

Beyond Control: Exploring Novel File System Objects for Data-Only Attacks on Linux Systems

Jinmeng Zhou
jinmengzhou@zju.edu.cn
Zhejiang University
China

Jiayi Hu
2020141530121@stu.scu.edu.cn
Sichuan University
China

Ziyue Pan
ziyuepan@zju.edu.cn
Zhejiang University
China

Jiaxun Zhu
sevenswords@zju.edu.cn
Zhejiang University
China

Guoren Li
gli076@ucr.edu
University of California, Riverside
USA

Wenbo Shen
shenwenbo@zju.edu.cn
Zhejiang University
China

Yulei Sui
y.sui@unsw.edu.au
University of New South Wales
Australia

Zhiyun Qian
zhiyunq@cs.ucr.edu
University of California, Riverside
USA

ABSTRACT

The widespread deployment of control-flow integrity has propelled non-control data attacks into the mainstream. In the domain of OS kernel exploits, by corrupting critical non-control data, local attackers can directly gain root access or privilege escalation without hijacking the control flow. As a result, OS kernels have been restricting the availability of such non-control data. This forces attackers to continue to search for more exploitable non-control data in OS kernels. However, discovering unknown non-control data can be daunting because they are often tied heavily to semantics and lack universal patterns.

We make two contributions in this paper: (1) discover critical non-control objects in the file subsystem and (2) analyze their exploitability. This work represents the first study, with minimal domain knowledge, to semi-automatically discover and evaluate exploitable non-control data within the file subsystem of the Linux kernel. Our solution utilizes a custom analysis and testing framework that statically and dynamically identifies promising candidate objects. Furthermore, we categorize these discovered objects into types that are suitable for various exploit strategies, including a novel strategy necessary to overcome the defense that isolates many of these objects. These objects have the advantage of being exploitable without requiring KASLR, thus making the exploits simpler and more reliable. We use 18 real-world CVEs to evaluate the exploitability of the file system objects using various exploit strategies. We develop 10 end-to-end exploits using a subset of CVEs against the kernel with all state-of-the-art mitigations enabled.

CCS CONCEPTS

• Security and privacy → Operating systems security.

KEYWORDS

Linux Kernel Security; Data-only Attack; File System

1 INTRODUCTION

Control flow hijacking used to be the mainstream attack method to exploit memory errors; however, it has become difficult to achieve after the Control Flow Integrity (CFI) defense was introduced. As a result, data-only attacks have become popular [55, 58, 67]. For example, a common data being targeted in the Linux kernel is `modprobe_path` [21, 25, 29, 79, 30, 47, 64, 63]. As long as an attacker can overwrite its value with a file path that points to the attacker-controlled executable file, it will be executed with the root privilege [57]. There is no need to collect and stitch together ROP gadgets in such data-only attacks. However, `modprobe_path` is a global variable, and overwriting it requires defeating the KASLR defense, which randomizes the address of all global variables. This is an extra complexity that can be a hurdle for exploits and exploit reliability.

Alternative to global variables, heap variables also contain critical data that can be leveraged to achieve privilege escalation. They are easier to exploit for common types of vulnerabilities, such as out-of-bound (OOB) memory access and use-after-free (UAF). This is because such vulnerabilities offer relative heap memory write, whereas KASLR does not help. For instance, attackers can spray `struct cred` objects and use OOB or UAF vulnerabilities to overwrite their `uid` fields to `0`, achieving the root privilege [44]. However, because of their popularity, such objects are restricted nowadays. For example, they are now isolated from the other heap objects and, therefore, much more difficult to exploit [68], requiring unstable cross-cache attacks. In Android kernels, they are even protected by the hypervisor [48, 61]. Because the well-known critical objects are specifically protected through advanced defenses [36, 3, 48, 61], attackers no longer target them even if they manage to achieve an arbitrary write primitive.

To fill the gap, we seek to identify additional heap objects suitable for data-only attacks. Unfortunately, discovering unknown non-control data (heap objects or not) can be daunting because they are tied heavily to semantics (e.g., user id) and lack universal patterns. In this paper, we perform a systematic investigation of objects associated with the file subsystem in the Linux kernel. This is

because file systems contain many critical files, e.g., `\etc\passwd`. If an attacker can manipulate a file system into overwriting such a critical file, an attacker can achieve privilege escalation directly [56].

To date, few exploitable objects have been identified in the file subsystem of the Linux kernel [40]. This is due to the lack of a systematic and automatic analysis of critical objects based on their semantics. To our knowledge, the only work that attempts to collect critical non-control data automatically in the Linux kernel is KENALI [62]. However, there are two serious limitations of the work: (1) its heuristic is overly generic (not specific to the file subsystem) and misses or falsely identifies many objects; (2) it does not analyze the exploitability of the identified objects.

Our solution is two-fold. First, we leverage minimal domain knowledge of the file subsystem in the Linux kernel to identify file system key objects (FSKO) with critical fields more accurately and completely. Specifically, the process starts with a few anchor objects in abstract layers of the file subsystem. Then, it identifies additional related objects in other layers through tailored propagation rules. We propose a framework named FSKO-AUTO that integrates static analysis with dynamic testing as a strategy to identify FSKOs and verify that their semantics are appropriate for exploitation. In total, FSKO-AUTO finds critical 23 fields within heap objects whose corruption directly achieves privilege escalation. Second, we map these objects into three targeted exploit strategies suitable for achieving privilege escalation using FSKO objects. Since many FSKOs are located in isolated slab caches, we leverage a novel exploit strategy that uses “bridge objects” to construct a powerful and general primitive of page UAF, which can reliably exploit FSKOs without having to resort to cross-cache attacks.

Out of 26 recent CVEs, our solution could demonstrate exploitability for 18 of them, implying its real-world impact and generality. For a subset of the 18 CVEs, we managed to write 10 end-to-end exploits, achieving privilege escalation by corrupting the FSKO objects with various exploit strategies. The results show that we can write diverse exploits without having to bypass KASLR, which makes the exploits simpler and more reliable than previous methods.

In summary, the paper makes three contributions.

- **Semi-automatic FSKO identification and confirmation.** With the help of static analysis and dynamic confirmation, we are able to drastically reduce the manual effort in identifying and confirming FSKOs. This led to 23 unique FSKO fields, covering all publicly exploited objects in the file subsystem.
- **Exploitability analysis and novel exploit strategy.** We analyze the exploitability of these FSKOs when paired with vulnerabilities of different capabilities and slab cache requirements. We propose a novel strategy that enables the many FSKOs to be exploitable, despite the fact that they are in isolated slab caches.
- **Practical evaluation and findings.** We confirmed that 18 out of 26 CVEs are exploitable using the FSKOs and further developed 10 end-to-end exploits for a subset of them. The identified objects allow our exploits to avoid bypassing KSLARs and be future-proof against CFI and advanced protection of well-known data.

2 MOTIVATION AND THREAT MODEL

Motivation. After control flow attacks are mitigated, data-only attacks have gained popularity. These attacks specifically target

and corrupt a program’s data rather than its control flow [33]. In the context of the Linux kernel, data-only attacks often target two types of data: (1) global variables and (2) heap variables. For global variables, the most widely targeted variable is `modprobe_path` [21, 25, 29, 79, 30, 47, 64, 63]. The variable is a string, the corruption of which can lead to the kernel executing an attacker-controlled executable file as root [57]. For heap variables, the most common ones are the `cred` [40] and the page table [23]. Corrupting them can lead to privilege escalation directly.

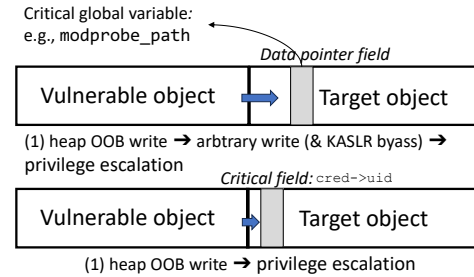


Figure 2: Corrupting critical heap object vs. corrupting critical global variable.

Using a common type of kernel vulnerability — heap out-of-bound write — as an example illustrated in Figure 2, achieving privilege escalation by corrupting critical global variables is significantly more cumbersome than corrupting critical heap variables. As shown in Figure 2(a), there are two extra steps: (1) deriving an arbitrary write primitive by corrupting a data pointer in a heap object and (2) bypassing KASLR to obtain the correct address of the global variable. In contrast, Figure 2(b) illustrates how overwriting a critical heap variable can directly lead to privilege escalation.

The example assumes the heap variables can be co-located with the vulnerable object, which may not always be possible. This is especially the case when defenses increasingly target the commonly exploited objects, e.g., isolation of `struct cred` in dedicated slab caches. It usually requires `infoleak` to obtain their address and arbitrary write primitive to corrupt them [19]. More importantly, commercial kernels apply advanced defenses, e.g., changing kernel configs and using hypervisors to deny illegitimate accesses, to prevent corruption of the critical objects in the dedicated slab cache [36, 3, 48, 61]. This motivates the discovery of additional critical heap objects.

Threat Model. We assume an unprivileged user aims to achieve local privilege escalation by exploiting a memory error in the kernel. Specifically, the error allows limited memory writes (e.g., OOB write or UAF write), possibly as limited as a single-bit write. We *do not* assume the availability of arbitrary memory write primitives, i.e., write anywhere in memory.

We assume the Linux kernel enables all modern mitigations, including CFI [41, 49], W@X [71], KASLR [34], SMAP/SMEP (x86) [51, 42, 43], PAN/PXN (ARM) [43], and KPTI [50]. An unprivileged user cannot modify kernel code or inject code into the kernel data segment, and the kernel address is randomized. The kernel is restricted from accessing the data or executing the code in user space.



Figure 1: The work overview of identifying FSKOs spreading all layers in the file system.

Control flow can not be hijacked. Besides, we assume advanced defenses of well-known non-control data are safeguarded, including `modprobe_path`, `cred`, and page table [36, 48, 61].

Attack Goal. Our goal is to achieve privilege escalation directly by corrupting non-control heap objects in the file subsystem of the Linux kernel. We envision two high-level approaches to achieve the goal. First, an attacker can corrupt certain heap objects (e.g., representing file permissions or owners) to gain write access to files that are read-only to the attacker [4, 22, 40], e.g., `/etc/passwd`. Second, an attacker can corrupt certain heap objects to turn non-setuid-root executable files into *setuid root* [52]; the attacker can succeed when the executables are either controlled by the attacker or already has sufficient functionalities, e.g., `/usr/bin/vi`.

In this paper, we focus on extracting heap objects in the file subsystem that match the above semantic descriptions in §3 suitable for exploitation and evaluating the practical uses them in §4.

3 FSKO IDENTIFICATION

In this section, we focus on identifying objects in the file subsystems that can lead to privilege escalation when specific fields are corrupted. We refer to them as **FSKOs**, i.e., **F**ile **S**ystem **K**ey **O**bjects. There are two main challenges in FSKO identification:

Challenge 1: There are a large number of unique objects in the file subsystem of the Linux kernel, i.e., 3,553 structs and 554,161 fields in kernel v5.14. Enumerating them manually is infeasible.

Challenge 2: These objects are scattered across several layers, encompassing both abstract layers (e.g., VFS) and concrete file system implementations (e.g., `ext2`, `ext3`). This makes it difficult to have a clean heuristic that works across layers.

To this end, we propose a new framework, FSKO-AUTO, that combines static analysis and dynamic verification in a structured pipeline, as shown in Figure 1. For the first challenge, we propose to focus on identifying a small set of anchor objects with specific semantics and in abstract layers only to bootstrap the identification process.

For the second challenge, from the anchor objects in abstract layers, we develop custom static analysis rules to search across other layers automatically. Finally, we develop a dynamic verification process to confirm the statically identified objects are, in fact, with the right semantics, and that corruption of them leads to privilege escalation.

3.1 Anchor Object Discovery

We define *anchor objects* to be FSKOs with fields that (1) hold the semantic that has security implications – we consider specifically three classes of semantics, and (2) reside in the abstract layer, making it more manageable to analyze, compatible with different implementations.

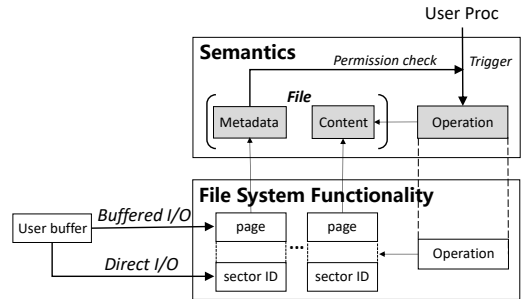


Figure 3: File system functionality in the abstract.

Classes of Semantics. As shown in Figure 3, file systems have a few basic design elements. In general, a file consists of two parts: *file content* that stores the actual data; *file metadata* that encompasses properties like creation time, name, ownership, and permission. Furthermore, there are *operations* performed over files, e.g., read, write, and execute. Naturally, they correspond to the three classes of semantics that matter to privilege escalation. We consider this minimal domain knowledge necessary for identifying critical objects in the file subsystem.

(1) *Metadata:* Unprivileged users can corrupt permission or ownership objects to obtain authorization over files that are not owned by the attacker. For example, an attacker can weaken the permission or change the ownership so that a sensitive read-only file is writable by the attacker [4]. They can also create a new file with malicious payload and turn it into a *setuid root* executable which will execute as root [52].

(2) *Content:* Corrupting content of sensitive files can obviously lead to privilege escalation, e.g., `/etc/passwd`. In terms of heap objects, we know that file content can be represented as caches in memory [22], for buffered I/O. In the case of direct I/O, the write will occur directly on the hard drive.

(3) *Operation:* An attacker can corrupt a file operation (represented by some heap object) from a read or an execute into a write, resulting in the corruption of file content indirectly (we have not seen any instance of similar strategies used in practice). This requires the attacker to have read or execute access already but not the write access. Note that these operations are performed only after successful permission checks. This means that an attacker cannot prematurely corrupt the operation and has to wait until the permission check is passed [40].

Abstract File System Layers. As part of the domain knowledge, the Linux file subsystem is designed to be layered, much like any complex system. As shown in Figure 4, there are four layers from the top to the bottom: virtual file system (VFS), file system implementations, generic block, and drivers. Specifically, we consider

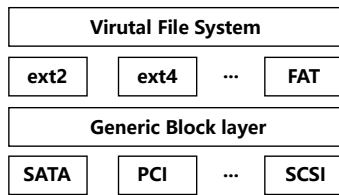


Figure 4: Simplified file system layers in Linux kernel.

```

1 static int acl_permission_check(struct inode *inode, int mask) {
2   unsigned int mode = inode->i_mode;
3   if (likely(uid_eq(current_fsuid(), inode->i_uid))) {
4     mask &= 7;
5     mode >>= 6;
6     return (mask & ~mode) ? -EACCES : 0;
7   }
8   ...
9 }
10
11 int generic_permission(struct inode *inode, int mask) {
12   int ret;
13   ret = acl_permission_check(inode, mask);
14   ...
15   mask &= MAY_READ | MAY_WRITE | MAY_EXEC;
16 }
  
```

Figure 5: DAC permission check validates the user-intended permission using defined macro (v5.14).

two well-known *abstract FS layers* in the Linux kernel: the *VFS layer* and the *generic block layer*. The former sits immediately below the syscall interface, e.g., invoked by `open`, `read`, and `write`. The latter sits below the layer of file system implementations that calculate file location on disk, so it connects the memory page and the disk sector. Both layers are agnostic to different underlying file system implementations (e.g., `ext2`, `ext4`) and various disk drivers.

We choose abstract layers as opposed to the concrete file system implementation layers because of the smaller and more manageable scope. It is also because sensitive FSKOs must already be defined in the abstract layers. Furthermore, such FSKOs may be propagated into various concrete file system implementation layers, which we will capture automatically through static analysis.

Object Identification. To identify objects relating to the three classes of semantics described earlier, we first identify the syscalls that can potentially affect the corresponding semantics (e.g., `open()` relates to metadata). We then review the functions involved in the syscalls that fall under the two abstract layers. For a specific class of semantic, we also follow certain heuristics to help us narrow down the search scope.

(1) *Metadata:* We rely on well-defined keywords and macros that identify the permission objects in the VFS layer. Specifically, we grep variable names and fields with “uid” and “gid” that identify `inode->i_uid` and `inode->i_gid` as the user ID and group ID of the owner of the file. Regarding permission objects, we reuse macros from a previous work aiming to identify permission-check-related kernel functions [78]. We limit our scope to macros related to discretionary access control (DAC) (e.g., no SELinux policy considered). We include the three key macros: `MAY_EXEC`, `MAY_READ`,

and `MAY_WRITE`, as well as `S_ISUID` and `S_ISGID` representing `setuid` and `setgid` permissions. From these macros, we then statically track what variables are “tainted” by them using a simple data flow analysis. The `mask` variable in line 15 in Figure 5 illustrates an example variable identified this way. We additionally track what variables are involved in bitwise operations with the previously identified variables. The `mode` variable in line 6 corresponds to this. All the objects that contain such variables are then recorded as anchor objects, e.g., `inode`, whose `i_mode` field corresponds to the mode local variable.

(2) *Content:* As alluded to earlier, file content can be represented in memory (buffered I/O) or on disk (direct I/O). Specifically, the in-memory caches are accessible via data pointers; the on-disk file content is addressable by sector numbers. In reviewing the syscall documentation, we realize that a flag in the `open()` syscall controls the access mode, and therefore, we consider the access mode separately for `read()` and `write()` syscalls, i.e., we review their call graphs conditioned on the access mode. In summary, we find page caches accessed through `inode->i_mapping` and `file->f_mapping` in the VFS layer, which is one anchor object. We find an anchor object in the generic block layer that contains the sector numbers and interacts with the disk, `bio->bi_iter->bi_sector`.

(3) *Operation:* We primarily focus on reviewing the `read()` and `write()` syscalls as the file operation is likely encoded in some form internally in the two abstract layers. To be more efficient, we compare the call graph of the two syscalls and focus on the functions that are shared by both and stop when the call graphs diverge. The intuition is that the divergence should exist before the divergence that leads to the separate functions being invoked. For buffered I/O, we find that there is no object representing file operations as the kernel quickly diverges by calling `vfs_read()` and `vfs_write()` to distinguish the read and write operations. On the other hand, for direct I/O, we find there is an object called `bio->bi_opf` in the generic block layer.

3.2 Cross-Layer FSKOs Discovery

In this step, we develop custom static analysis rules that track the propagation of anchor objects into other concrete implementation layers. In order to track the flow of these objects properly, these rules operate forward and backward. Ultimately, the key is to identify other objects that carry the same or similar semantics. Our static analysis takes critical fields identified in anchor objects as input and outputs the derived fields in other objects. At a high level, the rules explore the data dependency and control dependency of the anchor objects with specific considerations of file system characteristics. The static analysis consists of 3K LoC using LLVM/Clang 14.0.0.

Definitions. We categorize file metadata, content, and operation into two types of data, *flags* and *references*, to simplify the description of our algorithm. *Flags* represent file metadata and operations; overwriting the corresponding fields of such objects or object fields with constants (permission of 777 allows anyone to read/write a file) can directly lead to privilege escalation. *References* represent file content, i.e., file cache (pointers) and disk locations (sector numbers), which correspond to buffered I/O and direct I/O, respectively. Modifying the reference from pointing to a regular writable file (low-privilege) to a read-only (high-privilege) file can achieve root.

Table 1: Custom Propagation Rules.

	Data dependency	Control dependency
Flag	Assignment (optionally involving logical ops)	Value influence by control
Reference	Assignment (optionally involving arithmetic ops)	Reference assignment influenced by a flag variable

Scope. In practice, both data types above are represented as fields encapsulated in FSKOs (structs). Therefore, our analysis starts tracking the propagation of such fields throughout the file subsystem – to see how they end up being associated with fields in other structs. Note that we intentionally exclude the tracking of FSKO structs themselves, which may be address-taken and assigned to a pointer field of another non-FSKO struct. We call such non-FSKO structs *wrapper* objects, i.e., an object that has a pointer pointing to an FSKO. Even though an attack can corrupt the pointer in a wrapper object and force it to point to a new object, we find that oftentimes such a whole object replacement can lead to kernel instability or broken functionalities. For example, one can replace a pointer field in a wrapper object of file: `file->f_inode`, to attempt to write to a sensitive file. However, we find the kernel loses the ability to write to this file in the future regardless of who opens it (including root users). These anomalies could jeopardize the attack’s stealth, leading to its detection.

Analysis. First of all, to analyze a complex subsystem like the file system, we already face a significant challenge of pointer analysis, which is a foundation of data flow analysis. This is because, especially, there are multiple syscalls that are dependent on each other, e.g., `open()`, `read()`, `write()`, and `mmap()`. One syscall may store something on the heap, and another syscall may retrieve it subsequently (i.e., aliases across syscalls). We borrow the technique proposed in prior work, which uses access-path-based approaches (considering multiple levels of types) to identify aliases across syscalls [77]. In addition, we use the state-of-the-art indirect call analysis [45] that prunes many false indirect call edges.

Rule 1: Flag Propagation. We summarize the rules for tracking flags and references in Table 1. In the case of tracking a flag variable (e.g., a field within an anchor object), it is considered to potentially influence or be influenced by another variable through data or control dependency. Specifically, we consider both forward and backward dependencies. We customize the data dependency by tracing the direct assignment (e.g., `a = b`) and the assignment involving logical expressions (e.g., `a = b & STATE`), including `AND(&)`, `OR()`, `NOT(~)`. This is based on the observation that flags typically involve bitwise operations. Arithmetic operations are not included because flags are unlikely to be used in addition or subtraction operations. This helps us eliminate false positives.

Rule 2: Reference Propagation. The rules tracking a reference variable are also shown in Table 1. First, we consider a reference object propagates its particular value to another reference through solely data dependency (an exception will be discussed later). It is obvious if other variables also store sector numbers and page pointers, e.g., due to assignments, then clearly, these other variables will also be considered references. Here, we also customize

the data dependency rules, by allowing arithmetic operations because a reference can be used to add a relative value to find the neighbor reference. This customization excludes other types of data dependency, such as those bitwise operations in logical expressions.

Second, in addition to data dependencies, we also consider control dependencies in a very limited capacity. This is because there can be *flag* variables that determine a reference’s value through control dependency (e.g., whether a reference will now point to a new page cache, e.g., DirtyCow [22]). We generalize this pattern to consider whether any variables in a condition can lead to two branches where (1) at least one of the branches sees an assignment to a reference, (2) the assignment differs in the two branches – one branch missing an assignment or the reference is assigned two different values in the two branches.

3.3 Dynamic Verification

After identifying all candidate FSKOs and their specific fields, FSKO-AUTO dynamically verifies if they are true positives by forcefully corrupting their values at runtime and observing the effects relating to privilege escalation. Dynamic verification is non-trivial and necessary before writing the exploit of FSKOs. This is because we observe some objects seem to contain critical fields, but they can not lead to privilege escalation due to various reasons. One reason is that a path may be infeasible or require privileges to reach. For example, our static analysis observes that `iattr->ia_mode` is assigned to `inode->i_mode` where the latter field represents the file permission. This makes the former object seemingly exploitable since we have an opportunity to manipulate the permission. However, it turns out that only the file’s owner can trigger the assignment when setting the file attributes, making the object non-exploitable.

Another reason is that the dependency relationship does not guarantee the same semantics, especially in the case of objects identified by control dependency. As shown in Figure 7, `filp->f_flags` is identified as a candidate FSKO field (with the FSKO being `struct file`) because it can influence the value of `file->f_mode` through control dependency. However, the former has a completely different semantic than the latter. Furthermore, even if the former can influence the value of the latter, the values assigned in lines 2 and 4 are not useful in making the file writable.

We develop a process to perform the dynamic verification. To begin the process, we first collect userspace test cases to cover the file-related operations such as `read`, `write`, and `mmap`. This is necessary for us to observe and modify the values of a specific FSKO field in the kernel. We rely on two sources: (1) syzkaller traces collected during our fuzzing campaign and (2) manually curated test cases. We need to curate test cases because the test cases created by the vanilla syzkaller are insufficient. For example, syzkaller does not create multiple file systems to exercise the data structures specific to each file system (e.g., `ext2`, `ext3`). Fortunately, the number of file system-related syscalls is limited, and we managed to cover more than half of FSKO field candidates in the end (see details in §5).

Next, for the uses of specific FSKO fields covered by test cases at runtime, we simulate a successful attack by corrupting specific fields in FSKOs. At first glance, it is challenging to simulate a successful attack because we do not know what value to write into a specific field of an FSKO. To overcome this issue, we apply differential

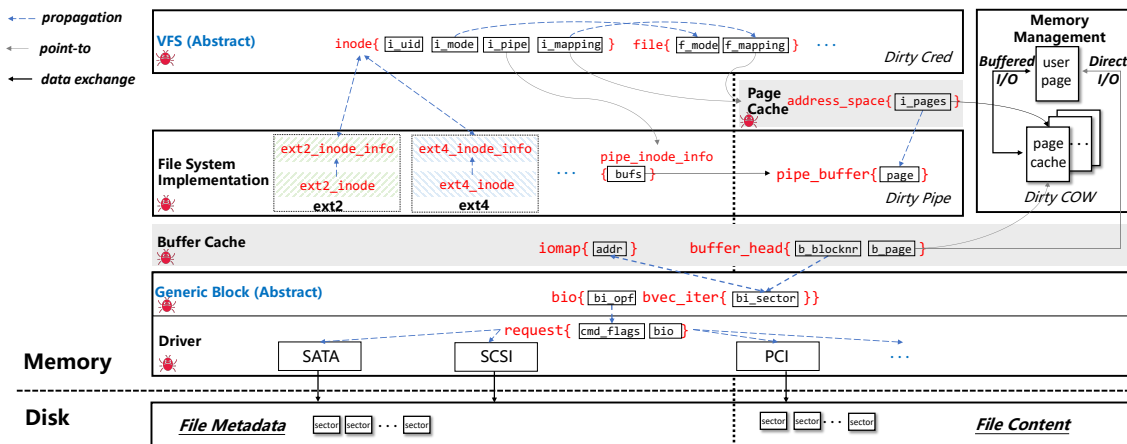


Figure 6: Some key objects identified by FSKO-AUTO are highlighted in red, propagating all file system layers.

```

1 if (filp->f_flags & O_EXCL)
2   filp->f_mode |= FMODE_EXCL;
3 if ((filp->f_flags & O_ACCMODE) == 3)
4   filp->f_mode |= FMODE_WRITE_IOCTL;

```

Figure 7: The false positive introduced by control dependency

analysis to automatically retrieve the appropriate value that can likely lead to privilege escalation. We apply this strategy to the two conceived attack approaches described in the attack goals of §2. Using the approach of “turning a read-only file into a writable file” as an example, below is how we perform the differential testing:

- (1) Record the value of a specific FSKO field of interest (e.g., `file->f_mode`) in a successful write operation (against a test file we create).
- (2) Replace the value of the same FSKO field in a read operation (e.g., against `\etc\passwd`) with the value recorded in the write operation. In the case of `file->f_mode`, this effectively changes the file permission from read-only to writable.
- (3) Observe whether the read operation is now turned into a successful write operation. If so, the FSKO field candidate is deemed a valid one.

3.4 Summary of Identified Objects

FSKO-AUTO initially identified 258 candidate FSKOs with static analysis. After we apply dynamic verification, 25 heap FSKOs and 8 stack FSKOs are confirmed to be effective for privilege escalation, we focus on the heap objects in this paper and will discuss them in more detail in §5.

In Figure 6, we illustrate and categorize some representative objects that are identified across various layers of the Linux kernel file subsystem, from the VFS layer at the top all the way to the driver layer at the bottom. We can see the left part represents the metadata, while the right side represents the content, with a dotted vertical line separating them. Taking `inode->i_mode` as an example anchor representing file permission, it propagated to another object `file->f_mode` which represents permission as well.

Of course, to corrupt these objects, we need to first be able to allocate them. For struct `inode`, we find it is possible to allocate it individually and as part of bigger objects, e.g., `ext2_inode_info`, that are specific to file system implementations. At the lower layers of the file subsystem, we started with the anchor of sector number `bio->bi_iter->bi_sector` but found other object fields with similar semantics. For example, `buffer_head->b_blocknr` represents sector numbers as well in the buffer cache-related data structure (they are used to map higher-level concepts such as pages into lower-layer concepts such as sector numbers). Ultimately, these additionally discovered objects and fields give the attack significantly more flexibility in the object to choose in heap fengshui. We will discuss a more complete list of FSKOs in Table 4.

It is worth noting that all key objects leveraged in prior exploits (which also target file system-related objects), including DirtyPipe, DirtyCow, and DirtyCred, are discovered in our work. For example, we find `pipe_buffer->flags` and `vm_fault->flags` which are used by DirtyPipe and DirtyCOW to manipulate flags that control which reference object is used. We also discovered `file->f_mapping` which represents a file page. One exception is that DirtyCred also targeted wrapper objects of struct `file`, i.e., objects that contain a pointer to struct `file` objects. As mentioned, we intentionally exclude any wrapper object in our workflow. As discussed in §3.2, swapping wrapper objects can lead to system stability issues. Nevertheless, in the case of struct `file` specifically, it turns out that it is safe to swap. We leave the investigation of wrapper objects to future work.

4 EXPLOITABILITY ANALYSIS

So far, we have gathered a list of FSKOs with specific critical fields for which we know their corruption will lead to privilege escalation. The dynamic verification step, however, only verified the semantics of FSKO fields by forcefully overwriting them. In the real world, we need to pair these objects with specific vulnerabilities (e.g., OOB and UAF). In general, the more FSKOs, the more likely we can pair one with a given vulnerability successfully. We show the exploitability analysis overview in Figure 8, including flag and reference corruption. Besides, we propose a new page-UAF strategy

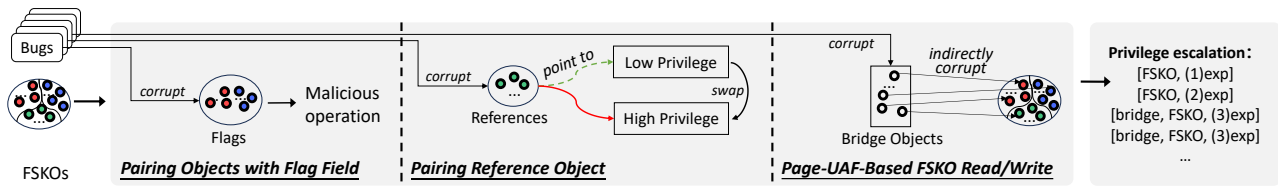


Figure 8: The work overview of exploitability analysis to find effective objects.

Table 2: The requirements to corrupt FSKOs.

Object Types	Information Leakage	Memory Corruption Capability
Flags	/	Limited write
File cache pointers	Required	Write controllable value
File cache pointers	/	Write last few bits/bytes
Disk sector numbers	/	Write controllable value

to pivot OOB/UAF bugs into a page UAF to then subsequently corrupt both flag and reference.

Pairing FSKOs with vulnerabilities. There are two key dimensions related to pairing objects with vulnerabilities:

(1) Capability of a vulnerability. As studied in [9], different vulnerabilities can exhibit different write capabilities, e.g., offset and length of an OOB write, as well as the possible values of the write (see Figure 2). Depending on the specific capability, we might choose to pair an FSKO with fields at the right offset.

(2) Slab cache of the vulnerable object. In heap exploits, a target object (e.g., FSKO) needs to be placed in the same slab cache (e.g., `kmalloc-512`) that the vulnerable object resides in [11]. The more FSKOs, the more likely we can find an object that co-locates with the vulnerable object. Otherwise, an attacker would resort to cross-cache attacks, which is considered much less reliable [31].

4.1 Pairing Requirement #1: Write Capability

Pairing Flag Objects. Objects with flag fields are straightforward. As long as they can be overwritten as pre-determined constant values, an attack would succeed. The only requirement is that the vulnerability (e.g., OOB write and UAF write) can overwrite at the specific field offset with an expected value. In the case of permission and `setuid`, it is sufficient to set a single bit (e.g., in `inode->imode`) to turn a file into writable by all users or `setuid` root executable (created by the attacker). Similarly, in the case of file operations, it is sufficient to set a single bit (e.g., in `bio->bi_opf`) to turn a file access request from read to write. In the case of owners, an attacker can write 0 to the owner field (e.g., `inode->i_uid`) of a `setuid` file created by the attacker, turning it into a `setuid` root executable; alternatively, an attacker can write 1000 to the owner field so that the attacker user becomes the owner of a critical file.

Note that for the FSKO fields derived from anchor objects, their semantics could change and we cannot write the same constants as if they are anchor objects. Nevertheless, as mentioned in §3.3, we can retrieve the proper values for such objects at runtime (during a successful write operation) and reuse them. We list the exact

constant values expected for a variety of flag objects later in Table 4. In summary, *flag* objects are generally easy to pair, requiring only limited write capabilities, e.g., setting a specific bit or writing a few bytes with constants or attacker-controlled values.

Pairing Reference Objects. As mentioned earlier, reference objects include data that represent file content, i.e., file page cache pointers or sector numbers. Generally speaking, we envision the proper exploit strategy to be swapping the references so that they point to a sensitive file for corruption.

To overwrite file cache pointers, the most straightforward approach is to use a specific file page address corresponding to a sensitive file. However, this approach requires an extra step of infoleak. In addition, it also requires the value of the write primitive to be controlled by the attacker because the address of a file page is dynamically determined on the heap. Alternatively, one can spray a large number of file pages corresponding to sensitive files so that they are co-located with the file page of an attacker-owned file. Overwriting the lower bits of the file page pointer (e.g., setting the lower two bytes to 0) will cause it to point to a nearby page. This method does not require any infoleak, and the write capability is also less demanding. This method has been used in prior work, e.g., Dirtycred [40]. However, this method places requirements on vulnerability writing capabilities, and some do not support only overwriting the lower bits. However, we need to place the page caches of the low-privilege and high-privilege files adjacent to each other. This highly requires manipulation of the memory layout, which is difficult and unstable. We can use page UAF to overcome this problem and corrupt references.

To overwrite sector numbers, we find that it does not require any infoleak because unprivileged users in Linux can obtain the sector numbers of any file through an `ioctl()` syscall with `FS_IOC_FIEMAP`. The syscall can retrieve the logical and physical block numbers [26], and the sector numbers can be computed by shifting several bits of the block numbers. This method requires a write capability with attacker-controlled values, as the sector numbers of sensitive files can be various integer values on different systems.

We summarize the write capability requirement of different types of objects in Table 2.

4.2 Pairing Requirement #2: Slab Cache

As mentioned before, we would like a diverse set of FSKOs that exist in various slab caches so that they can be co-located with a vulnerable object in a given vulnerability. As will be shown later in Table 4, most FSKOs are unfortunately located in dedicated slab caches, e.g., `ext4_inode_cachep` that hosts only a single type of objects. This makes them unlikely to be co-located with a vulnerable

object. To overcome this, we come up with a novel and general exploit strategy that can indirectly achieve complete FSKO read and write without resorting to cross-cache attacks, called **Page UAF strategy**. Besides, it also does not rely on a pre-existing information leakage primitive to bypass KASLR.

Page UAF Strategy. Given that most FSKOs are located in dedicated slab caches, we will pursue other non-FSKOs that are located in more standard slab caches (e.g., `kmalloc-512`), and by corrupting them we can reach a new primitive, i.e., page-level UAF, sufficiently strong to set up future FSKO read/write (and even other types of objects that are not protected by any advanced defenses). We term these non-FSKOs *bridge objects*. It turns out that there are many suitable bridge objects with variable sizes and can be placed in many different standard slab caches. We illustrate the two main steps in Figure 9 and explain them in detail below:

1. Page-level UAF construction. We consider *bridge objects* to be those that contain `struct page` pointers, e.g., `struct pipe_buffer`. This is because a `struct page` (64 bytes) object corresponds to a 4KB physical memory page. Through the bridge object, an attacker can manipulate (read or write) the physical page. If we can cause a UAF of a page object, the kernel would not only free the 64-byte object but also effectively release the corresponding physical page, creating a page-level UAF.

As shown in Figure 9, to trigger an invalid page free, we can first allocate multiple bridge objects (at least two) that co-locate with the vulnerable object (not shown in the figure). For bugs with standard OOB or UAF write primitives, we can use the primitive to corrupt the page pointer field in a bridge object (e.g., the first field of `pipe_buffer` and `configs_buf`) such that it points to another 64-byte page object nearby. This effectively causes two pointers to point to the same object. A user-space program can trigger `free_pages()` on one of the objects (e.g., by calling `close()`), which will create a dangling pointer to the freed page object and the corresponding physical page. In other words, we can read/write the corresponding physical page that is now considered freed by the OS kernel. For example, one can write to a pipe, which will lead to a write of the physical page, via the `pipe_buffer` object.

For bugs that have double free primitives, which often can be achieved from UAF by triggering the `free()` operation twice. For the first free, we spray a harmless object (e.g., `msg_msg`) to take the freed slot. The object should take attacker-controlled value from user space. Then, we trigger the kernel code to write the harmless object until reaching a certain offset, and use the FUSE technique [27] to stop the writing right before a planned offset – corresponding to the page pointer field of a planned bridge object. Now, we trigger the free for a second time to spray the planned bridge object to take the slot. The writing process is restarted to continue overwriting the lower bits of the page pointer field, which leads to a page UAF. Previously, to trigger page-level frees from double frees, one had to free an entire slab and then do cross-cache technique [40] whereas we do not need to.

2. FSKO corruption through page-level UAF. Now that we have a freed physical page and a dangling pointer that can read/write it. It is fairly easy to then allocate and corrupt FSKO objects. This is because the flag and reference objects are eventually allocated through the buddy allocator at the page granularity. Therefore,

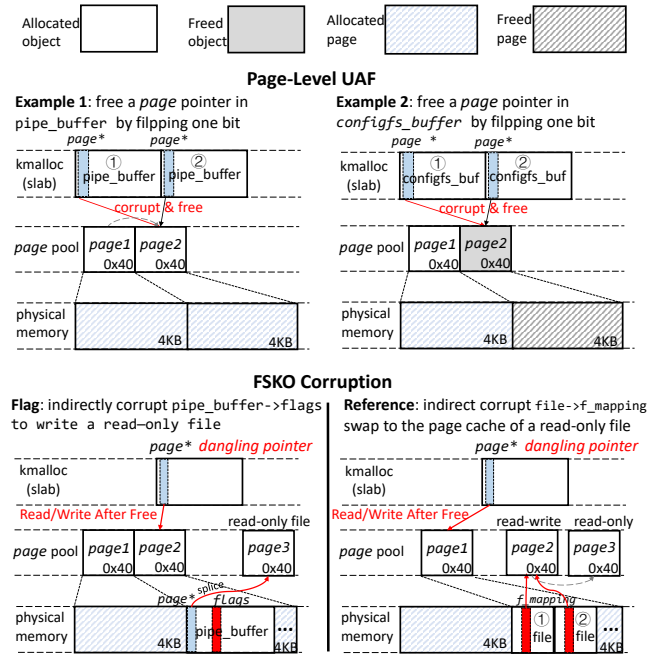


Figure 9: Use page-UAF strategy to corrupt flag and reference object, resulting in writing read-only files.

attackers can spray FSKOs into the freed page, as shown in Figure 9, as long as they have already exhausted all existing slab caches.

Figure 9 showed two example FSKOs that we successfully exploited. In the first example, we allocate `struct pipe_buffer` objects (they are both bridge objects and FSKOs). We then connect the pipe to a read-only sensitive file via `splice()` syscall, and achieve write permission by corrupting `pipe_buffer->flags`. In the second example, we allocate two `struct file` objects. The first object corresponds to a read-only sensitive file, and we then read its `file->f_mapping`, which points to a data structure that represents its file content. We then allocate a second `struct file` object that corresponds to an attacker-created file. Next, we write its `f_mapping` field with the value of the first file object. This allows the attacker to use the second file to write into the first.

Bridge object discovery. We simply query structs that contain `struct page` fields and record those that are found in our fuzzing campaign and manually curated test cases. Note that there are some bridge objects where a user cannot perform read and write operations over the page. We exclude such objects. Finally, we list the exploitable bridge objects in Table 5. Even though the number of such bridge objects is not very large, there are 8 objects that are elastic, i.e., with variable sizes. This makes them ideal for satisfying the slab cache requirement.

To our knowledge, we have not seen widespread use of bridge objects to achieve page-level UAF in real-world exploits. The only example we are aware of is a CTF competition [5] that uses the `struct pipe_buffer` as a bridge object. However, we note that `struct pipe_buffer` is being isolated into `kmalloc-cg` slab using flag `GFP_ACCOUNT` after Linux kernel v5.14. This slab is segregated

Table 3: Statistical results of identified objects.

	Category	#	Security Impact		
			PE	NPE	No test
Flag	Metadata	148	9	61	78
	Operation	29	8	11	10
	COW	5	2	0	3
Reference	Sector	54	10	26	18
	Page Cache	22	4	3	15
Sum	/	258	33	101	124

PE: can lead to Privilege Escalation after dynamic validation.
 NPE: cannot lead to Privilege Escalation after dynamic validation.
 No test: can not dynamically verify.

from the commonly used slab with flag GFP_KERNEL. It may require cross-cache corruption in future exploits.

Summary. The core of the page UAF strategy leverages insufficient page isolation. Unlike the slab caches with over a hundred isolated groups which offer strong isolation (as is the case for most FSKOs), there are only at most six page groups that are isolated from each other [46]. Because of this, the pages whose content is fully controlled by users (e.g., pipe pages) are not isolated from the page used in the slab. In other words, the slabs (that host FSKOs) can reuse the freed user-controlled pages, rendering the page-UAF strategy successful. This strategy also works for not just FSKOs but any other types of objects in isolated slab caches. For example, cred in dedicated caches can be used as well if they are not prevented by advanced protections (e.g., hypervisor monitoring).

5 EVALUATION

Experiment setup. (1) FSKO identification: We performed our static analysis and dynamic testing against the Linux kernel v5.14. (2) Exploitability analysis: We sample 26 recent Linux kernel CVEs that are either OOB, UAF, or double-free from 2020 to 2023. This includes 24 vulnerabilities from prior work [40], and 2 additional OOB ones missed by the prior work.

5.1 Exploitable Objects

Overall, as shown in Table 3, our static analysis found 258 key fields within 142 structure objects. Then, 134 fields are verified dynamically with test cases, out of which 23 fields belonging to 17 FSKOs are confirmed to lead to privilege escalation. We list all these 23 FSKO fields in Table 4 as well as a few other seemingly promising ones that did not get confirmed successfully. Note that we have deduplicated the objects by counting file-system-specific objects and fields only once, e.g., ext2_inode and ext4_inode are deduplicated. Even though not the focus of our research, we also find extra 8 in stack objects and list them in Table 7 in the appendix.

The objects dynamically verified to be exploitable for privilege escalation are marked with ✓. For the remaining ones, there are a few FSKO fields marked with ✗ that looked like promising fields (based on their names and usage patterns in source code). However, when we verify them dynamically, overwriting their values did not lead to privilege escalation. We also listed a few cases where we do not have the right test cases to reach them dynamically (—). We

Table 4: Identified FSKOs. The attack succeeds (✓) as long as corrupting them with the targeted values.

Data	Category	Target value	Memory Cache
ext4_inode_info ->vfs_inode->i_uid->val	metadata (A)	user id	ext4_inode_cachep ✓
posix_acl ->a_entries->e_perm	metadata (A)	MAY_WRITE	variable size ✓
ext4_inode_info ->vfs_inode->i_mode	metadata (A)	S_ISUID, S_IWOTH	ext4_inode_cachep ✓
ext4_inode->i_uid	metadata	user id	page ✓
ext4_inode->i_mode	metadata	S_ISUID, S_IWOTH	page ✓
file->f_mode	metadata	FMODE_WRITE	filp_cachep ✓
iattr->ia_mode	metadata	S_ISUID, S_IWOTH	✗
vm_area_struct ->vm_flags	metadata	VM_WRITE	vm_area_cachep ✓
nfs_pgio_header->rw_mode	metadata	FMODE_WRITE	✗
nfs_fattr->mode	metadata	—	—
nfs_open_context->mode	metadata	FMODE_WRITE	✗
nfs_openargs->mode	metadata	FMODE_WRITE	✗
pnfs_layout_range ->iomode	metadata	IOMODE_RW	—
xfs_buf->b_flags	operation	XBF_WRITE	✗
iomap_dio->flags	operation	IOMAP_DIO_WRITE	kmalloc-96 ✓
dio->op	operation	REQ_OP_WRITE	dio_cache ✓
request->cmd_flags	operation	REQ_OP_WRITE	kmalloc; driver-specific ✓
bio->bi_opf	operation (A)	REQ_OP_WRITE	slab cache; driver-specific ✓
scsi_cmnd->cmd[0]	operation	WRITE_32	sd_cdb_pool ✓
aio_kiocb->kiocb.flags	operation	IOCB_WRITE	—
dm_io_request->bi_op	operation	REQ_OP_WRITE	✗
pipe_buffer->flags	COW	PIPE_BUF_FLAG_CAN_MERGE	variable size ✓
fuse_copy_state->write	COW	0	—
bio->bi_iter->bi_sector	sector (A)	—	slab cache; driver-specific ✓
ext4_extent->ee_block	sector	—	✗
extent_status->es_pblk	sector	—	ext4_es_cachep ✓
extent_map->block_start	sector	—	extent_map_cachep ✓
extent_state->start	sector	—	✗
buffer_head->b_blocknr	sector	—	bh_cachep ✓
xfs_buf_map->bm_bn	sector	—	✗
request->__sector	sector	—	kmalloc; driver-specific ✓
nvme_request->result	sector	—	kmalloc; driver-specific ✓
ext4_inode_info->vfs_inode->i_mapping	page cache tree	—	FS-specific ✓
file->f_mapping	page cache tree (A)	—	filp_cachep ✓
pipe_buffer->page	page cache	—	variable size ✓
address_space->i_pages (xarray)	page cache	—	radix_tree_node_cachep ✓
bio_vec->bv_page*	page cache	—	✗

also show these unexploitable cases to demonstrate the necessity of dynamic verification because they have similar semantics to FSKOs.

Note that we exclude some redundant data from the statistics results. Specifically, the object inode is an embedded field and accessible through ext4_inode_info->vfs_inode of inode type. Therefore, the allocation of inode goes with the allocation of the parent object, i.e., ext4_inode_info. Each file system implementation has its own parent object, like ext2_inode_info. FSKO-AUTO identifies 48 parent objects which are not listed in the paper.

Now we explain the columns in Table 4. We show all FSKO fields (including the struct name) in the first “Data” column. The “Category” column lists the semantics of the FSKO fields, and (A) represents an anchor object. The “Target value” column represents the constants used to escalate privilege after corrupting an FSKO field during dynamic verification. Finally, as for the last column, the memory cache, a few objects are allocated in different caches depending on where they are allocated (drivers or network protocols), and their exact sizes vary accordingly. We label them with a suffix “-specific” at the end (e.g., “driver-specific”). In addition, some objects are allocated as arrays with variable sizes, and we denote them as “variable size”. Additionally, some objects are directly

Table 5: The bridge objects facilitating the page-UAF strategy.

Data	Memory Cache	Read/Write
address_space->i_pages (xarray)	radix_tree_node_cachep	* ★
pipe_buffer->page*	variable size	* ★
bio_vec->bv_page*	variable size	* ★
mptcp_data_frag->page*	page	★
sock->sk_frag.page*	slab cache/kmalloc; net-specific	★
wait_page_queue->page*	variable size	* ★
xfrm_state->xfrag->page*	kmalloc-1k	* ★
pipe_inode_info->tmp_page*	kmalloc-192	* ★
lbuf->l_page*	kmalloc-128	* ★
skb_shared_info->frags->bv_page*	variable size	* ★
cifs_writedata->pages**	variable size	★
cifs_readdata->pages**	variable size	*
fuse_args_pages->pages**	variable size	★
ceph_sync_write() **	variable size	★
ceph_sync_read() **	variable size	*
process_vm_rw_core() **	variable size	* ★
orangepfs_bufmap_desc->page_array**	variable size	* ★
vm_struct->pages**	variable size	*
agp_memory->pages**	variable size	★
ttmp->pages**	variable size	★
Sg_scatter_hold->pages**	variable size	*
st_buffer->reserved_pages**	variable size	* ★
configs_buffer->page	kmalloc-128	* ★

*: the pointer can read information of the page content.
 ★: the pointer can write the page content through page pointer.

allocated via the page allocator instead of the slab/slub allocator, and we use “page” in the table.

Flag objects. The beginning of Table 4 presents 13 flags that result in privilege escalation as long as they are corrupted with the targeted values. Again, we deduplicated FSKOs that are file system specific and use objects under ext4 that represent others. The COW in Table 3 and Table 4 means a special type of flag that controls Copy-On-Write (COW). These flags are in conditional branches controlling whether to obtain the existing page cache or allocate a new page. It is worth mentioning that FSKO-AUTO also identifies other COW objects for disk operations (instead of memory). Unfortunately, we find that a normal user is not permitted to trigger a disk COW of a read-only file; thus, such objects are not included in our results.

Reference objects. We list 10 objects for sector numbers and 4 objects representing the page-cache tree in Table 4. Substituting them with the reference for a read-only file enables unprivileged users to write it. We also use ext4 as an example, ignoring similar objects in other file systems.

Bridge Objects in page-UAF Strategy. Table 5 show 23 exploitable objects facilitating the page-UAF strategy. All these objects can be corrupted to point to other pages, thus creating page UAFs when they are freed. Then, we observe that 16 objects can be used to read (*), and 20 write (★) the page content through their inside pointer. 12 can achieve both reading and writing. In the first column in Table 5, * signifies the field being of type page *. As for **, it means the field is of type page **, i.e., a heap-allocated array of page *. These arrays are often variable in size, and their elements can be corrupted to construct page UAF. There is one entry (the first) marked with xarray, which is a specific data structure that is similar to page **. Finally, some entries in Table 5 are listed

with function names ending with () and ** (in the first column). They represent cases where an array of page pointers are allocated on the heap with kmalloc_array() without being assigned to a field of any heap object (instead if it is assigned to a stack variable). Nevertheless, the page pointer array can still be corrupted.

Comparison with KENALI [62]. We run the open-source code in GitHub [35], which uses generic heuristics to identify as many non-control data as possible across the Linux kernel, regardless of their semantics and exploitability. The idea is to look at various that can affect the control flow of the program, where one path leads to an error code in return and the other does not. After running it on the same kernel version v5.14, we found 189 struct types, and only 10 of them belong to the file subsystem. 2 of them structs were also found, namely inode and address_space. We investigated the remaining 8, and none of them can lead to privilege escalation using the three types of semantics we considered.

5.2 Exploitability Evaluation

Out of the 26 CVEs mentioned earlier, we confirmed that 18 of them are exploitable, by attempting to pair FSKOs with the vulnerability in the two dimensions: (1) write capabilities, and (2) slab cache requirements. The 10 CVEs are not suitable for FSKOs because they occur in specific subsystems that are not general to heap objects, e.g., eBPF subsystem.

The results are shown in Table 6. We list the exploit strategies suitable for each CVE. To further demonstrate the exploitability, we develop 10 end-to-end exploits where the underlying Linux kernel enables all the exploit mitigation mechanisms: CVE-2021-22600 (Double Free), CVE-2021-22555 (OOB), CVE-2022-0995 (OOB), and CVE-2022-0185 (OOB), CVE-2023-5345 (UAF). These CVEs cover all three kinds of vulnerabilities: OOB, UAF, and double-free. For the remaining CVEs, we have determined that they can also be exploited similarly based on their write capability and slab cache requirement. However, writing them is time-consuming because of tasks such as creating stable heap fengshui.

Corrupting flag objects. Corrupting flag objects is straightforward as long as the write capability is suitable for the specific flag field. For this exercise, we write an exploit using CVE-2022-0995 (an OOB vulnerability) to corrupt the pipe_buffer->flags. This vulnerability happens to have the right capability to write this flag. This CVE provides a capability that sets an arbitrary bit to 1 based on the vulnerable object of size 0x38. After placing the FSKO adjacent to the vulnerable object, we can directly corrupt the flag.

Corrupting reference objects. We write an exploit using CVE-2021-22600 to corrupt the reference objects with the target value after infoleak. Specifically, this double-free vulnerability enables an attack to read the value of a page cache pointer to a read-only file (by spraying pipe_buffer objects). Subsequently, we trigger another double free to replace a read-write file’s page cache pointer (again in pipe_buffer) to have it point to the read-only page cache. Note that this exploit does require infoleak. Fortunately, we can construct page UAF to avoid the need for infoleak for this CVE, as discussed next.

Corrupting FSKOs via page-level UAF. We developed 8 exploits against CVE-2023-5345 (one exploit), CVE-2022-0995 (two

Table 6: Exploitability demonstrated on 18 real-world vulnerabilities using different exploit strategies. ○ means the exploits can not bypass CFI or data protection; ◐ means the exploits can bypass CFI but fail if data protection is enforced; ● means the exploits can bypass both CFI and data protection.

CVE ID	CVE Type	Strategy	Existing Exploit		Exploit FSKO	
			Bypass CFI & Data Protection	No Need InfoLeak	Bypass CFI & Data Protection	No Need InfoLeak
CVE-2023-5345	Use-After-Free	Page UAF	◐	✗	●	✓
CVE-2022-0995	Out-Of-Bounds	Corrupt Flag, Page UAF	○	✗	●	✓
CVE-2022-0185	Out-Of-Bounds	Page UAF	●	✗	●	✓
CVE-2021-22555	Out-Of-Bounds	Page UAF	●	✗	●	✓
CVE-2021-22600	Double-Free	Corrupt Flag, Page UAF	●	✗	●	✓
CVE-2022-27666	Out-Of-Bounds	Page UAF	◐	✗	●	✗
CVE-2022-25636	Out-Of-Bounds	Corrupt Flag, Page UAF	●	✗	●	✓
CVE-2022-2639	Out-Of-Bounds	Page UAF	●	✗	●	✓
CVE-2022-2588	Use-After-Free	Corrupt Flag, Page UAF	◐	✓	●	✓
CVE-2021-3492	Double Free	Corrupt Flag, Page UAF	○	✗	●	✓
CVE-2021-43267	Out-Of-Bounds	Page UAF	○	✗	●	✓
CVE-2021-41073	Use-After-Free	Corrupt Flag, Page UAF	○	✗	●	✓
CVE-2021-4154	Use-After-Free	Corrupt Flag, Page UAF	◐	✓	●	✓
CVE-2021-42008	Out-Of-Bounds	Page UAF	●	✗	●	✓
CVE-2021-27365	Out-Of-Bounds	Page UAF	○	✗	●	✓
CVE-2021-26708	Use-After-Free	Page UAF	○	✗	●	✗
CVE-2021-23134	Use-After-Free	Corrupt Flag, Page UAF	◐	✗	●	✓
CVE-2020-14386	Out-Of-Bounds	Page UAF	○	✗	●	✓

exploits), CVE-2022-0185 (three exploits), CVE-2021-22555 (one exploit), and CVE-2021-22600 (one exploit), to corrupt FSKO fields, `pipe_buffer->flags`, `file->f_mode`, and `file->f_mapping`. The exploits spray bridge objects to construct page UAF – we specifically target `pipe_buffer->page` and `configfs_buffer->page`. These exploits are generally easier and more stable because there is no need for infoleak, cross-cache attacks, and page-level fengshui.

Note that there are two exceptions: CVE-2022-27666 and CVE-2021-26708, where we still need to construct an infoleak together with the page-UAF strategy. This is due to their limited out-of-bound write capabilities. Specifically, after performing an out-of-bound write, the vulnerability of CVE-2022-27666 adds uncontrollable junk data at the end, which makes it difficult to corrupt only the lower bits of the page pointer. The vulnerability capability of CVE-2021-26708 allows us to write four bytes at a fixed 0x40 offset out-of-bounds, which do not match the offset of the lower bits of page pointers in our bridge objects. Therefore, we have to construct an infoleak to disclose a valid page pointer and then completely replace the original page pointer.

Complexity Reduction. Given that our exploits, which use the page-UAF strategy and directly overwrite flag, do not require infoleak or KASLR bypassing, we are interested in measuring the exploit complexity by the number of lines of code (LoC). For the exploits to corrupt flags, our exploits achieve 20% and 70% reduction in exploit LoC compared to the original publicly available exploits. The case with 70% reduction is where the capability of the CVE allows a simple corruption of the flag at the correct offset (without requiring the page-level UAF first). Even with the page-level UAF,

we still managed to see a 20% reduction in LoC. For the other exploits, we observe approximately a 30% reduction in LoC compared to the original exploits. Besides, some of the original exploits are based on return-oriented programming (ROP) and require an extra offline step of gadget extraction. This step is not necessary in our exploits. In summary, the complexity reduction is mainly due to the absence of infoleak requirements and ROP (collect gadgets and get their offsets).

Stability and Usability. The page-level-UAF-based exploits have a high success rate, around 90% - 100%. We ran our exploit of five exploits using CVE-2022-0995 and CVE-2022-0185. Each exploit is run 10 times, and 9 or 10 out of 10 are successful without any crashes. This is attributed to the fact that the exploits have a built-in feedback mechanism (i.e., read of the physical page) that helps make sure the final write is performed on the right target. Using the feedback mechanism, we can restart the page-level UAF if the expected target is not detected.

Future-proof. We compare our exploits with the publicly available exploits in Table 6. Note that a CVE can have multiple public exploits, and we select the most powerful one that can bypass as many protections as possible. Some exploits are rendered ineffective under the deployment of CFI, and we use ○ to indicate the exploits that cannot bypass CFI. Besides CFI, it is a trend that some critical non-control data will be protected, e.g., `struct cred`, `modprobe_path`, and page table are already protected in Android kernels [48, 61]. This will make a number of exploits obsolete in the near future. We mark existing exploits with ◐ if they choose to corrupt the non-control data that are protected today in Android

kernels (because they are likely going to be protected in Linux kernels in the future). We use ● to represent the exploits that bypass CFI and data protection.

Notably, we have verified that our exploits are able to succeed, even when the kernel enables the protection against DirtyCred, i.e., providing isolation of file objects by separating them into high-privileged and low-privileged ones. The reason this defense is ineffective against our exploit strategy is that we do not rely on co-locating high-privileged and low-privileged objects. Instead, once a page is freed, we can spray all high-privileged objects, e.g., spraying `\etc\passwd` files, and modify the permission `file->f_mode`.

To the best of our knowledge, only five open-source existing exploits corrupt unprotected non-control data. But they all require info leak ✗; our exploits do not ✓. Four of them corrupt the flag `pipe_buffer->flags`, which is found in our paper. One (CVE-2021-22600) uses the User-Space-Mapping-Attack (USMA) technique that requires the target code address through infoleak. However, it typically leaks a specific address, where attackers calculate the target address by adding an offset based on the leaked address. The kernel code may differ version by version, resulting in the offset not being the same and the need to adapt for different kernels.

6 DISCUSSION

Ethical consideration. This work does not find any new vulnerabilities but exploits only existing ones with newly identified data. All exploited CVEs are patched already.

Defense. The newly identified objects in this work can enhance existing protection mechanisms, fostering further research. Many protection mechanisms can safeguard critical non-control objects, such as data integrity using shadow memory [62] and supervisor-based protection [48]. Many researchers propose various techniques to prevent specific vulnerabilities from being exploited, e.g., eBPF-based monitor [65], slab allocator redesign [24, 59, 1, 66, 54].

To mitigate the page UAF strategy that we proposed, we first observe that it works because the kernel allows mixes of two types of pages: (1) those freely accessible to user-space, e.g., `pipe_buffer` pages where arbitrary read and write can be performed from the user-space; (2) those that are used as slab caches. Our proposal is simple: we can isolate these two types of page groups so that they will never be overlapping in a page UAF situation. In fact, the Linux kernel already has mechanisms of page isolation [46], where between 4 and 6 page groups are enabled (depending on the kernel config). We can simply add another page group to address this concern.

Generality. Our methodology is general to other operating systems because the classification is based on three essential abstract semantics of files. The elements are general not only to the file systems of Linux but also to the ones of others. For example, Android inherited much of the Linux kernel file subsystem. FreeBSD [6] has an object `vnode` in VFS similar to `inode` in Linux VFS. They represent and manage a file on disk, indicating the generality. We leave the exploration of other OS file subsystems as future work.

Other objects in the file subsystem. An interesting question is how many more undiscovered objects are exploitable in the file subsystem. We hypothesize there must be some for two reasons.

First, we have left out a large number of wrapper objects due to potential challenges in overcoming the side effects of swapping complex objects embodying multiple semantics. By our estimate, there are about 500 wrapper objects that have pointers to one of the FSKOs we identified. Second, we have considered only three types of semantics relevant to privilege escalation. There could be other types of semantics such as namespace, which potentially can have implications on container escape and privilege escalation. We leave the exploration of these additional objects as future work.

Other OS subsystems. In this paper, we have performed a systematic analysis of objects in the file subsystem suitable for non-control attacks. We expect the experience can serve as a guideline for researching other subsystems. For example, memory management is a complex subsystem that likely will have many exploitable objects. So far, the page table has been the primary example. Another example is the driver subsystem where many specialized drivers exist. They are capable of managing their own memory and issuing direct memory read/write from the devices (e.g., GPU).

7 RELATED WORK

Kernel Vulnerability Exploitation. Control flow hijacking is a powerful exploit strategy [7]. FUZE [70] and KOUBE [9] utilize techniques like kernel fuzzing, symbolic execution, and heap manipulation for automatic exploit generation, targeting use-after-free and out-of-bound vulnerabilities, respectively. These exploits corrupt control data to achieve control flow hijacking. KEPLER [69] leverages kernel-user interactions to convert control-flow hijacking into stack overflow and ROP attacks. GREBE [39] and SyzScope [80] investigate fuzzer-exposed bugs and reveal the transformation of low-risk bugs into high-risk ones, such as control flow hijacking for kernel exploitation. The prevalence of control flow hijacking and code reuse attacks [76, 10] in real-world incidents has spurred the development of various defense mechanisms [28, 73, 20, 75, 41, 49], making control-data related exploitation more challenging.

Non-Control Data Attacks. Turing-complete DOP [32] shows code execution can be achieved by chaining data gadgets without altering the control flow. However, it needs an arbitrary write primitive, which is not always available in real-world vulnerabilities. VIPER [74] identifies the syscall-guard variables in user space programs that determine to invoke security-related system calls. However, existing works on kernel space data-only attacks merely concentrate on a limited set of known objects. Ret2page [55] corrupts the page table entry to launch attacks. DirtyCred [40] offers a new exploit strategy to grant high-privilege credentials to normal users. DirtyCOW [22] exploits a race condition bug that exists in Copy-on-Write. DirtyPipe [4, 15] leverages the use-before-initialization flag of a pipe to write to the page cache, bypassing permission checks. Previous work [8] shows that user identity, configuration, input, and decision-making data can be exploited for attacks. The variables controlling the Linux auditing framework, AppArmor, and NULL pointer dereference mitigation can be bypassed through data attacks [72]. Our work aims to discover more critical non-control data and expose unknown exploitable objects in the Linux kernel.

Protection of Critical Data. As for control data, control-flow integrity mechanisms are available in the mainline Linux kernel [13] and enforced in Android by default [2], making hijacking control flow increasingly difficult. Moreover, commercial kernels use advanced techniques to safeguard the commonly exploited non-control data. These advanced defenses can protect specific data from corruption even in the face of an attack achieving arbitrary kernel-memory-write primitive. Specifically, Android uses the hypervisor to monitor `cred` and page table via real-time protection [48, 61]. Meanwhile, the Page Protection Layer (PPL)[37], an Apple-introduced feature tailored to protect certain parts of the kernel from itself, is also utilized to protect critical data, including `cred` and page tables [60, 3]. The global variable `modprobe_path` [38] can be marked as read-only when enabling `config` introduced in Linux kernel 4.11 (`CONFIG_STATIC_USERMODEHELPER`), and is not available by default in Android kernels. In research studies, other ideas have been proposed, e.g., software-based shadow memory [35] and hardware-based Extended Page Tables (EPT) [53] that target a wide range of non-control objects. Unfortunately, their effectiveness highly depends on their ability to identify meaningful non-control data. For instance, since KENALI [35] failed to identify the majority of FSKOs that we reported, our objects would not be put into the protection domain.

8 CONCLUSION

This paper offers a comprehensive investigation into the exploitability of non-control data in the Linux file system. We systematically summarize three types of FSKOs for privilege escalation with the help of automated analyses. We analyze their exploitability by understanding their requirements of pairing with vulnerabilities of different capabilities. Along the way, we develop a novel strategy that can indirectly read and write FSKOs with high reliability. Finally, using the discovered FSKOs, we develop end-to-end exploits using 18 recent real-world CVEs.

REFERENCES

- [1] Sam Ainsworth and Timothy M. Jones. 2020. Markus: drop-in use-after-free prevention for low-level languages. In *2020 IEEE Symposium on Security and Privacy (SP)*, 578–591. doi: 10.1109/SP40000.2020.00058.
- [2] Android. [n. d.] Control flow integrity. <https://source.android.com/docs/security/test/cfi>. ().
- [3] [n. d.] Apple-oss-distributions/xnu. <https://github.com/apple-oss-distributions/xnu>. (Accessed on 09/18/2023). ().
- [4] Aaron Esau. [n. d.] Arinerron/cve-2022-0847-dirtypipe-exploit: a root exploit for cve-2022-0847 (dirty pipe). <https://github.com/Arinerron/CVE-2022-0847-DirtyPipe-Exploit>. (Accessed on 12/05/2023). ().
- [5] [n. d.] Arttnba3/d3ctf2023_d3kcache: attachment and write up for d^3ctf 2023's pwn challenge - d3kcache. https://github.com/arttnba3/D3CTF2023_d3kcache. (Accessed on 09/30/2023). ().
- [6] [n. d.] Basicvfsconcepts - freebsd wiki. <https://wiki.freebsd.org/BasicVfsConcepts>. (Accessed on 12/11/2023). ().
- [7] Nicholas Carlini, Antonio Barresi, Mathias Payer, David Wagner, and Thomas R Gross. 2015. {Control-flow} bending: on the effectiveness of {control-flow} integrity. In *24th USENIX Security Symposium (USENIX Security 15)*, 161–176.
- [8] Shuo Chen, Jun Xu, Emre Can Sezer, Prachi Gauriar, and Ravishankar K Iyer. 2005. Non-control-data attacks are realistic threats. In *USENIX security symposium*. Vol. 5, 146.
- [9] Weiteng Chen, Xiaochen Zou, Guoren Li, and Zhiyun Qian. 2020. {Koobe}: towards facilitating exploit generation of kernel {out-of-bounds} write vulnerabilities. In *29th USENIX Security Symposium (USENIX Security 20)*, 1093–1110.
- [10] Yueqi Chen, Zhenpeng Lin, and Xinyu Xing. 2020. A systematic study of elastic objects in kernel exploitation. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 1165–1184.
- [11] Yueqi Chen and Xinyu Xing. 2019. Slake: facilitating slab manipulation for exploiting vulnerabilities in the linux kernel. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 1707–1722.
- [12] Haehyun Cho, Jinbum Park, Joonwon Kang, Tiffany Bao, Ruoyu Wang, Yan Shoshitaishvili, Adam Doupe, and Gail-Joon Ahn. 2020. Exploiting uses of uninitialized stack variables in linux kernels to leak kernel pointers. In *14th USENIX Workshop on Offensive Technologies (WOOT 20)*.
- [13] [n. d.] Clang cfi support upstreamed for linux 5.13 - but only on arm64 for now - phoronix. <https://www.phoronix.com/news/Clang-CFI-Linux-5.13>. (Accessed on 09/28/2022). ().
- [14] [n. d.] Cve - cve-2022-0435. <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2022-0435>. (Accessed on 08/03/2023). ().
- [15] [n. d.] Cve - cve-2022-0847. <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2022-0847>. (Accessed on 09/28/2022). ().
- [16] [n. d.] Cve - cve-2022-1015. <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2022-1015>. (Accessed on 08/03/2023). ().
- [17] [n. d.] Cve - cve-2022-1016. <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2022-1016>. (Accessed on 08/03/2023). ().
- [18] [n. d.] Cve - cve-2022-4378. <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2022-4378>. (Accessed on 08/03/2023). ().
- [19] [n. d.] Cve-2021-26708 exploit. <https://github.com/azpema/CVE-2021-26708/blob/master/vuln.c>. (Accessed on 01/25/2024). ().
- [20] Rémi Denis-Courmont, Hans Liljestrand, Carlos Chinea, and Jan-Erik Ekberg. 2020. Camouflage: hardware-assisted cfi for the arm linux kernel. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 1–6.
- [21] DEVIL. 2021. [CVE-2021-42008] Exploiting A 16-Year-Old Vulnerability In The Linux 6pack Driver. <https://syst3mfailure.io/sixpack-slab-out-of-bounds>. (2021).
- [22] [n. d.] Dirty cow (cve-2016-5195). <https://dirtycow.ninja/>. (Accessed on 09/28/2022). ().
- [23] 2023. Dirty Pagetable: A Novel Exploitation Technique To Rule Linux Kernel. https://yanglingxi1993.github.io/dirty_pagetable/dirty_pagetable.html.
- [24] Márton Erdős, Sam Ainsworth, and Timothy M. Jones. 2022. Minesweeper: a “clean sweep” for drop-in use-after-free prevention. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '22)*. Association for Computing Machinery, Lausanne, Switzerland, 212–225. isbn: 9781450392051. doi: 10.1145/3503222.3507712.
- [25] 2021. Exploiting CVE-2021-43267. <https://haxx.in/posts/pwning-tipc/>. (2021).
- [26] [n. d.] Fiemap, an extent mapping ioctl. <https://lwn.net/Articles/297696/>. (Accessed on 01/20/2024). ().
- [27] [n. d.] Fuse for linux exploitation 101. <https://exploiter.dev/blog/2022/FUSE-exploit.html>. ().
- [28] Xinyang Ge, Nirupama Talele, Mathias Payer, and Trent Jaeger. 2016. Fine-grained control-flow integrity for kernel software. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 179–194.
- [29] Hackthebox. 2022. CVE-2022-0185 Writeup. https://www.hackthebox.com/blog/CVE-2022-0185_A_case_study. (2022).
- [30] Cedric Halbronn. 2022. SETTLERS OF NETLINK: Exploiting a limited UAF in nf_tables (CVE-2022-32250). https://research.nccgroup.com/2022/09/01/settlers-of-netlink-exploiting-a-limited-uaf-in-nf_tables-cve-2022-32250. (2022).
- [31] [n. d.] How autoslab changes the memory unsafety game. https://grsecurity.net/how_autoslab_changes_the_memory_unsafety_game. (Accessed on 01/20/2024). ().
- [32] Hong Hu, Shweta Shinde, Sendriou Adrian, Zheng Leong Chua, Prateek Saxena, and Zhenkai Liang. 2016. Data-oriented programming: on the expressiveness of non-control data attacks. In *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, 969–986.
- [33] Kyriakos K. Ispoglou, Bader Albassam, Trent Jaeger, and Mathias Payer. 2018. Block oriented programming: automating data-only attacks. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security (CCS '18)*.
- [34] [n. d.] Kernel address space layout randomization [lwn.net]. <https://lwn.net/Articles/569635/>. (Accessed on 07/13/2023). ().
- [35] [n. d.] Sslab-gatech/kernel-analyzer. <https://github.com/sslslab-gatech/kernel-analyzer>. (Accessed on 01/25/2024). ().
- [36] [n. d.] Kernel exploitation techniques: modprobe_path. https://sam4k.com/like-techniques-modprobe_path/. (Accessed on 01/08/2024). ().
- [37] Ivan Krstic. 2019. Behind the Scenes of iOS and Mac Security. In *(Blackhat USA)*. <https://i.blackhat.com/USA-19/Thursday/us-19-Krstic-Behind-The-Scenes-Of-iOS-And-Mac-Security.pdf>.
- [38] Guoren Li, Hang Zhang, Jimmeng Zhou, Wenbo Shen, Yulei Sui, and Zhiyun Qian. 2023. A hybrid alias analysis and its application to global variable protection in the linux kernel. In *Proceedings of the 32nd USENIX Security Symposium*.
- [39] Zhenpeng Lin, Yueqi Chen, Yuhang Wu, Dongliang Mu, Chensheng Yu, Xinyu Xing, and Kang Li. 2022. Grebe: unveiling exploitation potential for linux kernel bugs. In *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2078–2095.

- [40] Zhenpeng Lin, Yuhang Wu, and Xinyu Xing. 2022. Dirtycred: escalating privilege in linux kernel. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*.
- [41] [n. d.] Kernel support for control-flow enforcement [lwn.net]. <https://lwn.net/Articles/758245/>. (Accessed on 09/27/2022). ()
- [42] [n. d.] The status of kernel hardening [lwn.net]. <https://lwn.net/Articles/705262/>. (Accessed on 07/13/2023). ()
- [43] [n. d.] Kernel.org/doc/documentation/security/self-protection.txt. <https://www.kernel.org/doc/Documentation/security/self-protection.txt>. (Accessed on 07/13/2023). ()
- [44] [n. d.] Linuxflaw/cve-2017-7184 at master · vulnreproduction/linuxflaw. <https://github.com/VulnReproduction/LinuxFlaw/tree/master/CVE-2017-7184>. (Accessed on 01/07/2024). ()
- [45] Kangjie Lu. 2023. Practical program modularization with type-based dependence analysis. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 1610–1624.
- [46] [n. d.] Migrate type - linux source code (v5.14.21) - bootlin. https://elixir.bootlin.com/linux/v5.14.21/source/mm/page_alloc.c. (Accessed on 01/25/2024). ()
- [47] Arthur Mongodin. 2022. [CVE-2022-34918] A crack in the Linux firewall. <https://randorisec.fr/crack-linux-firewall>. (2022).
- [48] Peng Ning. 2014. Samsung Knox and enterprise mobile security. In *Proceedings of the 4th ACM Workshop on Security and Privacy in Smartphones & Mobile Devices*, 1–1.
- [49] [n. d.] X86: add support for clang cfi [lwn.net]. <https://lwn.net/Articles/869267/>. (Accessed on 09/27/2022). ()
- [50] [n. d.] The current state of kernel page-table isolation [lwn.net]. <https://lwn.net/Articles/741878/>. (Accessed on 07/13/2023). ()
- [51] [n. d.] Supervisor mode access prevention [lwn.net]. <https://lwn.net/Articles/517475/>. (Accessed on 07/13/2023). ()
- [52] [n. d.] Privilege escalation via setuid. <https://antonyt.com/blog/2020-03-22/privilege-escalation-via-setuid>. (Accessed on 12/05/2023). ()
- [53] Sergej Proskurin, Marius Momeu, Seyedhamed Ghavamnia, Vasileios P Kemerlis, and Michalis Polychronakis. 2020. Xmp: selective memory protection for kernel and user space. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 563–577.
- [54] [n. d.] Randomized slab caches for kcalloc() [lwn.net]. <https://lwn.net/Articles/938246/>. (Accessed on 12/19/2023). ()
- [55] [n. d.] Us-22-wang-ret2page-the-art-of-exploiting-use-after-free-vulnerabilities-in-the-dedicated-cache.pdf. <https://i.blackhat.com/USA-22/Thursday/US-22-WANG-Ret2page-The-Art-of-Exploiting-Use-After-Free-Vulnerabilities-in-the-Dedicated-Cache.pdf>. (Accessed on 07/20/2023). ()
- [56] sam4k. 2022. CVE-2023-29383: Abusing Linux chfn to Misrepresent /etc/passwd. <https://www.trustwave.com/en-us/resources/blogs/spiderlabs-blog/cve-2023-29383-abusing-linux-chfn-to-misrepresent-etc-passwd/>. (2022).
- [57] sam4k. 2022. Kernel Exploitation Techniques: modprobe_path. https://sam4k.com/like-techniques-modprobe_path. (2022).
- [58] Di Shen. 2017. Defeating samsung Knox with zero privilege. *BlackHat USA*, 13–14.
- [59] Zekun Shen and Brendan Dolan-Gavitt. 2020. Heapexpo: pinpointing promoted pointers to prevent use-after-free vulnerabilities. In *Proceedings of the 36th Annual Computer Security Applications Conference (ACSAC '20)*, Association for Computing Machinery, <conf-loc>, <city>Austin</city>, <country>USA</country>, </conf-loc>, 454–465. ISBN: 9781450388580. DOI: 10.1145/3427228.3427645.
- [60] Siguza. 2019. Evolution of iOS mitigations. In *(TyphoonCon)*. <https://raw.githubusercontent.com/ssd-secure-disclosure/typhooncon2019/master/Siguza%20-%20Mitigations.pdf>.
- [61] Samsung Enterprise Mobility Solutions. 2015. White paper: an overview of the samsung Knox platform. (2015).
- [62] Chengyu Song, Byoungyoung Lee, Kangjie Lu, William Harris, Taesoo Kim, and Wenke Lee. 2016. Enforcing kernel security invariants with data flow integrity. In *NDSS*.
- [63] S.T.A².R.S Team. 2023. CVE-2023-32233: Privilege escalation in Linux Kernel due to a Netfilter nf_tables vulnerability. <https://www.tarlogic.com/blog/cve-2023-32233-vulnerability>. (2023).
- [64] TurtleARM. 2023. CVE-2023-0179-PoC. <https://github.com/TurtleARM/CVE-2023-0179-PoC>. (2023).
- [65] Zicheng Wang, Yueqi Chen, and Qingkai Zeng. 2023. PET: prevent discovered errors from being triggered in the linux kernel. In *32nd USENIX Security Symposium (USENIX Security 23)*. USENIX Association, Anaheim, CA, (Aug, 2023), 4193–4210. ISBN: 978-1-939133-37-3. <https://www.usenix.org/conference/usenixsecurity23/presentation/wang-zicheng>.
- [66] Brian Wickman, Hong Hu, Insu Yun, DaeHee Jang, JungWon Lim, Sanidhya Kashyap, and Taesoo Kim. 2021. Preventing Use-After-Free attacks with fast forward allocation. In *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, (Aug, 2021), 2453–2470. ISBN: 978-1-939133-24-3. <https://www.usenix.org/conference/usenixsecurity21/presentation/wickman>.
- [67] [n. d.] Will’s root: reviving exploits against cred structs - six byte cross cache overflow to leakless data-oriented kernel pwnage. <https://www.willsroot.io/2022/08/reviving-exploits-against-cred-struct.html>. (Accessed on 12/07/2023). ()
- [68] [n. d.] Will’s root: reviving exploits against cred structs - six byte cross cache overflow to leakless data-oriented kernel pwnage. <https://www.willsroot.io/2022/08/reviving-exploits-against-cred-struct.html>. (Accessed on 01/07/2024). ()
- [69] Wei Wu, Yueqi Chen, Xinyu Xing, and Wei Zou. 2019. {Kepler}: facilitating control-flow hijacking primitive evaluation for linux kernel vulnerabilities. In *28th USENIX Security Symposium (USENIX Security 19)*, 1187–1204.
- [70] Wei Wu, Yueqi Chen, Jun Xu, Xinyu Xing, Xiaorui Gong, and Wei Zou. 2018. {Fuze}: towards facilitating exploit generation for kernel {use-after-free} vulnerabilities. In *27th USENIX Security Symposium (USENIX Security 18)*, 781–797.
- [71] [n. d.] Extending the use of ro and nx [lwn.net]. <https://lwn.net/Articles/422487/>. (Accessed on 07/13/2023). ()
- [72] Jidong Xiao, Hai Huang, and Haining Wang. 2015. Kernel data attack is a realistic security threat. In *International Conference on Security and Privacy in Communication Systems*. Springer, 135–154.
- [73] Yutian Yang, Songbo Zhu, Wenbo Shen, Yajin Zhou, Jiadong Sun, and Kui Ren. 2019. Arm pointer authentication based forward-edge and backward-edge control flow integrity for kernels. *arXiv preprint arXiv:1912.10666*.
- [74] Hengkai Ye, Song Liu, Zhechang Zhang, and Hong Hu. 2023. VIPER: spotting Syscall-Guard variables for Data-Only attacks. In *32nd USENIX Security Symposium (USENIX Security 23)*. USENIX Association, Anaheim, CA, (Aug, 2023), 1397–1414. ISBN: 978-1-939133-37-3. <https://www.usenix.org/conference/usenixsecurity23/presentation/ye>.
- [75] Sungbae Yoo, Jinbum Park, Seolheui Kim, Yeji Kim, and Taesoo Kim. 2022. {Inkernel} {control-flow} integrity on commodity {oses} using {arm} pointer authentication. In *31st USENIX Security Symposium (USENIX Security 22)*, 89–106.
- [76] Kyle Zeng, Yueqi Chen, Haehyun Cho, Xinyu Xing, Adam Doupe, Yan Shoshitaishvili, and Tiffany Bao. 2022. Playing for {k(h) eaps}: understanding and improving linux kernel exploit reliability. In *31st USENIX Security Symposium (USENIX Security 22)*, 71–88.
- [77] Hang Zhang, Weiteng Chen, Yu Hao, Guoren Li, Yizhuo Zhai, Xiaochen Zou, and Zhiyun Qian. 2021. Statically discovering high-order taint style vulnerabilities in os kernels. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 811–824.
- [78] Jimmeng Zhou, Tong Zhang, Wenbo Shen, Dongyoon Lee, Changhee Jung, Ahmed Azab, Ruowen Wang, Peng Ning, and Kui Ren. 2022. Automatic permission check analysis for linux kernel. *IEEE Transactions on Dependable and Secure Computing*.
- [79] Xiaochen Zou. 2022. CVE-2022-27666 Writeup. <https://etenal.me/archives/1825>. (2022).
- [80] Xiaochen Zou, Guoren Li, Weiteng Chen, Hang Zhang, and Zhiyun Qian. 2022. {Syzscope}: revealing {high-risk} security impacts of {fuzzer-exposed} bugs in linux kernel. In *31st USENIX Security Symposium (USENIX Security 22)*, 3201–3217.

A APPENDIX

Additionally, we also verified 8 FSKO fields are in the stack shown in Table 7, which may raise security concerns for stack-based corruption [16, 17, 18, 14, 69, 12]. Besides, DirtyCOW is a logic bug that manipulates the flag `vm_fault->flags` in the stack to achieve root. We believe these objects may also be valuable in stack-based vulnerabilities and logic bugs for the future research.

Table 7: Identified and verified critical FSKOs in the stack.

Data	Category	Target value	Verified
<code>iov_iter->data_source</code>	operation	WRITE	✓
<code>iomap_iter->flags</code>	operation	IOMAP_WRITE	✓
<code>blk_mq_alloc_data->cmd_flags</code>	operation	REQ_OP_WRITE	✓
<code>vm_fault->flags</code>	COW	FAULT_FLAG_WRITE	✓
<code>iomap->addr</code>	sector	sector number	✓
<code>ext4_map_blocks->m_pblk</code>	sector	sector number	✓
<code>xfs_bmbt_irec->br_startblock</code>	sector	sector number	✓
<code>erofs_map_blocks->m_pa</code>	sector	sector number	✓