ffff888005b992f0 --peek-skb

at 0xffff888005ba88c0: value 0xffff888005b99000. slab object: skbuff_head_cache+0xc0 VALUE 0: ffff888005dc1328 [slab object: TCP+0xa8] 8: ffff888005dc1328 [slab object: TCP+0xa8] : ffff888005dc1280 [slab object: TCP+0x0] : 00000000000000000 **:** 815c311c815c1d8a : 000000180000000 00000000be00254f 00000000000000003 000000000000000000 : 0000000000000000 0000000000000000000 : 00000000000000000 000012d400001392 00a000010000000e 0: 0000000010420040 8: 000000000003ae6 : 8fce04210000003 0000020200000000 000000000000000000 000000000000000000 0040004e00620008 000002c000000140 0: ffff888005b99000 ffff888005ba88c8: ffff888005b99082





Variation 2

- Let's check...
- Bingo!



>	>	>		S	k	b		=		0	b	•
>	>	>		S	k	b						
(S	t	r	u	С	t		S	k	_	b	ι
								•	n	e	X	1
								•	p	r	e	١
								•	d	e	V	
								•	d	e b	V	-
								•	Г	D	n	(
								l				
								۲ •	ί	i	s	1
								} •	, l	l		ŗ
								}	,			
								•	S	k		=
								•	t	S	t	ē
								•	S	k h	b	-
								•	C	0 c	k	= }
								•	d	э е	r S	1
								•	t	C	р	
								١				
								ን -	,	S	k	
								•		n	f	-
								•	l	e	n	
								•	d	а	t	ē
								•	m	a	C	
								•	n a	d	r	- -
								•	q	u	e c	ן ן
								•	– c	1	с О	ľ
								•	n	0	h	1
								•	f	С	ι	(
								•	р	e	e	ł
								•	h	e f	a	(
								•	p n	T n	M	e
									Р а	Р Р	Ť	
								•			p	ł
								•	р	k	t	=
									t	a	i]
									e	n	d	
								•	h	e	a	С
								•	d	a	t	ĉ
•	•	•	•		•	•	•	•	•	•	•	•
								•	t	r	U	e
								•	U	S	e	r
								}				
								J	ı e	X	t	e
}												

ject(prog, "struct sk_buff", address=0xffff888005ba8800)

```
uff){
: = (struct sk_buff *)0xffff888005dc1328,
v = (struct sk_buff *)0xffff888005dc1328,
= (struct net_device *)0x0,
_scratch = (unsigned long)0,
ode = (struct rb_node){
   .__rb_parent_color = (unsigned long)18446612682168341288,
   .rb_right = (struct rb_node *)0xffff8888005dc1328,
   .rb_left = (struct rb_node *)0x0,
 = (struct list_head){
   .next = (struct list_head *)0xffff888005dc1328,
   .prev = (struct list_head *)0xffff888005dc1328,
node = (struct llist_node){
   .next = (struct llist_node *)0xffff8888005dc1328,
 (struct sock *)0xffff888005dc1280,
amp = (ktime_t)0,
_mstamp_ns = (u64)0,
 (char [48])'' \times 8a \times 1d \times x81 \times 1c1 \times x81'',
_refdst = (unsigned long)0,
tructor = (void (*)(struct sk_buff *))sock_rfree+0x0 = 0xfffffffff81902740,
_tsorted_anchor = (struct list_head){
   .next = (struct list_head *)0x0,
   _redir = (unsigned long)0,
ct = (unsigned long)0,
= (unsigned int)5010,
a_len = (unsigned int)4820,
_{len} = (\__u16)14,
__len = (__u16)0,
ue_mapping = (__u16)1,
loned_offset = (__u8 [0]){},
ned = (__u8)0,
dr = (__u8)0,
one = (\__u8)0,
ked = (\__u8)0,
d_frag = (__u8)1,
emalloc = (__u8)0,
recycle = (__u8)1,
ive_extensions = (__u8)0,
kt_type_offset = (__u8 [0]){},
_type = (\__u8)0,
L = (sk_buff_data_t)320,
= (sk_buff_data_t)704,
esize = (unsigned int)6144,
rs = (refcount_t){
   .refs = (atomic_t){
           .counter = (int)1,
   },
ensions = (struct skb_ext *)0x0,
```



Conclusions

- Found page leaks in SKB fragments.
- pp_ref_count = 1
- Diagnosis: false positive.



• Last one.



+ /w/linux/tools/net/ynl/cli.py --no-schema --spec /w/linux/Documentation/netlink/specs/netdev.yaml --dump page-pool-get [{'detach-time': 40, 'id': 4, 'ifindex': 3, 'ifindex': 3, 'inflight': 1, 'inflight-mem': 4096, 'napi-id': 515}]



Variation 1

 Search sockets for SKBs with leaked page_pool pages.



Example 3 Variation 1

• No results!





Variation 2

 Scan all pages for leaked page pool pages.



(drgner Leaked Page co ADDRESS ffff888 ffff888

<pre>nv) dtatulea page: 0xff</pre>	@virtme_ng /x> ./ls_pp_leaks. ffea00004207c0
ontent:	
5	VALUE
301081f000:	8f11037a6bb52ca9
301081f008:	abcb9c9244b629ff
301081f010:	a59158bdadd9d0b9
301081f018:	37ab6a2e4a4cda09
301081f020:	a015038cf21050e7
301081f028:	13854b48b3f9be10
301081f030:	412c61dd7f680af2
301081f038:	8a3f965101b1de67
301081f040:	1beeb874d7cb7f42
301081f048:	beaa8d6c9fd31501
801081f050:	2e8bbf0b8452c9bc
801081f058:	c7baf6a920a9b5eb
301081f060:	00000000787bcaa <u>f</u>



Variation 2

- Let's peek in the page.
- Large pp_ref_count.
- HW GRO: header or data page?



>		р	а	g	e		=
> t	r	p u	a c	g t	e	р •	a f l
						} • •	, m b
						} •	, p
						} • •	, misppp dp
						} • •	, c p z c
						} •	, p
						} • •	, , m

```
Object(prog, "struct page", address=0xffffea00004207c0)
age){
flags = (unsigned long)36028797018963968,
Lru = (struct list_head){
      .next = (struct list_head *)0xdead0000000000040,
      .prev = (struct list_head *)0xffff888008586000,
_filler = (void *)0xdead000000000040,
nlock_count = (unsigned int)140009472,
buddy_list = (struct list_head){
      .next = (struct list_head *)0xdead0000000000040,
      .prev = (struct list_head *)0xffff8888008586000,
>cp_list = (struct list_head){
      .next = (struct list_head *)0xdead0000000000040,
      .prev = (struct list_head *)0xffff888008586000,
napping = (struct address_space *)0x2,
Index = (unsigned long)276951040,
hare = (unsigned long)276951040,
orivate = (unsigned long)64,
p_magic = (unsigned long)16045481047390945344,
op = (struct page_pool *)0xffff888008586000,
_pp_mapping_pad = (unsigned long)2,
lma_addr = (unsigned long)276951040,
>p_ref_count = (atomic_long_t){
      .counter = (s64)64,
compound_head = (unsigned long)16045481047390945344,
ogmap = (struct dev_pagemap *)0xdead000000000000,
cone_device_data = (void *)0xffff8888008586000,
callback_head = (struct callback_head){
      .next = (struct callback_head *)0xdead0000000000040,
      .func = (void (*)(struct callback_head *))0xffff888008586000,
bage_type = (unsigned int)4294967295,
_mapcount = (atomic_t){
      .counter = (int) - 1,
_refcount = (atomic_t){
      .counter = (int)1,
nemcg_data = (unsigned long)0,
```



Conclusions

Found leaking page with high pp_ref_count

- Driver leak.
- Leaked page was in header page.



Ending Notes

- How about XDP?
- Examples are on a <u>drgn fork</u> for now. Hoping to get some of it accepted in drgn.





Thank you! Questions?



