

Container Security

Chih-Hsuan Yang
National Sun-Yet-San University, Taiwan
Bachelor's degree graduation project
Advisor: Chun-I Fan

1. Abstract

Recently, many companies use containers to run their microservices, since containers could make their hardware resources be used efficiently. For example, GCP(Google Cloud Platform), AWS(Amazon Web Services), and Microsoft Azure are using this technique to separate subscribers' resources and services. However, if the hacker attacks the kernel or gets privilege escalation of containers, then such attacks would influence the host or the other containers. Therefore, this research would analyze, implement, and protect the container escalation. The container escalation could inspire in a container and influence the host or the other containers. This project would use the medical system with FHIR(Fast Healthcare Interoperability Resources) to simulate the real world threat and purpose a secure solution to protect patient's privacy.

2. Motivation

The Container is a virtualization technique to package applications and dependencies to run in an isolated environment. Containers are faster to start-up, lighter in memory/storage usage at run time and easier to deploy than virtual machines. Because the container shares the kernel with the host OS and other containers, and deploys by a configure file. First, we often used to run a docker container to host our services. For example: as-

signments, servers and some services in Information security club at NSYSU(National Sun Yat-sen University). But there are some threats about container technique. Like "Dirty CoW[1]" and "Escape vulnerabilities".

"Dirty CoW is a vulnerability in the Linux kernel. It is a local privilege escalation bug that exploits a race condition in the implementation of the copy-on-write mechanism in the kernel's memory-management subsystem"[2]. It founded by Phil Oester. We were 16, the first year we had touched the docker container. We tried to use the Dirty CoW vulnerability to take the root privilege of my Android phone.

Escape vulnerability is a subcategory of sandbox security. At first, security researchers often need sandbox to help they analyze malware, which prevent the malware influence researcher's host OS. Nowadays, the sandbox not only be used in analyzing, but also used to execute a normal application for an isolated environment. However if the application could modify the outside resources without the kernel permission. That loses the purpose of isolation. That might cause the information leaked or the kernel be hacked.

Hence, there is a big problem about: "How to make sure my services isolated and secure?" The author of this paper is the leader of Information security club. He should maintain all the services working perfectly. Moreover we are information security club. Therefore, the security and performance issue is the top-

priority requirement.

Second, the medical system is the most famous part internationally. Including face the COVID-19 in Taiwan, we not have the local COVID-19 case in more than 250 days.[3] However, the medical profession needs to renew the exchanging EHR(electronic health records) system these years. In order to protect the privacy of patients, and producing the high performance system to exchange the EHR, we need a easy deployment, effective runtime, and secure system. We have to do this research in this project.

3. Papers Review

3.1. Study of the Dirty Copy On Write

In this paper[4] show the race condition, and the mechanism of "copy on write". "Copy on write" is " a resource-management technique used in computer programming to efficiently implement a "duplicate" or "copy" operation on modifiable resources." [5] It often be inspire when fork() or mmap().

mmap This is a system call of mapping files or devices in to memory, which creates a new mapping in the virtual address space of the caller process. Such that the process could operate the instance of file in memory directly. And some libraries are also mapped into the virtual address space to share and handle the function call. Therefore, the processes could take the same view of libraries in it's memory space.

Copy on write This mechanism purposes of a resource being duplicated but not modified, it is not necessary to create a new resource. Therefore the kernel can make callers share the same memory resources. The mmap is a system call could inspire this mechanism in above paragraph. Some processes request

same memory resources, and the kernel supplies the same memory page to callers.

Race condition That is processes or threads are racing the same modifiable resources. For example: There is a page(4KB as default) of memory initialized as 0. And there are 2 threads or processes share that page. Consider one of the tasks is assigning the page full of character 'A'. However, the scheduler context switches to the other task assigns that page full of 'B'. And the scheduler context switches again to the first task. There is a problem now. What is the page for the first task looked like? Obviously, it does not meet the expectation for the first task.

3.1.1. Dirty CoW demo code. Let's analyze the proof of concept(PoC) of dirty CoW.(Oester, 2016)[1] The key of inspiring this vulnerability is the mmaped memory space, which is mapped with the PROT_READ flag. The PROT_READ flag declares the page is read only.

```
87 f=open(argv[1],O_RDONLY);
88 fstat(f,&st);
89 name=argv[1];
90 map=mmap(NULL,st.st_size,PROT_READ
    ,MAP_PRIVATE,f,0);
```

src/dirtyc0w.c

It creates 2 threads, which would have a race condition of the mmaped memory space, [madviseThread](#) and [procselmemThread](#). threads in main

```
106 pthread_create(&pth1,NULL,
    madviseThread,argv[1]);
107 pthread_create(&pth2,NULL,
    procselmemThread,argv[2]);
```

src/dirtyc0w.c

In one thread, call a system call "madvise", would make the user thread gain the root privilege to operate the protected page temporary. And the flag MADV_DONTNEED

would tell the kernel: "Do not Expected access it in the near future.[6]" Moreover, this flag might not lead to immediate freeing of pages in the range. The kernel is free to delay free the pages until an appropriate moment.[6]

madviseThread

```
33 void *madviseThread(void *arg)
34 {
35     char *str;
36     str=(char*)arg;
37     int i,c=0;
38     for(i=0;i<1000000000;i++)
39     {
40         c+=madvise(map,100,MADV_DONTNEED
41         );
42     }
43     printf("madvise %d\n\n",c);
44 }
```

src/dirtycow.c

In another thread, open its memory resource file. This file is a special file, which allow the process reads its memory by itself. Than, we move the printer of file descriptor of the memory resource file to the mmaped space. And try to write it. But the mmaped space is a read only space. We expected the kernel would create a copy of the this space and write the copy[7]. procselfmemThread

```
50 void *procselfmemThread(void *arg)
51 {
52     char *str;
53     str=(char*)arg;
54     int f=open("/proc/self/mem",O_RDWR
55     );
56     int i,c=0;
57     for(i=0;i<1000000000;i++) {
58         lseek(f,(uintptr_t) map,SEEK_SET
59         );
60         c+=write(f,str,strlen(str));
61     }
62     printf("procselfmem %d\n\n", c);
63 }
```

src/dirtycow.c

But there is a problem! There is an another thread is racing this page with root privilege.

If the scheduler context switches the madviseThread to procselfmemThread, while the adviseThread is calling the "madvise" system call. It would cause the procselfmemThread gain the root privilege from madviseThread to control the mmaped file.

3.2. Container Security: Issues, Challenges, and the Road Ahead

This paper[8] has derived 4 generalized container security issues: (I) protecting a container from applications inside it, (II) inter-container protection, (III) protecting the host from containers, and (IV) protecting containers from a malicious or semi-honest host.[8]

The Dirty CoW vulnerability is a exploit from kernel. But the benefit of container and host OS are share the same kernel. This vulnerability can be used in container to attack the kernel, and gives this application root privilege, changes this containers as a privileged container or supervises the other containers. Therefore, we should protect the host form the container(which belongs to type (III) threat in this paper).

3.2.1. Virtual machine and container. // FIXME: Draw the architecture of VM and container.

3.2.2. Linux kernel features. // FIXME: Introduce these features for isolating processes in Linux.

namespaces

// FIXME: Namespaces perform the job of isolation and virtualization of system resources for a collection of processes.[8]

cgroups

// FIXME: Limits, accounts for, and isolates the resource usage of a collection of processes. [9]

capabilities

// FIXME: Divide the privileges traditionally associated with superuser into distinct units.

seccomp

// FIXME: Only some specified process could call some specified system calls.

3.3. Study of Security Flaws in the Linux Kernel by Fuzzing

This paper[10].

3.4. PINE: Optimizing Performance Isolation in Container Environments

This paper[11].

4. Methods

4.1. Study CVEs about the Linux kernel

This project will study several kernel vulnerabilities were discover in packet socket: CVE-2016-8655 [12], CVE-2017-7308[13], and CVE-2020-14386[14].

And study some kernel exploit techniques[15], because the container shares the kernel. If I could exploit the kernel in the suffering container, it might have more chance to influence the other containers or host.

4.2. Study related mechanisms

The Linux kernel is a monolithic kernel, which is over 28 million lines of codes now(2020). There are many mechanisms to solve the real life situations. Study those CVEs' related mechanism in the kernel, might have more chance to find new vulnerabilities.

4.3. Implement a simple container

The Linux kernel supply some system calls to clone a process(also in thread) in their own namespace and group. We could implement a simple container by ourself, so that we

can make a list of vulnerabilities may happen.

```
1  âĢĤ container git:(main) âĢĤ sudo
   ./c bash
2  Success on creating container
3  Start container: bash with clone id:
   375696
4  In container PID: 1
5  bash-5.0# ./test.sh
6  This is the self test script in
   container!
7  Support bash cat echo ls rm hostname
   tree, 7 commands.
8  /bin/bash
9  ./test.sh
10 -----FILE: ./test.sh -----
11     1  #!/bin/bash
12     2
13     3  echo "This is the self test
       script in container!"
14     4  echo "Support bash cat echo
       ls rm hostname tree, 7 commands."
15     5
16     6  echo "$SHELL"
17     7  echo $0
18     8
19     9  echo "-----FILE: $0
       -----"
20    10  cat -n $0
21    11  echo
       "-----"
22    12
23    13  echo $(hostname) >
       âĢĤ'l'ĉňžéi jăěžĹěžĹ
24    14  cat âĢĤ'l'ĉňžéi jăěžĹěžĹ
25    15  rm âĢĤ'l'ĉňžéi jăěžĹěžĹ
26    16  ls
27 -----
28 container
29 bin dev etc home lib lib64 mnt
       opt proc root run/sbin sys
       test.sh tmp usr var
30 bash-5.0# exit
31 exit
32 âĢĤ container git:(main) âĢĤ
```

src/lc_out.txt

4.4. List secure issues of the simple container

4.5. Aim a vulnerability

After studying the CVEs and related mechanisms. This project would find a vulnerability of privilege escalation in the container by attacking the kernel.

4.6. Implement the PoC

Implement the proof of concept of the aimed vulnerability.

4.7. Implement the patch and pull request

Being a security researcher, we cannot just only exploit the software, but also give patches to the maintainer. The Linux kernel is an open source project under GPL-2.0 license in GitHub. I would pull request to the repository. If my patch could be merged into the kernel to solve the container vulnerability.

5. Expected Outcome

Would research some related vulnerabilities, and implement the PoC code. Moreover this project will generate the patch of the vulnerability(s) to protect these attack(s).

6. References

References

- [1] Phil Oester. *Dirty CoW CVE-2016-5195*. URL: <https://dirtycow.ninja/>.
- [2] Wikipedia. *Dirty CoW*. URL: https://en.wikipedia.org/wiki/Dirty_COW.
- [3] James Griffiths Chandler Thornton. *Taiwan reports first local Covid-19 case in more than 250 days*. URL: <https://edition.cnn.com/2020/12/22/asia/taiwan-coronavirus-intl-hnk/index.html>.
- [4] Tanjila Farah Delwar Alam Moniruz Zaman. "Study of the Dirty Copy On Write, A Linux Kernel Memory Allocation Vulnerability". In: 2017. URL: <https://ieeexplore.ieee.org/abstract/document/7530217>.
- [5] Wikipedia. *Copy-on-write*. URL: <https://en.wikipedia.org/wiki/Copy-on-write>.
- [6] GNU. *Manpage of madvise*. URL: <https://www.man7.org/linux/man-pages/man2/madvise.2.html>.
- [7] Babak D. Beheshti A.P. Saleel Mohamed Nazeer. "Linux kernel OS local root exploit". In: 2017. URL: <https://ieeexplore.ieee.org/document/8001953>.
- [8] Tassos Dimitriou Sari Sultan Imtiaz Ahmad. "Container Security: Issues, Challenges, and the Road Ahead". In: *IEEE Access* 7.18620110 (2019).
- [9] Wikipedia. *cgroups*. URL: <https://en.wikipedia.org/wiki/Cgroups>.
- [10] E.V. Sharlaev P.A. Teplyuk A.G. Yakunin. "Study of Security Flaws in the Linux Kernel by Fuzzing". In: 2020. URL: <https://ieeexplore.ieee.org/document/9271516>.
- [11] Congfeng Jiang Youhuizi Li Jiancheng Zhang. "PINE: Optimizing Performance Isolation in Container Environments". In: *IEEE Access* 7.18526707 (2019).
- [12] Inc. Red Hat. *CVE-2016-8655*. URL: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2016-8655>.
- [13] MITRE Corporation. *CVE-2017-7308*. URL: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2017-7308>.

- [14] Inc. Red Hat. *CVE-2020-14386*. URL: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2020-14386>.
- [15] xairy. *Linux Kernel Exploitation*. URL: <https://github.com/xairy/linux-kernel-exploitation>.

7. Academic Advisor

- Organize to a complete structure.
- Extend to a formal paper, and publish.