



上海交通大学  
SHANGHAI JIAO TONG UNIVERSITY

网络空间安全学院  
SJTU-CSE

# Detecting Kernel Memory Leaks in Specialized Modules With Ownership Reasoning

NIS-8018 论文阅读

谌翔 杨光

2022 年 4 月



# 目录



① Introduction

② K-MELD

③ Ownership Reasoning

④ Implementation

⑤ Evaluation

⑥ Conclusion

## 第 1 节

# Introduction

# Background



- Kernel Memory
  - Shared by hardware and all processes
  - Hard to reclaim
  - Denial of Service
- Challenges
  - Specialized functions.
  - Complicated and lengthy data flow.

# Contributions



- An approach for identifying specialized allocation functions.
- A rule-mining approach for corresponding specialized deallocations.
- An ownership reasoning mechanism for kernel objects.
- A scalable implementation K-MELD(Kernel MEmory Leak Detector)
  - 218 new bugs, 41 CVEs assigned.

# CVE-2019-19062



```
1 /* File: crypto/crypto_user_base.c */
2 static int crypto_report(struct sk_buff *in_skb,
3                         struct nlmsghdr *in_nlh, struct nlattr **attrs)
4 {
5     struct crypto_dump_info info;
6     ...
7     alg = crypto_alg_match(p, 0);
8     if (!alg)
9         return -ENOENT;
10
11    err = -ENOMEM;
12    /* Memory is allocated here */
13    skb = nlmsg_new(NLMSG_DEFAULT_SIZE, GFP_KERNEL);
14    if (!skb)
15        goto drop_alg;
16    ...
17    info.out_skb = skb;
18    ...
19    err = crypto_report_alg(alg, &info);
20
21 drop_alg:
22    crypto_mod_put(alg);
23
24    if (err)
25        return err; /* Memory leaks here */
26    ...
27 }
```

## 第 2 节

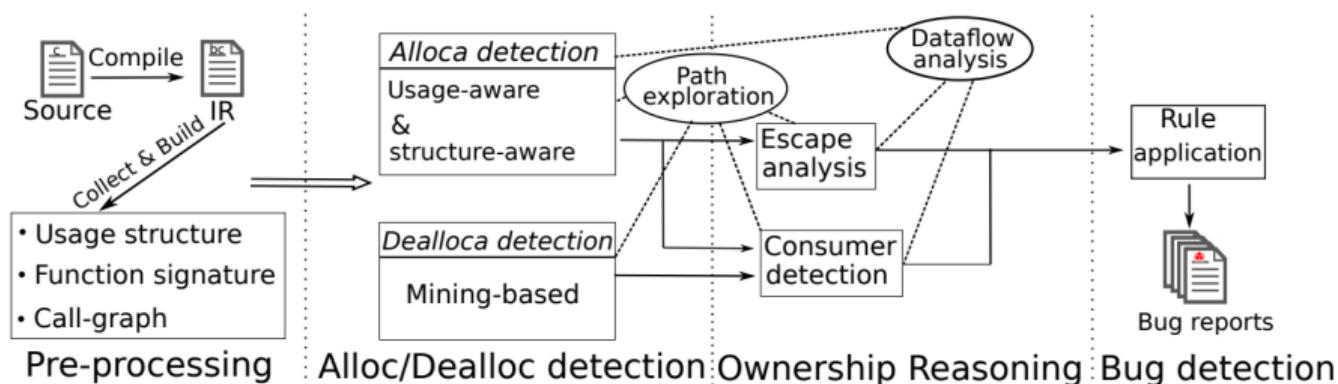
# K-MELD

# Kernel Memory Allocation and Leaks



- kmalloc & kfree
- vmalloc & vfree
- Specialized allocation
  - netlink: nlmsg\_new & nlmsg\_free

# K-MELD



# Allocation and Deallocation Identification



- Properties of allocation functions
  - Return a pointer
  - Followed by a NULL check
  - Not derived from another base pointer
  - Initialized before being used
- Methods
  - Use-finding
  - Source-finding

# Allocation and Deallocation Identification



- Detect deallocation functions
  - Error-Handling Paths
    - Explore the control flow graph(CFG)
  - Sequential Pattern Mining
    - ⟨call FOI, check, release, return⟩

## 第 3 节

# Ownership Reasoning

# Ownership



何谓 Ownership<sup>1</sup> → 对于动态分配内存的管理权

- 申请该动态内存的函数也负责了这块内存的回收
- 如果将动态分配的内存传递<sup>2</sup>给了 caller 或者 callee, ownerships 随之转移

---

<sup>1</sup><https://www.oreilly.com/library/view/programming-rust/9781491927274/ch04.html>

<sup>2</sup>通过引用或者全局变量

# Ownership



Ownership Reasoning 的作用 → 辅助过程间分析，追踪 ownership 的转移情况和内存的释放情况

推断方式	推断方向	作用
Enhanced escape analysis	CFG Upward	减少 caller 释放造成的误报
Consumer function detection	CFG Downward	减少 callee 造成的误报和漏报

## Enhanced escape analysis



**Escape:** ownership 被转移到了上层函数，当前函数不再需要负责内存的回收。

- **path-sensitive analysis:** 不需要再检测 escape 所在路径的内存泄漏情况
- **inter-procedure analysis:** 追踪 escape 以外的路径或遍历 escape 指针可能到达的上层函数，检测内存泄漏情况

# Case Study1 - Enhanced escape analysis



```
1 /* File: drivers/virt/vboxguest/vboxguest_utils.c */
2 static int hgcm_call_preprocess_linaddr(
3     const struct vmmdev_hgcm_function_parameter *src_parm,
4     void **bounce_buf_ret, size_t *extra)
5 {
6     void *buf, *bounce_buf;
7     bool copy_in;
8     u32 len;
9     int ret;
10
11    buf = (void *)src_parm->u.pointer.u.linear_addr;
12    len = src_parm->u.pointer.size;
13    copy_in = src_parm->type != VMMDEV_HGCM_PARM_TYPE_LINADDR_OUT;
14
15    if (len > VBG_MAX_HGCM_USER_PARM)
16        return -E2BIG;
17
18    /* Memory is allocated here */
19    bounce_buf = kmalloc(len, GFP_KERNEL);
20    /* Check for allocation success */
21    if (!bounce_buf)
22        return -ENOMEM;
23
24    if (copy_in) {
25        ret = copy_from_user(bounce_buf, (void __user *)buf, len);
26        if (ret)
27            return -EFAULT;
28    } else {
29        memset(bounce_buf, 0, len);
30    }
31
32    /* Allocation pointer is assigned to a reference argument */
33    *bounce_buf_ret = bounce_buf;
34    hgcm_call_add_pagelist_size(bounce_buf, len, extra);
35    return 0;
36 }
```

# Consumer function detection



**Consumer**: ownership 被转移到了下层函数，需要追踪 callee 的内存回收情况。如果 callee 的所有路径都对指针进行了逃逸/释放的操作，则 callee 不产生内存泄漏。

- **path-sensitive analysis**: 探索每条路径的内存泄漏情况
  - **Unconditional**: caller 的所有执行路径都会调用 consumer
  - **Conditional**: caller 的部分执行路径不会调用 consumer
- **inter-procedure analysis**: 需要继续追踪 consumer 中内存对象 escape 的情况

# Case Study2 - Unconditional consumer



```
1 /* File: drivers/net/ethernet/chelsio/cxgb4/srq.c */
2 int cxgb4_get_srq_entry(struct net_device *dev,
3                         int srq_idx, struct srq_entry *entryp)
4 {
5     ...
6     struct adapter *adap;
7     struct sk_buff *skb;
8     ...
9     adap = netdev2adap(dev);
10    ...
11    /* ALLOCATION */
12    skb = alloc_skb(sizeof(*req), GFP_KERNEL);
13    if (!skb)
14        return -ENOMEM;
15    ...
16    t4_mgmt_tx(adap, skb); /* CONSUMER */
17    ...
18    return rc;
19 }
20
21 /* File: drivers/net/ethernet/chelsio/cxgb4/sge.c */
```

```
21 /* File: drivers/net/ethernet/chelsio/cxgb4/sge.c */
22 int t4_mgmt_tx(struct adapter *adap, struct sk_buff *skb)
23 {
24     int ret;
25     ...
26     ret = ctrl_xmit(&adap->sge.ctrlq[0], skb); /* CONSUMER */
27     ...
28     return ret;
29 }
30
31 static int ctrl_xmit(struct sge_ctrl_txq *q, struct sk_buff *skb)
32 {
33     ...
34     if (unlikely(!is_immm(skb))) {
35         WARN_ON(1);
36         dev_kfree_skb(skb); /* RELEASE */
37         return NET_XMIT_DROP;
38     }
39     ...
40     if (unlikely(q->full)) {
41         ...
42         __skb_queue_tail(&q->sendq, skb); /* ESCAPE */
43         spin_unlock(&q->sendq.lock);
44         return NET_XMIT_CN;
45     }
46     ...
47     kfree_skb(skb); /* RELEASE */
48     return NET_XMIT_SUCCESS;
49 }
```

# Case Study3 - Conditional consumer



```
1 /* File: net/dsa/tag_ksz.c */
2 static struct sk_buff *ksz_common_xmit(struct sk_buff *skb,
3     struct net_device *dev, int len)
4 {
5     nskb = alloc_skb(NET_IP_ALIGN + skb->len +
6         padlen + len, GFP_ATOMIC);
7     if (!nskb)
8         return NULL;
9     ...
10    /* CONDITIONAL CONSUMER */
11    if (skb_put_padto(nskb, nskb->len + padlen))
12        return NULL;
13    ...
14    return nskb;
15 }
16 /* File: net/core/skbuff.c */
```

```
16 /* File: net/core/skbuff.c */
17 int __skb_pad(struct sk_buff *skb, int pad, bool free_on_error)
18 {
19     int err;
20     int ntail;
21     if (!skb_cloned(skb) && skb_tailroom(skb) >= pad) {
22         memset(skb->data+skb->len, 0, pad);
23         return 0;
24     }
25     ...
26     if (likely(skb_cloned(skb) || ntail > 0)) {
27         err = pskb_expand_head(skb, 0, ntail, GFP_ATOMIC);
28         if (unlikely(err))
29             goto free_skb;
30     }
31     err = skb_linearize(skb);
32     if (unlikely(err))
33         goto free_skb;
34     memset(skb->data + skb->len, 0, pad);
35     return 0;
36 free_skb:
37     ...
38     kfree_skb(skb);
39     return err;
40 }
```

## 第 4 节

# Implementation

# Implementation



- 工具链: LLVM Pass<sup>3</sup> 和 Python 脚本
- 思路: 排除掉上述两种情况造成的误报后, 对剩余的路径进行模式匹配, 判断是否正确回收内存
- 过程:
  - ① 分析 call 指令的参数, 追踪已分配内存对象的指针在函数之间的传递
  - ② 处理 ownership 的转移, 排除掉上述两种情况造成的误报
  - ③ 从剩余的路径提取 <call FOI, check, call release, return> 四元组, 进行模式匹配, 判断是否正确回收内存
  - ④ 对于漏洞样本进行手工分析, 向 Linux 上游提交 patch

---

<sup>3</sup><https://llvm.org/docs/WritingAnLLVMPass.html>

## 第 5 节

# Evaluation

# Scalability



- 实验对象：5.2.13 版本的 Linux kernel 进行，将其编译为 allyesconfig<sup>4</sup>版本的 bitcode
- 实验环境：48 核 Intel Xeon CPU, 256G 内存的服务器上进行

过程	耗时	备注
生成 bitcode	5 hours	可重用
收集内存分配函数	2 hours	可重用
内存回收函数模式匹配	1 hour	
对 FOI 的检查	几分钟到 4 个小时不等，平均 3 分钟	可以独立并行

<sup>4</sup>编译尽可能多的内核驱动模块

## Set of Allocations and Associated Deallocation



- 从 4621 个备选函数中，提取出了 807 个内存分配/回收函数对<sup>5</sup>，15 个为误报
- 其中 21 个为 primitive allocators，即实际操作内存的函数，其余大部分为 specialized allocator

*To the best of our knowledge, none of the previous detection techniques used such a rich set of allocation-deallocation functions.*

---

<sup>5</sup><https://github.com/Navidem/k-meld/blob/main/results/FOIs.txt>

## Bug finding



- 一共报出 458 个 warnings, 手工分析后认为 218 处存在内存泄漏漏洞, 即误报率为 52%
- 向 Linux 内核提交了 107 个 patch, 申请到了 41 个 CVE
- Bug 分布在各类函数中, 例如 kmalloc (9244 call-sites), sync\_file\_alloc (2 call-sites)
- 有 115 个 bugs 来自于 specialized modules ( $\text{call-sites} \leq 400$ ) 的函数

# Exploitability Analysis



- 使用 fuzzing 或者符号执行的方法花销很大
- Linux 内核中存在三类用户可控输入，即 entry points
- 构建 CFG，找到 entry point 和 leak 所在函数的最短路径
- 83.9% 的 bugs 能够被本地用户输入触发

Entry functions	Attacker	Count
System calls	Userspace	137
Ioctl handlers	Userspace	173
IRQ handlers	Hardware	173
Reachable from any entry		182 (83.9%)

TABLE I: The number of reachable bugs for different types of entry functions.  
IRQ = Interrupt request, Ioctl = I/O control.

# False Positive&Negative Analysis



CVE #	Reproducible	K-MELD Success/Failure
CVE-2019-8980	✓	Success
CVE-2019-9857	✓	Success
CVE-2019-16995	✓	Failure
CVE-2019-16994	✓	Success
CVE-2019-15916	✓	Success
CVE-2019-15807	✓	Success
CVE-2019-12379	✓	Success (correctly rejected)
CVE-2018-8087	✓	Success
CVE-2018-7757	✓	Success
CVE-2018-6554	X	—
CVE-2016-9685	✓	Success
CVE-2016-5400	✓	Success
CVE-2015-1339	X	—
CVE-2015-1333	X	—
CVE-2014-8369	✓	out-of-scope
CVE-2014-3601	✓	out-of-scope
CVE-2010-4250	X	—

## • False Positive

- 不可行的执行路径，可以通过符号执行来缓解
- 对于 specified 内存回收函数的识别错误

## • False Negative

- 过于复杂的指针传递
- Linux 内核版本不同

# Effectiveness of Escape and Consumer Analysis



- 在分别禁用 escape analysis 和 consumer analysis 两种分析的情况下，工具报告的数量上升到 **1292** 和 **1386** 个
- 随机抽取 20 个新增的报告手工分析均为误报，证明两种方法确实有效

## 第 6 节

# Conclusion

# 简单复现



- 论文代码: <https://github.com/Navidem/k-meld>
- 从大型项目构建 bitcode: <https://github.com/travitch/whole-program-llvm>
- 尝试分析一下 FreeBSD: [https://github.com/travitch/whole-program-llvm/blob/master/doc/tutorial-freeBSD.md<sup>6</sup>](https://github.com/travitch/whole-program-llvm/blob/master/doc/tutorial-freeBSD.md)

---

<sup>6</sup>K-MELD is also extendable to other OS kernels like FreeBSD.

# 阅读感想



- Linux 内核非常复杂，漏洞涉及到的调用链长，漏洞验证很麻烦
- 专注于某一类漏洞（memory leak, memory corruption, data race, lock 等）的研究
- 作者在 Linux 内核安全研究领域引发的一些学术伦理讨论<sup>7</sup>

---

<sup>7</sup><https://www.inforsec.org/wp/?p=4893>



上海交通大学  
SHANGHAI JIAO TONG UNIVERSITY

谢谢

谌翔 杨光 · Detecting Kernel Memory Leaks in Specialized Modules With Ownership Reasoning