SOSCON Attack and Defense on Linux kernel

SAMSUNG Research | Security Team | Jinbum Park 2018.10.18

SOSCON 2018

Contents

Full steps of attack on Linux kernel

Attack-1: Modify sensitive RW data

Defense-1: ro-after-init

Attack-2: Modify process credential

Defense-2: PrivWatcher

Attack-3: addr_limit bug

Defense-3: Add checking for addr_limit

Attack-4: Modify addr_limit via stack-based attack

Defense-4: Split addr_limit from stack

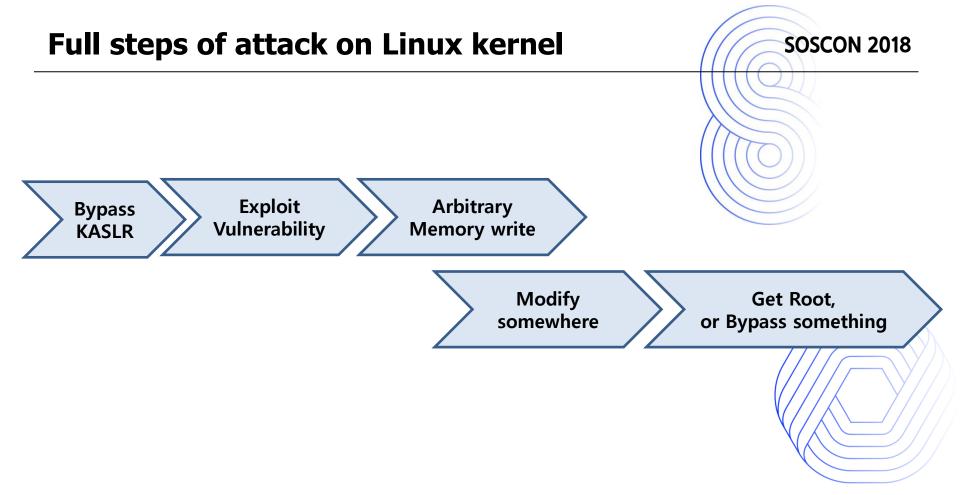
+ Bonus: Advanced attacks

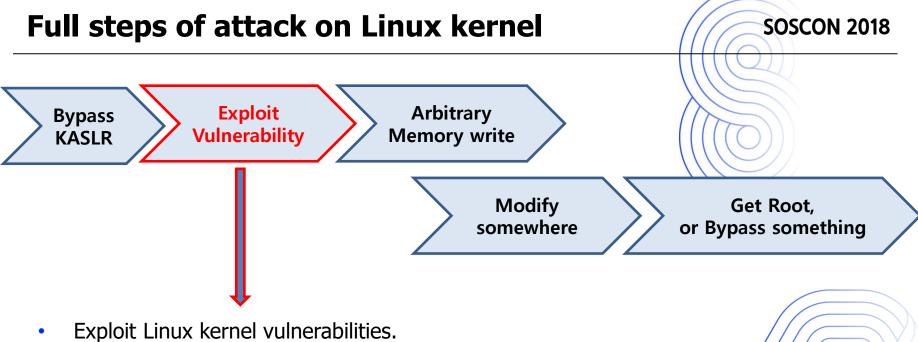
SOSCON 2018



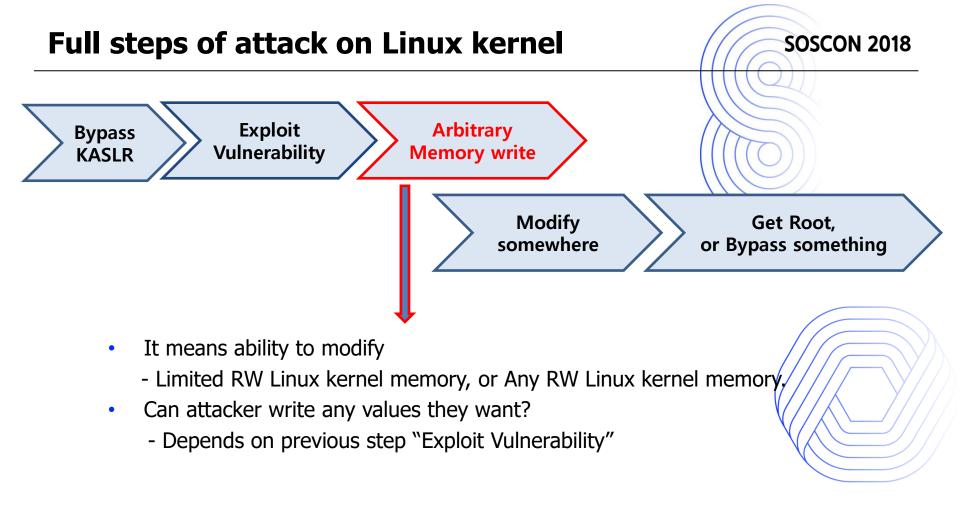
Full steps of attack on Linux kernel

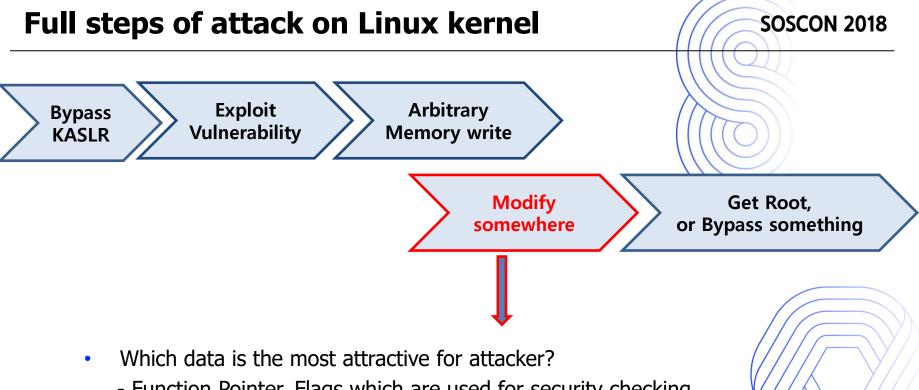
SOSCON 2018





- By exploiting them, Attacker can
 - Modify control flow,
 - Do arbitrary memory write.





- Function Pointer, Flags which are used for security checking.
- Attacker can get root by modifying function pointer.
- Attacker can bypass security mechanism by modifying some flags.

Full steps of attack on Linux kernel



Exploit Bypass Vulnerability KASLR

Arbitrary Memory write

> Modify somewhere

> > Final Goal!!



Modify sensitive RW data SOSCON 2018 Expioit Arbitrary Bypass Vulnerability KASLR Memory write Modify **Get Root,** or Bypass something somewhere

Sensitive RW data: Function Pointer

Why function pointer is critical? Let's look at Linux kernel 3.10.

If attacker can modify function pointer - .secmark_refcount_inc, What can attacker do?



Sensitive RW data: Function Pointer

```
static void selinux secmark refcount inc(void)
                                                      Normal
static struct security operations selinux ops = {
                      "selinux",
   .name =
secmark relabel packet =
                       selinux secmark relabel packet
                                                                           oid reset security ops(void)
.secmark refcount inc =
                       selinux secmark refcount inc,
                       selinux secmark refcount dec.
secmark refcount dec =
                                                                              security ops = &default security ops;
                                                        Malicious!!
       *secmark refcount inc)
                                                                             security capget (struct task struct
      (*secmark refcount dec) (void);
                                                                                       kernel cap t *effective,
                                                    Possible?
                                                                                       kernel cap t *inheritable,
                                                                                       kernel cap t *permitted)
```

Attacker can call other security-critical function which has same function type. "reset_security_ops()" disables Linux security module such as Smack, SELinux, ... So that, Attacker can bypass Linux security module!!

Sensitive RW data: Flags which are used for security checking

Let's look at Linux kernel 3.10.

```
Flag to represent whether SELinux is initialized or not.

Used for security checking!

If (!ss_initialized) {
    avtab_cache_init();
    rc = policydb_read(&policydb, fp);

If attacker can set this to 0,
    Reinitializing SELinux policy is possible!!

And other operations too!!
```

Sensitive RW data: Flags which are used for security checking

```
static struct sidtab sidtab;
struct policydb policydb;
int ss_initialized;

Flag to represent whether
SELinux is initialized or not.
```

Defeating Samsung KNOX with zero privilege, Di shen, Blackhat USA 2017







ro-after-init

Read only after initialization

- What is a key insight inside ro-after-init?
 - A lot of RW data are used to be written only one time.
 - When?? → Kernel Initialization time!!
 - Then?? → The RW data can be marked as read-only after initialization!
 - It reduces a lot of attack surface with no performance overhead!!





ro-after-init SOSCON 2018

Read only after initialization

- How to apply ro-after-init?

- Just add keyword "_ro_after_init" to variables which you want to protect.
- Limitation: Developer should know which variables can be marked as ro-after-init.

 Automatic process for marking them has not been appeared yet.

ro-after-init SOSCON 2018

Read only after initialization

- Real-world cases for protecting function pointers

```
static struct security_hook_list smack_hooks[] = {
   LSM_HOOK_INIT(ptrace_access_check, smack_ptrace_access_check),
   LSM_HOOK_INIT(ptrace_traceme, smack_ptrace_traceme),
   LSM_HOOK_INIT(syslog, smack_syslog),
```

Linux kernel 4.8

Linux kernel 4.12



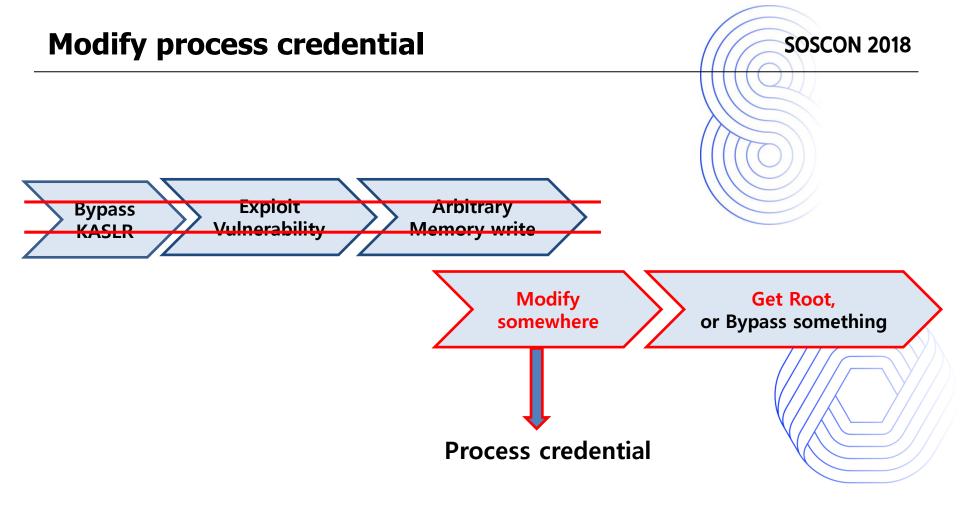


Reduce attack surface as much as possible!!





Attack-2: Modify process credential

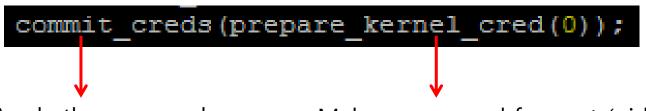


cred

```
Kernel structure to represent one process
truct task struct
  volatile long state;
  void *stack;
 process credentials */
  const struct cred
                     rcu *real cred;
                                                   Credential for this process.
                    credentials (COW)
                                                   We will modify this!
  const struct cred
                    credentials (COW)
             usage;
             subscribers;
                             /* number of processes s
  atomic t
             *put addr;
             magic;
                                                          Credential is tightly related to
                 0x43736564
lefine CRED MAGIC DEAD 0x44656144
                                                          permission of process!!
  kuid t
             uid:
                         /* real UID of the task */
                            real GID of the task
  kaid t
             gid:
```

Type1: Function calls to modify cred for root

- Attacker executes below two function calls. (kernel function)



Apply the new cred to current process

Make a new cred for root (uid=0)

truct task struct

These function calls makes attacker to get root!!

A lot of real-world attacks use this technique,

- CVE-2016-0728, ...

Type2: Reuse init_cred

```
struct task struct {
                                                         Original cred for user permission
   volatile long state;
   void *stack;
  process credentials */
                     rcu *real cred;
   const struct cred
                   * credentials (COW
                                                            init_cred for root permission
                      rcu *cred;
   const struct cred
                   * credentials (COW) */
                                                         The initial credentials for the initial task
                                                        ruct cred init cred = {
                                                                         = ATOMIC INIT(4),
                                                       ifdef CONFIG DEBUG CREDENTIALS
                                                           .subscribers
                                                                             = ATOMIC INIT(2),
                                                                         = CRED MAGIC,
                                                           .magic
                                                                          = GLOBAL ROOT UID,
                                                          .uid
                                                                          = GLOBAL ROOT GID.
```

Type3: Modify cred itself

struct task struct {

```
volatile long state;
  void *stack:
 process credentials */
                       rcu *real cred;
  const struct cred
                     credentials (COW)
                       rcu *cred;
  const struct cred
                      credentials (COW) */
ruct cred
               usage;
               subscribers;
                               /* number of processes s
  atomic t
               *put addr;
              magic;
define CRED MAGIC 0x43736564
define CRED MAGIC DEAD 0x44656144
                           /* real UID of the task */
   kuid t
               uid:
   kgid t
                              real GID of the task
               gid:
```



Modify these directly!!





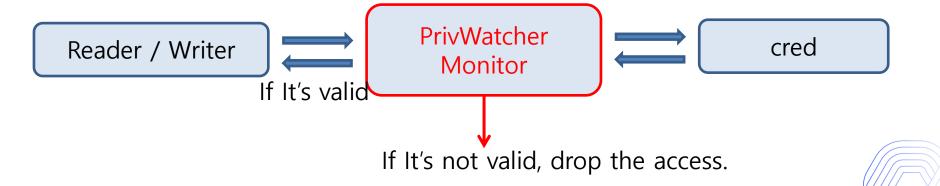


 PrivWatcher: Non-bypassable Monitoring and Protection of Process Credentials from Memory Corruption Attacks,
 AsiaCCS 2017, Samsung Research America

- → This is a paper proposed by Samsung Research America!!
- → This is not merged in Linux kernel mainline.
- → Is this merged in Linux kernel for Galaxy??



Simple principle for defense



Attack Type1: Function calls to modify cred

Attack Type2 : Reuse init_cred

Attack Type3: Manipulate cred itself

PrivWatcher can prevent all attack types!! Prevent privilege escalation through cred.

PrivWatcher SOSCON 2018

Is this merged in Linux kernel for Galaxy?

```
#ifdef CONFIG_RKP_KDF
atomic_t *use_cnt;
    struct task_struct *bp_task;
    void *bp_pgd;
    unsigned long long type;
#endif /*CONFIG_RKP_KDP*/
```

```
/* Main function to verify cred security context of a process */
int security_integrity_current(void)
{
    if ( rkp_cred_enable &&
        (rkp_is_valid_cred_sp((u64)current_cred(),(u64)current_cred()->security)||
        cmp_sec_integrity(current_cred(),current->mm)||
        cmp_ns_integrity())) {
        rkp_print_debug();
        panic("RKP_CRED_PROTECTION_VIOLATION\n");
    }
    return 0;
}
```

Not same solution to PrivWatcher.
 But, There is a similar solution in after Galaxy S7.



Galaxy Note9 Kernel Code

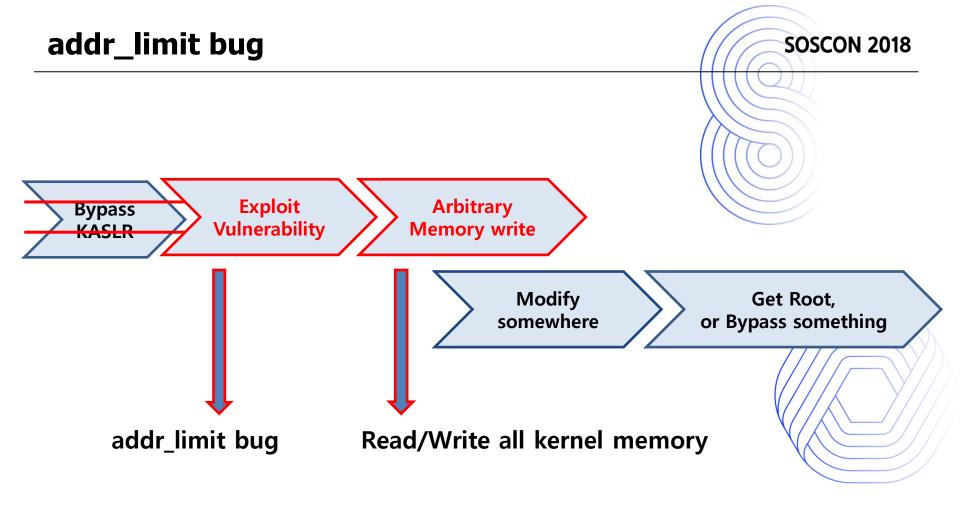


You can add your security solution into your product!









What is addr_limit?

- Look at "struct thread_info" which is generated per process.
- It's different per CPU type. Below one is for arm64.

- addr_limit have a role like partition between user and kernel space.



Normal state-flow of addr_limit

User : addr_limit == USER_DS

Can access user space only



Updated by Kernel or Kernel driver

Kernel: addr_limit == KERNEL_DS

Can access user+kernel space



Restored by Kernel or Kernel driver

User: addr_limit == USER_DS

Can access user space only

Mistaken state-flow of addr_limit (mistakes from developer)

User: addr_limit == USER_DS

Can access user space only



Kernel: addr_limit == KERNEL_DS

Can access user+kernel space



Miss restore!! (human error)

User: addr_limit == KERNEL_DS

Can access user+kernel space!! Read/Write all Kernel memory!!

Real-world vulnerability

```
int _write_log(char *filename, char *data)
  strict file *file:
if (f54_window_crack || f54_window_crack_check_mode == 0) {
 <u>mm_segment_t_old_fs</u> = get_fs();
 set_fs(KERNEL_DS);
                                         → addr limit == KERNEL DS
 flags = U_WRUNLY | O_CREAT;
if (filename) {
 file = filp_open(filename, flags, 0666);
 sys_chmod(filename, 0666);
 TOUCH_E("%s : filename is NULL, can not open FILE\n",
                                                           Not restored!!
   __f unc__);
 return -1;
```



This is one of real-world vulnerabilities, which in LG G4 touch screen driver in Android.

addr_limit bug

How can modify kernel memory actually??

memcpy(kernel_addr, buf, len);

< User >

User : addr_limit == KERNEL_DS

- Then, Can an attacker modify kernel memory like above?? (after addr_limit bug) Definitely No...



- How to modify??
 - Exploiting pipe subsystem (http://blog.daum.net/tlos6733/184)



Add checking for addr_limit

What is the most critical problem for handling addr_limit?

- Possibility of human error!!

User: addr_limit == USER_DS



Kernel : addr_limit == KERNEL_DS



Human error point!!

User: addr_limit == KERNEL_DS





Add checking for addr_limit

Solution

- Enforce security-checking when returned from Kernel to User.

User : addr_limit == USER_DS



Kernel : addr_limit == KERNEL_DS



Checking!! Reporting error!! Process will be killed!!

User: addr_limit == KERNEL_DS





SOSCON 2018

Add checking for addr_limit

Solution

```
static inline void set_fs(mm_segment_t fs)
{
    current_thread_info()->addr_limit = fs;

    /* On user-mode return, check fs is correct */
    set_thread_flag(TIF_FSCHECK);
```

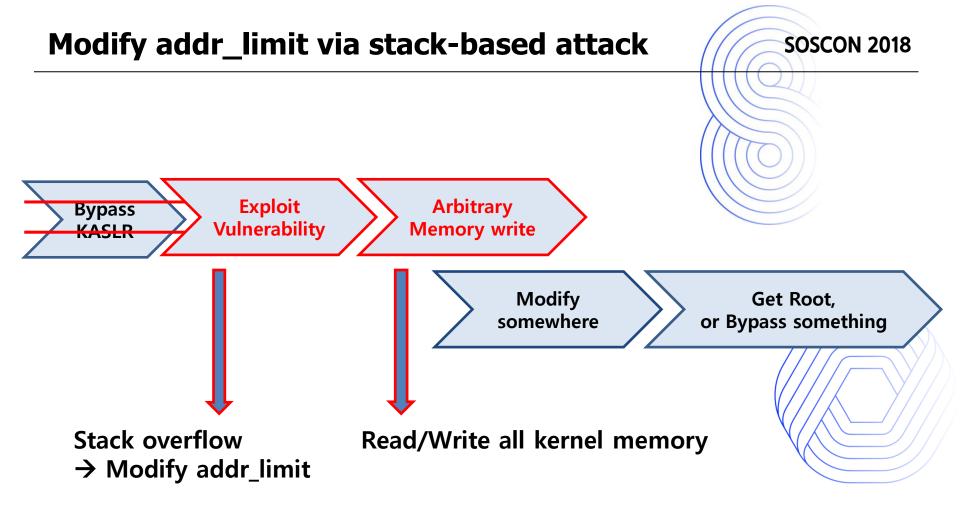




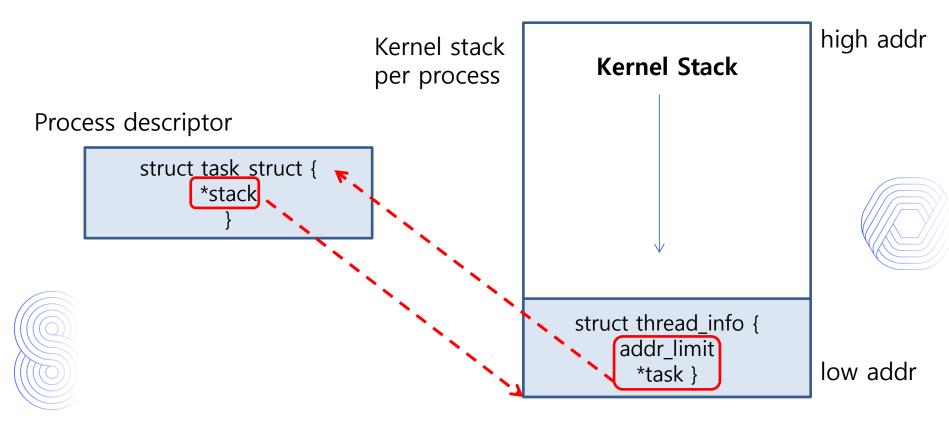
Enforce security checking to eliminate human errors!!



Attack-4: Modify addr_limit via stack-based attack

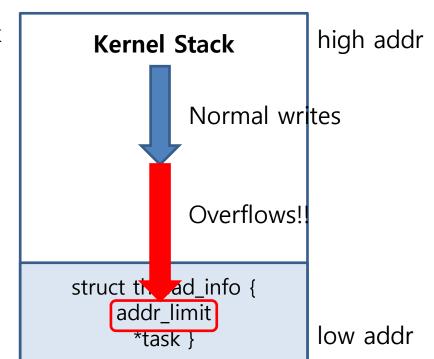


Where "addr_limit" be stored? In kernel stack!!



How about trying stack overflow attack as a classic?

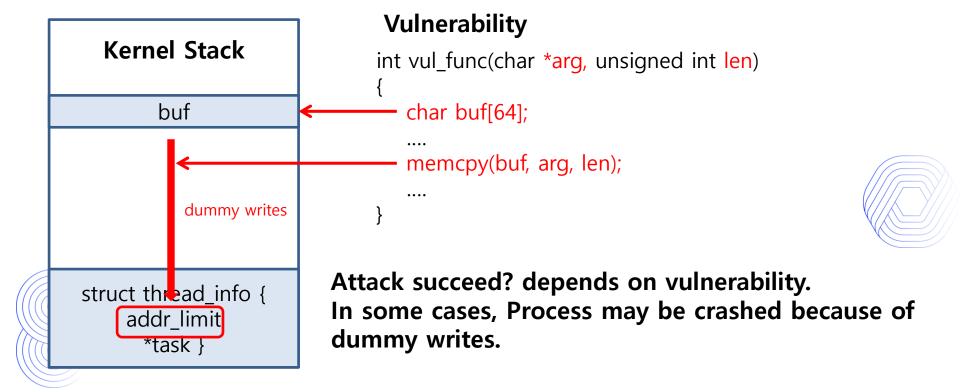
Kernel stack per process



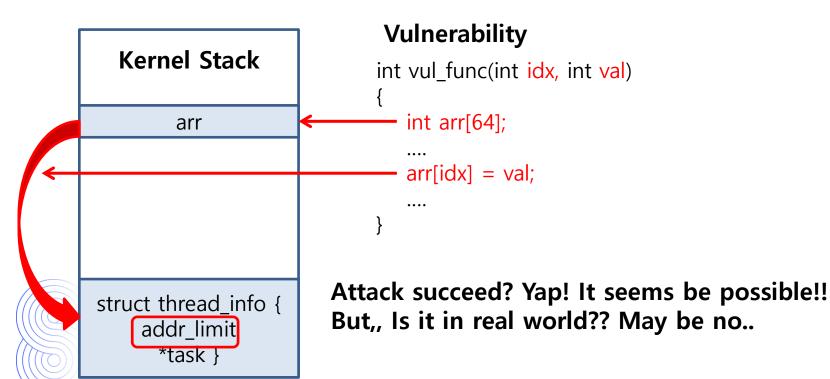




Stack overflow – Type 1: classic buffer overflow



Stack overflow – Type 2: out-of-bound index



Stack overflow – Type 3: VLA (Variable Length Array)

Kernel Stack

arr

struct thread_info {
 addr_limit
 *task }

Vulnerability

```
int vul_func(int size, int off, int val)
{
    int arr[size];
    ....
    for (i=0; i<size; i++)
        arr[i] = val;
    ....
}</pre>
```

Attack succeed? Depends on vulnerability. Is it in real world?

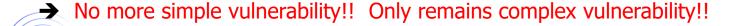
→ CVE-2010-3848, CVE-2010-3850

Stack overflow – Type 4: Recursion

```
Vulnerability
  Kernel Stack
                             int vul func(char *str)
                                 char buf[64];
                                if (~~)
                                   vul func(str);
        buf
                                strcpy(buf, str);
struct thread info
                           Attack succeed? Too difficult...
    addr_limit
                           Is it in real world?
      *task }
                           → CVE-2016-1583
```

Stack overflow – Summary

- Type1: Classic buffer overflow, Simple, No vulnerability these days
- Type2: Out-of-bound index, Simple, No vulnerability these days
- Type3: VLA (Variable Length Array), Complex, Real-world vulnerability
- Type4: Recursion, Complex, Real-world vulnerability







Why "struct thread_info" be in kernel stack??

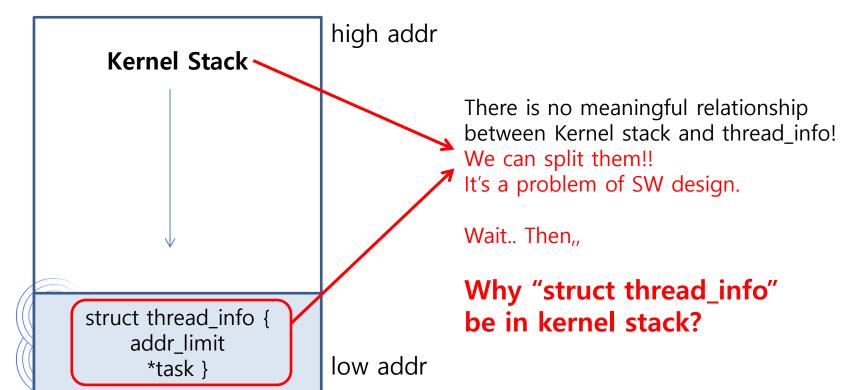
Kernel Stack struct thread info { addr limit *task }

high addr

If "struct thread_info" can be stored somewhere not related to kernel stack,,
Safe against the previous stack-based attack!!

low addr

Why "struct thread_info" be in kernel stack??



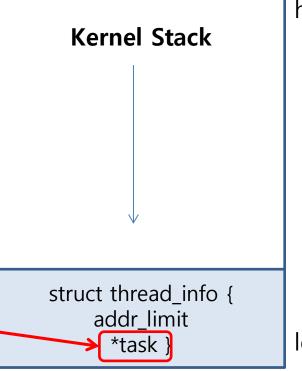
Stack pointer to point task_struct

- Access from register is faster than from memory.
- There is a register to point kernel stack, called SP.
- There are a lot of accesses to task.
- If thread_info is in kernel stack,
 We can access task through SP reg.
- So that,, Performance is improved!

SP (Stack Pointer)

Access from register!!

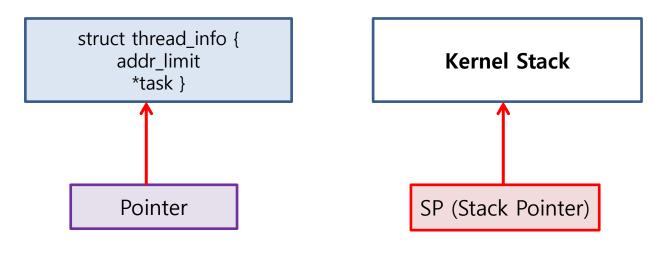
Too fast!!



high addr

low addr

Split addr_limit from stack



"Pointer" is needed for pointing thread_info instead of SP.

Performance:

Register

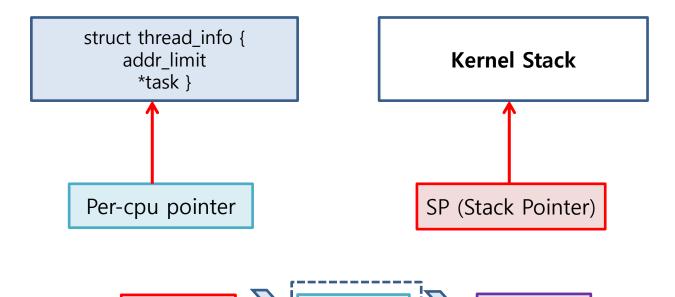


Security is ok, But..
Performance overhead here!!

Register

Optimization on Intel x86_64

Performance:

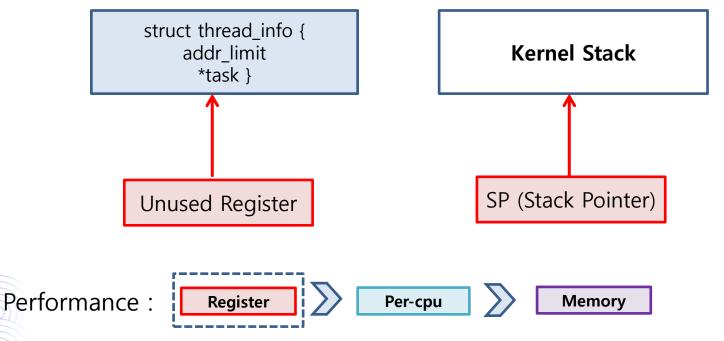


Per-cpu

Security is ok, and Overhead is not bad!!

Memory

Optimization on ARM 64



Security is ok, and Overhead is near zero!!

Tradeoff between Security and Performance

Performance : Register >> Per-cpu >> Memory

Security: Memory > Per-cpu > Register

Always there is a tradeoff between Performance and Security..

Are "Register" and "Per-cpu" really safe?? Hmm...



Defense solution for fixing SW design problem have to satisfy both Security and Performance!!

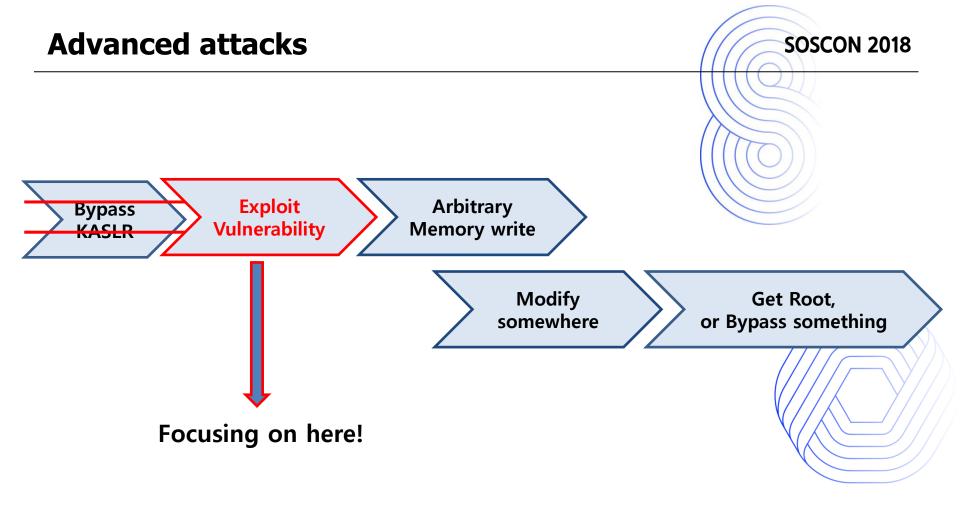








SAMSUNG OPEN SOURCE CONFERENCE 2018



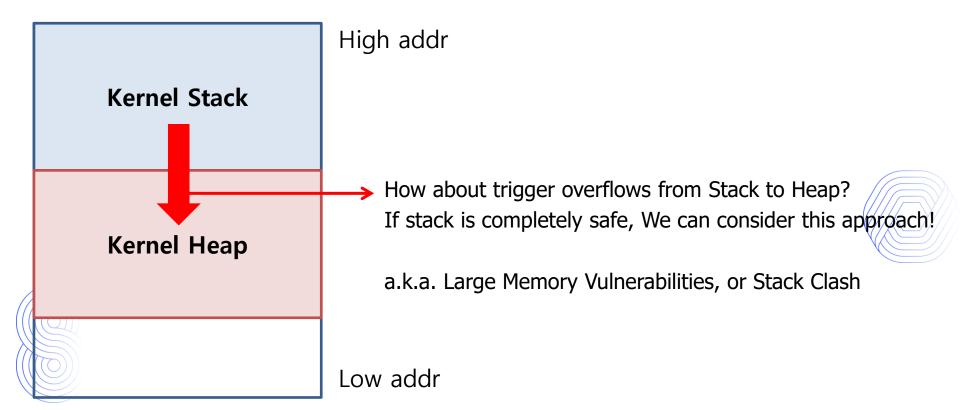
Advanced attacks Pick two keywords of advanced attacks

Adjacent / Spraying

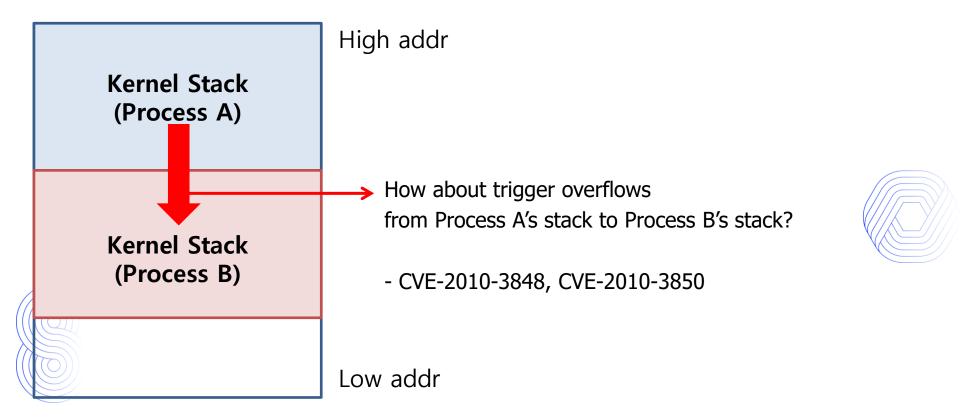




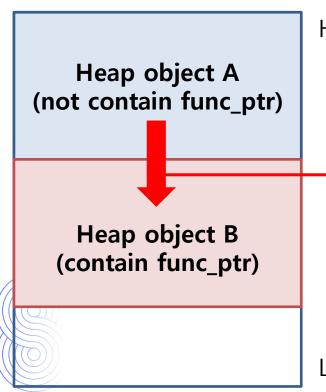
Adjacent, Type1: Heap / Stack



Adjacent, Type2: Stack of Process A / Stack of Process B



Adjacent, Type3: Heap object A / Heap object B



High addr

How about trigger overflows from Heap object A to Heap object B?

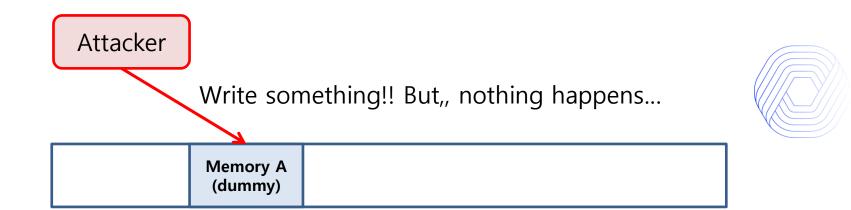
Attacker can't modify func_ptr in object A, But, Can modify func_ptr in object B!!

Low addr

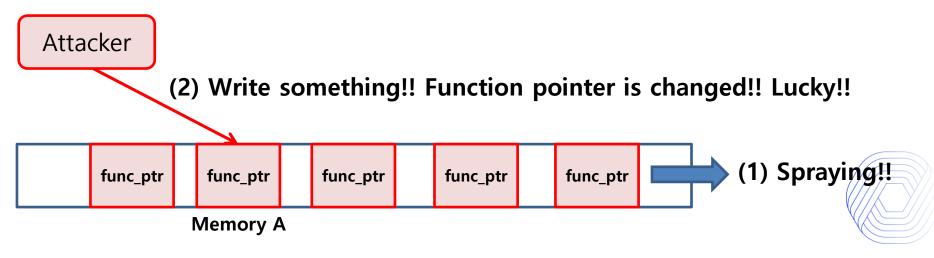


Advanced attacks Spraying

- Assume that attacker get an ability to write value to kernel memory A.
- Kernel memory A is random address. Attacker doesn't know what here it is.



Spraying





THANK YOU

Sample exploit code is at

https://github.com/jinb-park/linux-exploit/tree/master/samples/adjacent-kstacks

SOSCON 2018

SAMSUNG OPEN SOURCE CONFERENCE 2018