

# **CS315 Computer Security Term Project**

## **Docker Escape - Final Report**

11712121 胡玉斌/11811902 顾同舟/11812102 黄炜杰/11812103 张文灏/11812104 石文轩



# Docker Escape?

Docker escape is to **escape from the preset limitations of containerization**.

Linux provides multitudinous of limitations for containers:

- Permission limitations
  - Isolation among containers
  - Isolation between containers and host
- Resource limitations
  - CPU cores
  - Memory usage
  - Derived process numbers

# Why break the permission limitations?

Hackers can take control of the host machine!

e.g. Docker containers farm (Amazon, Ali Cloud, Tecent Cloud...)

~~So it's super hard.~~

# Why break the resource limitations?

In a multi-tenant container environment, an adversarial container is able to:

1. Amplify the amount of consumed resources.
2. Significantly slow-down containers on the same host. **(DoS Attack)**
3. Gain extra unfair advantages on the system resources.



# Escape from permission limitation

CVE-2019-5736 (RunC Escape)

# RunC Escape

Overwriting the host system's RunC binary within a container.

Starting containers depends on RunC.

# Approach

1. Overwritten RunC binary by hack's code with *CVE-2019-5736 RunC vulnerability*.
2. Waiting for a new container starting.
3. 

Demo in week 10.

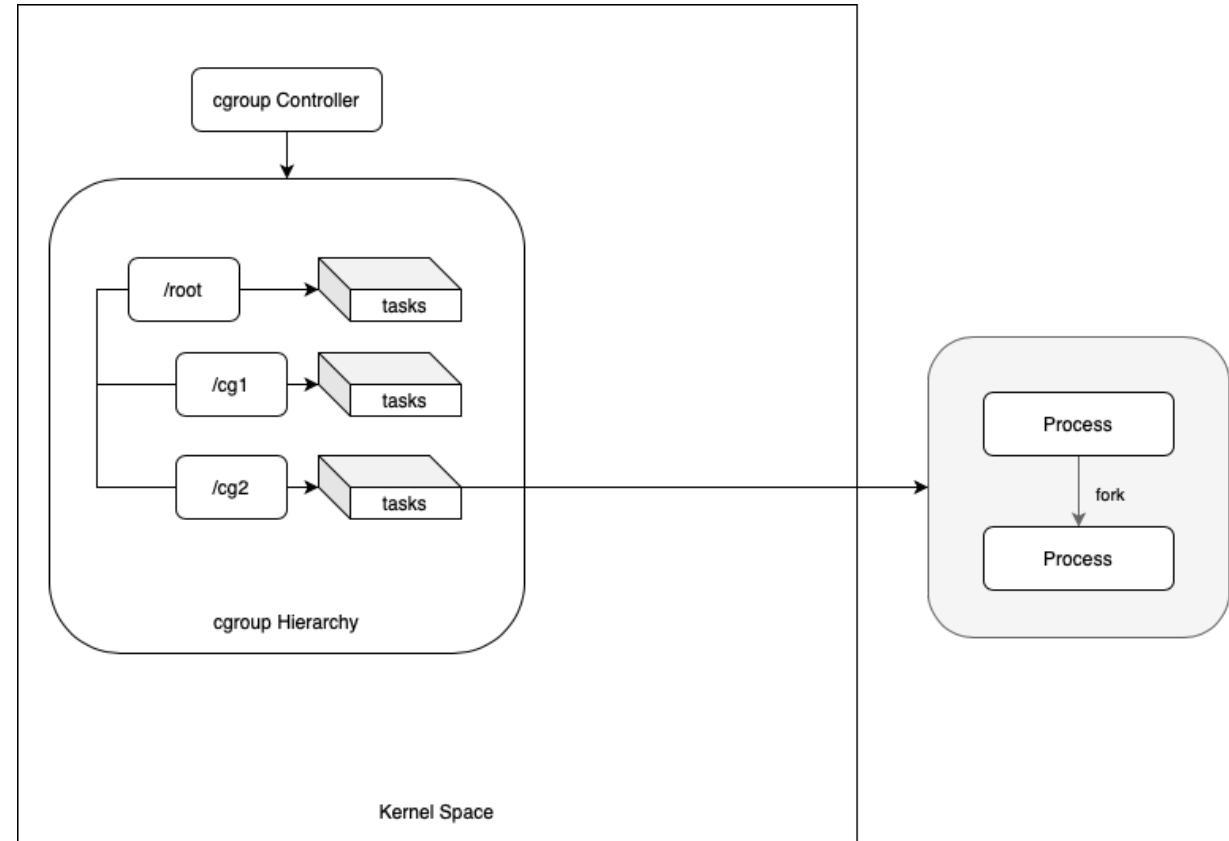


## Escape from resource limitation

# Background

Docker containers depend on **cgroups** for resource management.

**Partitions** a group of processes and their children into **hierarchical groups**, applies different controllers to manage and limit various system resources (CPU time, computer memory, block I/O, etc.)



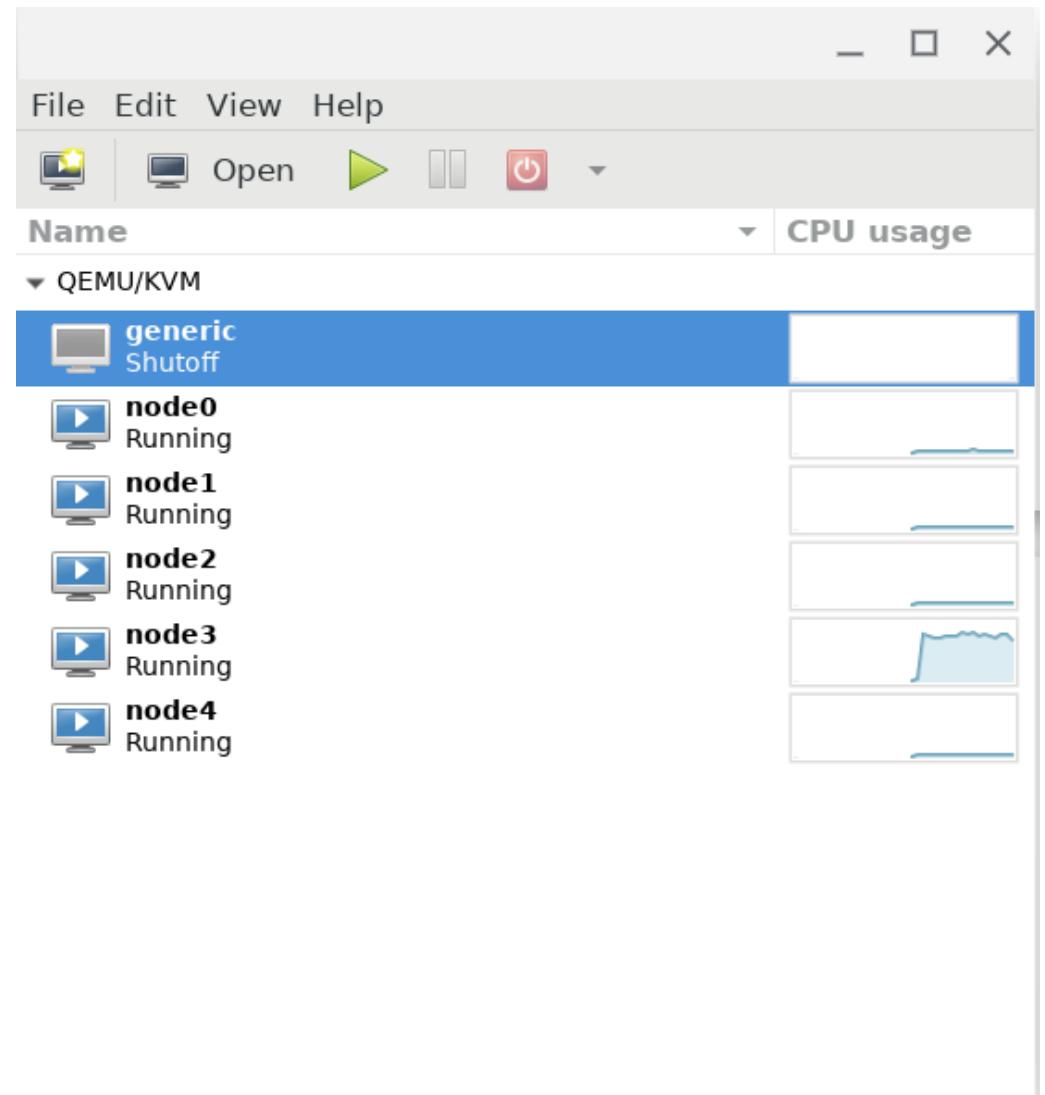
# Base Principle: cgroups miscounting

We will give 5 different fancy ways to trigger cgroups miscounting. But first...

# Environment Setup

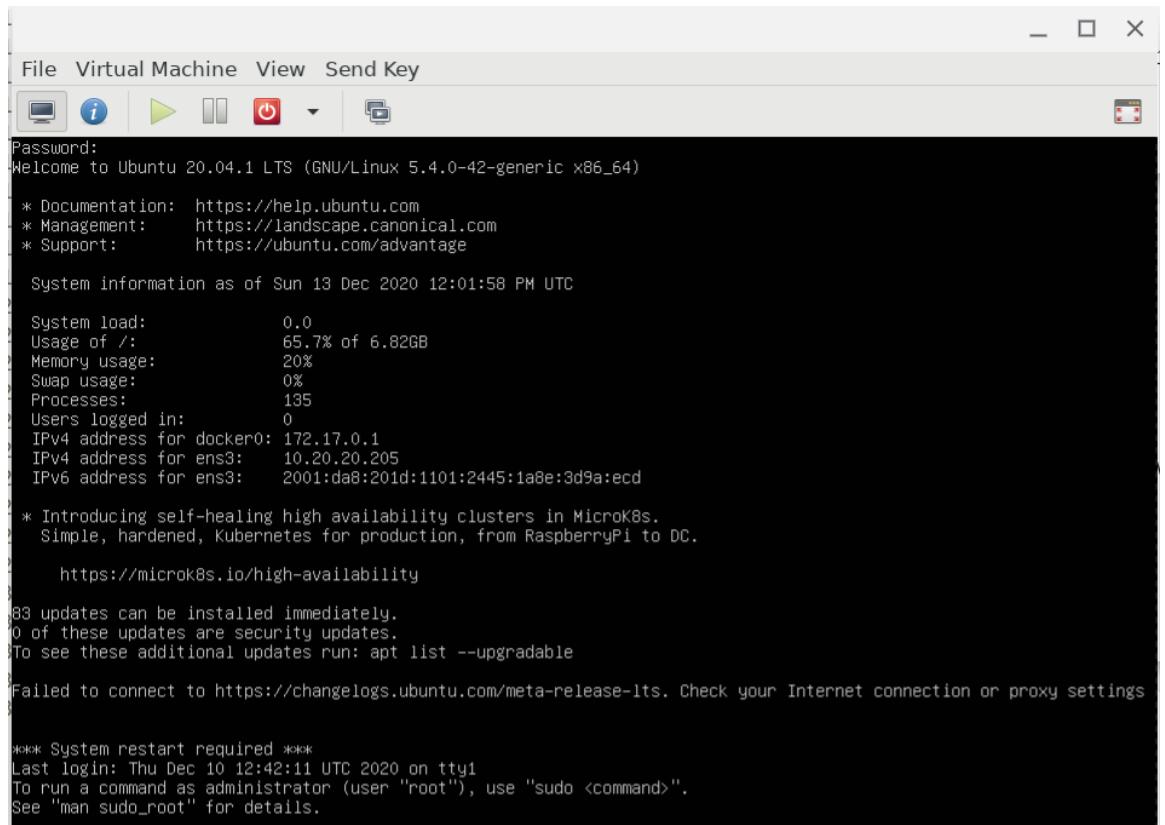
# Environment Setup

- Host Configuration
  - CPU: Intel Xeon E5-2697 v4  
(36 x 2.3GHz)
  - Mem: 512GB
  - Disk: 1TB SSD
  - OS: Ubuntu 18.04
  - Kernel: Linux 4.15.0



# Environment Setup

- Guest Configuration
  - Hypervisor: QEMU + KVM
  - CPU: 4 Cores
  - Mem: 2GB
  - Disk: VirtIO + QCOW2
  - Net: VirtIO + NAT
  - OS: Ubuntu 20.04
  - Docker: 19.03.11

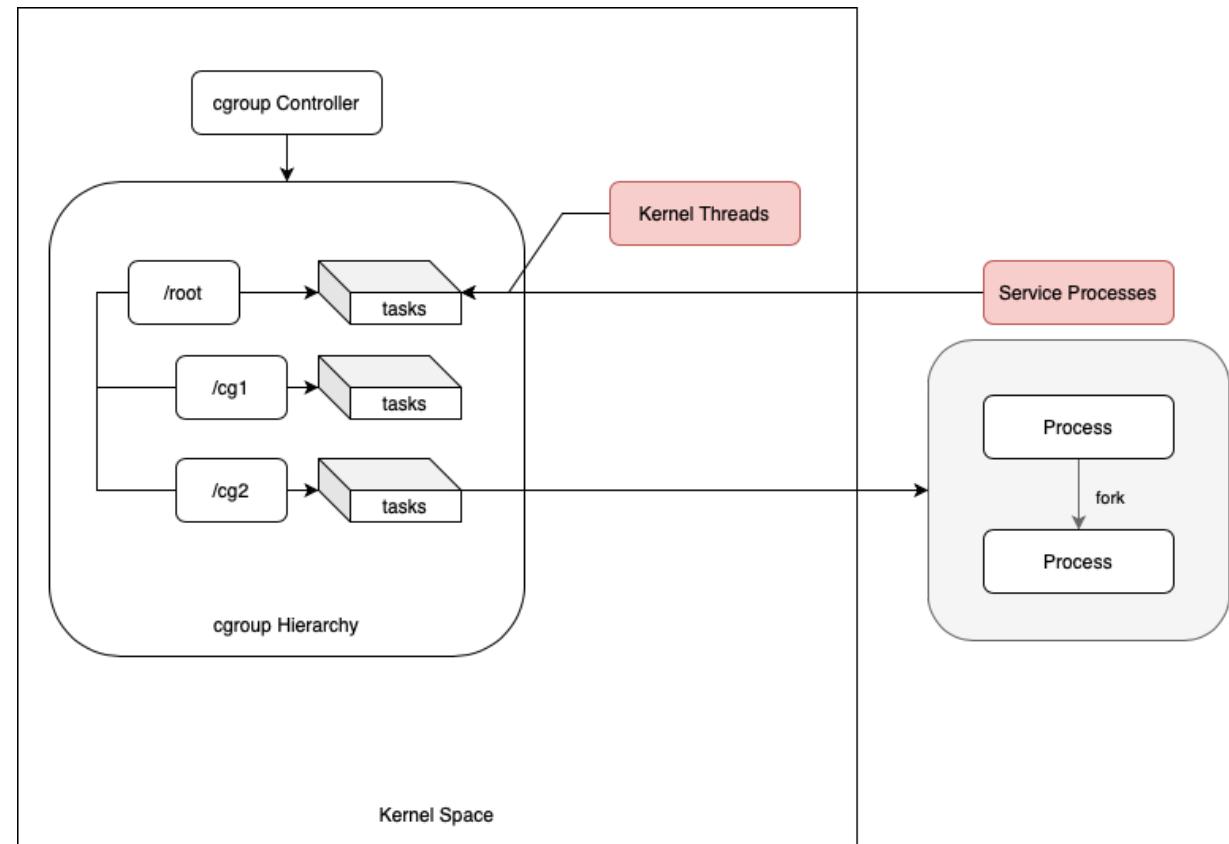


# **Part 1: Exploiting Service Processes**

## **Conainer Engine**

# What will happen when sending command to container?

- → CLI Client (i.e. docker )
- → Docker REST API
- → Docker Daemon (Service Processes)  
(i.e. dockerd )
- → tty Driver (Kernel Threads)
  - → Executing work queue (*kworker* kernel thread)
  - → LDISC



## How to exploit it?

Single Line Command `for((;)); do lsmod; done`

Repeatedly generate tons of meaningless output to stress Docker Daemon and [tty Driver](#)

# Evaluation

- Create container
  - Command: `docker run -itd --name atk4 --cpus=1 ubuntu:20.04 bash`
  - CPU Limit: 100% Single Core
- Launch attack
  - Command: `docker exec -it atk4 bash , for((;)); do lsmod; done`
  - Comparison: `for((;)); do lsmod > /dev/null; done`

```
1 [██████████|          55.5%] Tasks: 41, 90 thr; 4 running
2 [███████████|         62.2%] Load average: 3.29 2.55 2.15
3 [███████████|         59.5%] Uptime: 3 days, 09:51:16
4 [███████████|         57.3%]
Mem[███████████|         310M/1.94G]
Swp[██████████|         3.25M/1.19G]
```

PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
2302	root	20	0	991M	98176	45660	S	42.0	4.8	3h10:26	dockerd --group docker
4116956	root	20	0	395M	57812	33508	S	32.4	2.8	1:21.90	docker exec -it atk4 ba
4114570	t	20	0	13900	5516	4024	S	32.4	0.3	1:13.10	sshd: t@pts/1
303970	root	20	0	106M	6940	4280	S	44.1	0.3	5h13:02	containerd-shim -namesp
4117042	root	20	0	395M	57812	33508	S	0.0	2.8	0:23.12	docker exec -it atk4 ba
2376	root	20	0	991M	98176	45660	S	9.0	4.8	14:14.18	dockerd --group docker
4116991	root	20	0	395M	57812	33508	S	9.6	2.8	0:25.16	docker exec -it atk4 ba
303978	root	20	0	106M	6940	4280	S	0.0	0.3	46:32.74	containerd-shim -namesp
2371	root	20	0	991M	98176	45660	S	10.3	4.8	52:25.74	dockerd --group docker
2450	root	20	0	991M	98176	45660	S	5.5	4.8	14:32.68	dockerd --group docker
2434	root	20	0	991M	98176	45660	S	0.0	4.8	15:41.41	dockerd --group docker
304024	root	20	0	106M	6940	4280	S	10.3	0.3	39:13.10	containerd-shim -namesp
2372	root	20	0	991M	98176	45660	S	6.2	4.8	13:56.40	dockerd --group docker

55

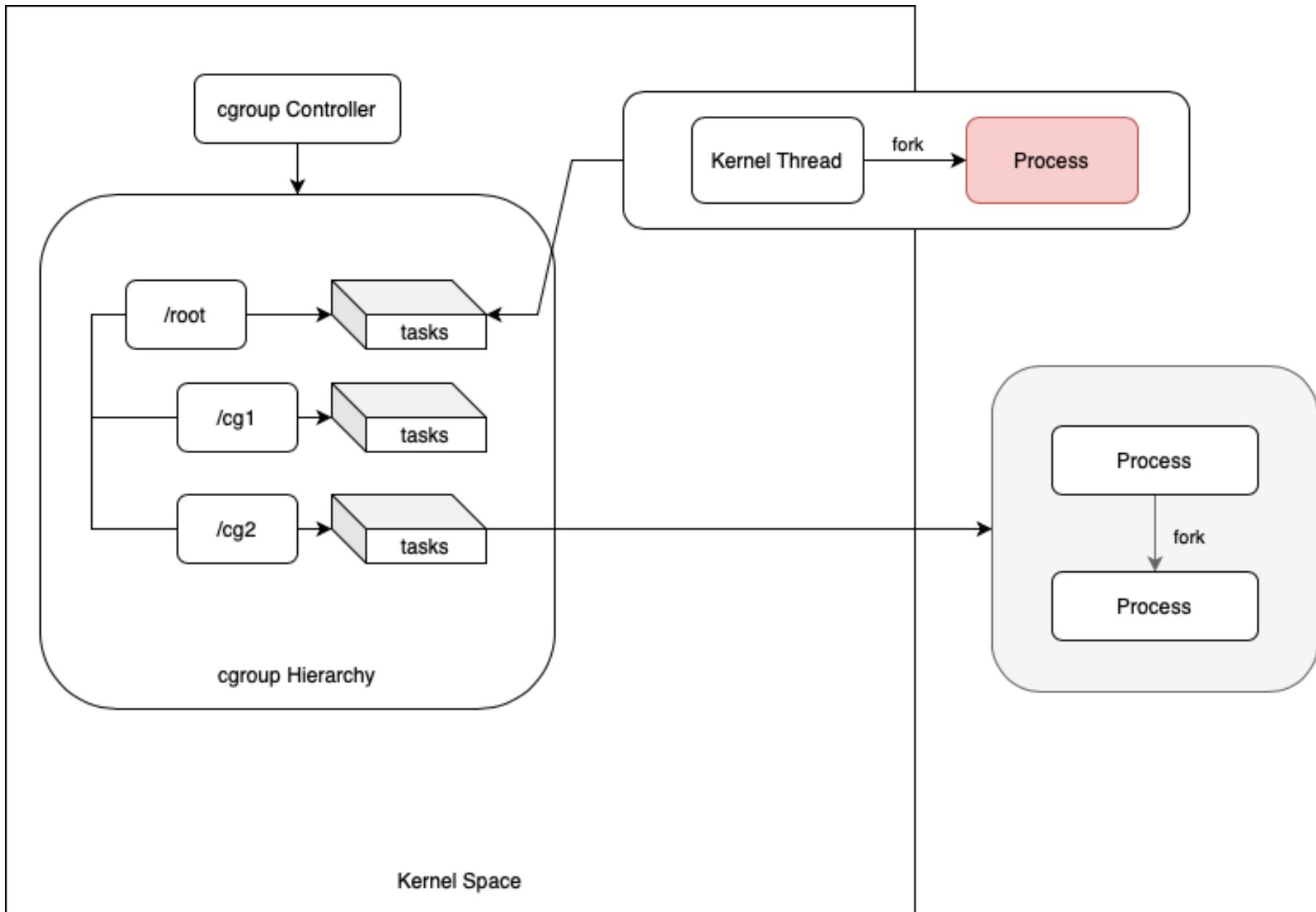
PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
4093965	t	20	0	8276	5156	3460	S	10.0	0.3	4:13.55	-bash
431652	t	20	0	7948	3984	3368	R	0.7	0.2	0:00.13	htop
1716870	t	20	0	8272	4304	3368	S	0.7	0.2	2:36.09	htop
2380	root	20	0	901M	46180	24444	S	0.0	2.3	25:54.64	containerd --config /var
2302	root	20	0	991M	101M	45660	S	0.0	5.1	3h11:52	dockerd --group docker
2371	root	20	0	991M	101M	45660	S	0.0	5.1	52:43.11	dockerd --group docker
122513	root	20	0	901M	46180	24444	S	0.0	2.3	0:09.79	containerd --config /var
819	root	20	0	232M	7820	6796	S	0.0	0.4	0:06.00	/usr/lib/accountsservic
817	root	20	0	232M	7820	6796	S	0.0	0.4	0:06.06	/usr/lib/accountsservic
137953	root	20	0	901M	46180	24444	S	0.0	2.3	0:09.66	containerd --config /var
139171	root	20	0	901M	46180	24444	S	0.0	2.3	0:10.75	containerd --config /var
554	root	RT	0	273M	17952	8184	S	0.0	0.9	0:27.15	/sbin/multipathd -d -s
2415	root	20	0	901M	46180	24444	S	0.0	2.3	4:48.45	containerd --config /Va
2417	root	20	0	901M	46180	24444	S	0.0	2.3	4:48.45	containerd --config /Va

# Result

Occupy about 250% resource

## **Part 2: Exploiting Processes Forked by Kernel Thread**

### **Exception Handling: Core Dump**



# Exception and Core Dump

- Core Dump is a kernel function.
- When exception happening, the process will halt, and invoke the Core Dump function.
- The Core Dump kernel function will invoke the core dump application in **usermode**.

# How to use it?

Core Dump application is created by kernel thread, rather than container. Thus it won't be controlled by cgroup's resource limit.

In Ubuntu, the default core dump application is Apport. Apport produces big dump files.

# Attack Strategy

1. Enable Core Dump in container.
2. Run `div 0` continuously in malicious container.

Apport will continuously generate core dump files.

# Experiment

- One container attached to CPU core 0, limited with CPU usage 5%.
- Measure CPU usage.



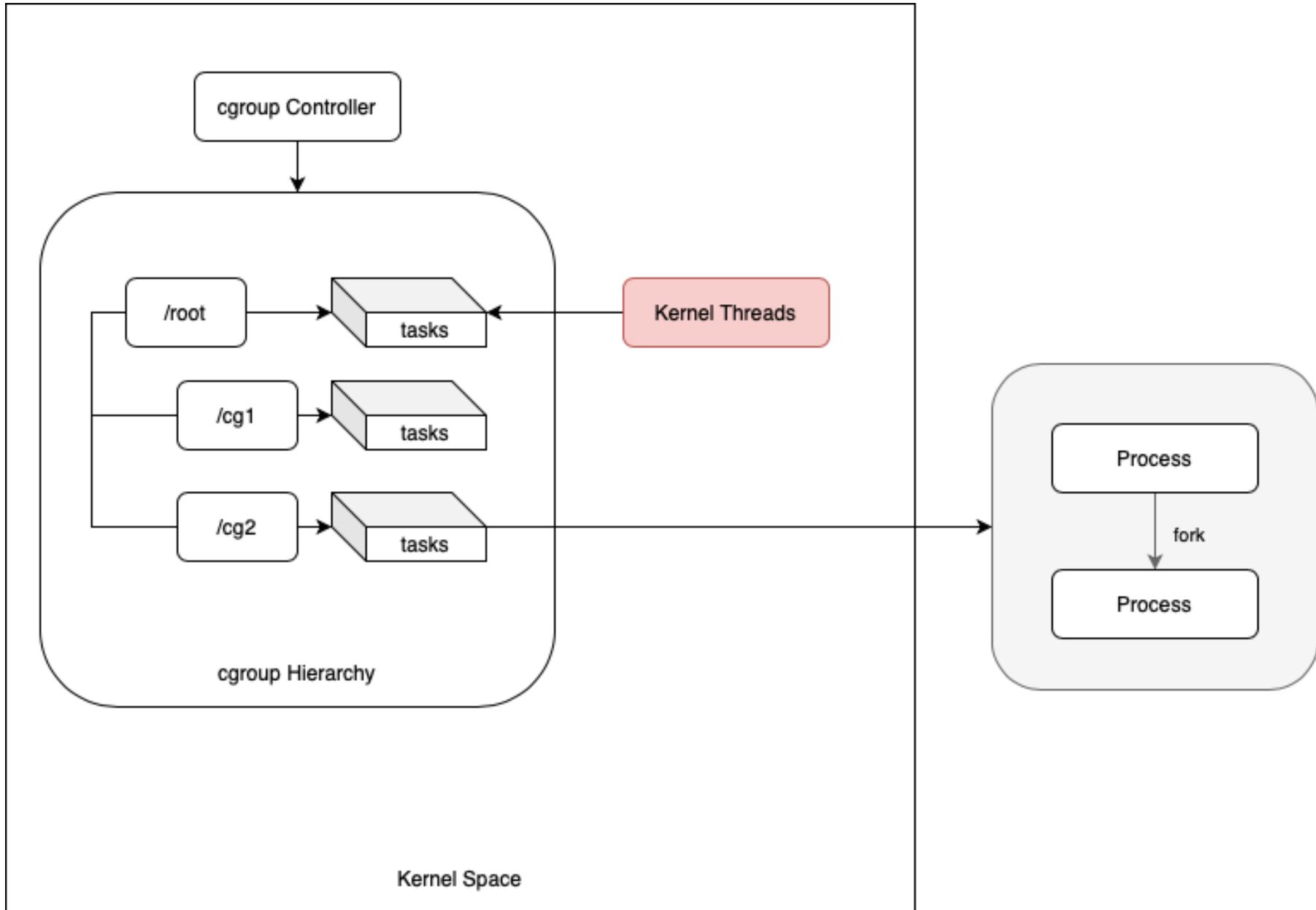
PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
10758	root	20	0	2004M	106M	53256	S	0.0	0.3	0:07.45	/usr/bin/dockerd -H fd:// --containerd
10759	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.00	/usr/bin/dockerd -H fd:// --containerd
10760	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.76	/usr/bin/dockerd -H fd:// --containerd
10761	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.44	/usr/bin/dockerd -H fd:// --containerd
10762	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.00	/usr/bin/dockerd -H fd:// --containerd
10764	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.18	/usr/bin/dockerd -H fd:// --containerd
10767	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.29	/usr/bin/dockerd -H fd:// --containerd
10768	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.26	/usr/bin/dockerd -H fd:// --containerd
10769	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.44	/usr/bin/dockerd -H fd:// --containerd
10770	root	20	0	2004M	106M	53256	S	1.3	0.3	0:01.66	/usr/bin/dockerd -H fd:// --containerd
10771	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.71	/usr/bin/dockerd -H fd:// --containerd
10772	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.29	/usr/bin/dockerd -H fd:// --containerd
10773	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.76	/usr/bin/dockerd -H fd:// --containerd
10774	root	20	0	2004M	106M	53256	S	0.0	0.3	0:08.79	/usr/bin/dockerd -H fd:// --containerd
10775	root	20	0	2004M	106M	53256	S	0.7	0.3	0:01.04	/usr/bin/dockerd -H fd:// --containerd
10776	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.63	/usr/bin/dockerd -H fd:// --containerd
10777	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.04	/usr/bin/dockerd -H fd:// --containerd
10778	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.67	/usr/bin/dockerd -H fd:// --containerd
10779	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.54	/usr/bin/dockerd -H fd:// --containerd
10780	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.56	/usr/bin/dockerd -H fd:// --containerd
10781	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.19	/usr/bin/dockerd -H fd:// --containerd
10782	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.81	/usr/bin/dockerd -H fd:// --containerd
10783	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.19	/usr/bin/dockerd -H fd:// --containerd
10784	root	20	0	2004M	106M	53256	S	0.7	0.3	0:02.46	/usr/bin/dockerd -H fd:// --containerd
11040	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.72	/usr/bin/dockerd -H fd:// --containerd
11045	root	20	0	2004M	106M	53256	S	0.0	0.3	0:02.22	/usr/bin/dockerd -H fd:// --containerd
11046	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.11	/usr/bin/dockerd -H fd:// --containerd
11751	root	20	0	2004M	106M	53256	S	0.0	0.3	0:04.09	/usr/bin/dockerd -H fd:// --containerd
12113	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.61	/usr/bin/dockerd -H fd:// --containerd
12114	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.05	/usr/bin/dockerd -H fd:// --containerd
12115	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.68	/usr/bin/dockerd -H fd:// --containerd
12116	root	20	0	2004M	106M	53256	S	0.0	0.3	0:01.12	/usr/bin/dockerd -H fd:// --containerd
12117	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.97	/usr/bin/dockerd -H fd:// --containerd
12118	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.65	/usr/bin/dockerd -H fd:// --containerd
12119	root	20	0	2004M	106M	53256	S	0.0	0.3	0:00.82	/usr/bin/dockerd -H fd:// --containerd

# Result

CPU usage increased: 2000%

## **Part 3: Exploiting Kernel Thread**

### **Data Synchronization**



# Writeback Strategy

- Cached data are written back only when the page is evicted
- Separate processes for **synchronization** and I/O  
(Could be utilized to escape from `cgroup`)

# System Call: `sync`

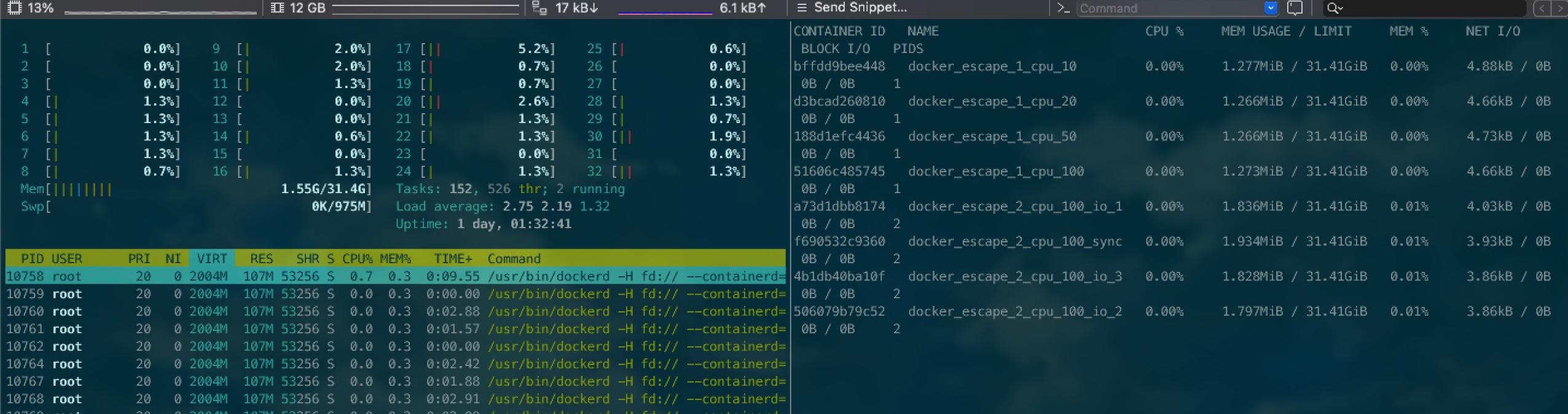
- Create kernel thread
- Scan the entire file system (not only the container, but the **host OS**)
- Force writing all **dirty** pages back into disk

# Attack Strategy

- Call `sync` continuously in malicious container
- I/O sensitive tasks will be affected

# Experiment

- Three containers performing file I/O
- One malicious container
- Each container restricted to a disparate CPU core
- Measure time consumption and CPU usage with/without the malicious system call  
`( time , htop )`



F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice -F8Nice +F9Kill F10Quit

```
huangweijie@lab-vm ~>
docker exec -it docker_escape_2_cpu_100_io_1 "/bin/
bash"
root@a73d1dbb8174:/# cd /data/docker_share/
root@a73d1dbb8174:/data/docker_share# time ./io

real    3m45.704s
user    3m45.554s
sys     0m0.033s
root@a73d1dbb8174:/data/docker_share#
```

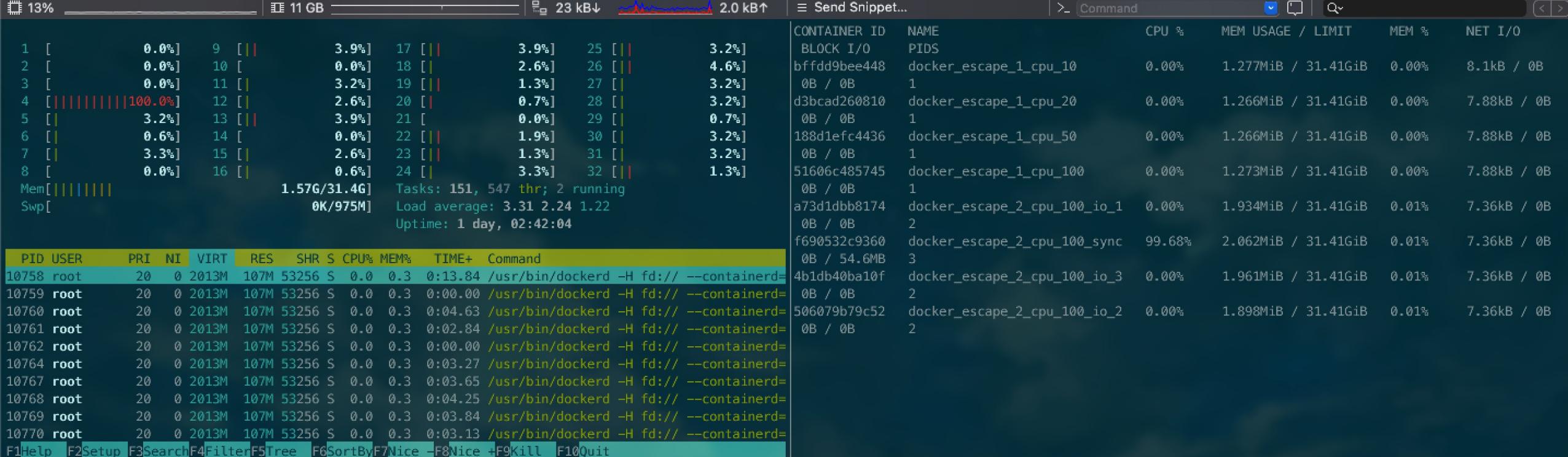
```
huangweijie@lab-vm ~>
docker exec -it docker_escape_2_cpu_100_io_2 "/bin/
bash"
root@506079b79c52:/# cd /data/docker_share/
root@506079b79c52:/data/docker_share# time ./io

real    3m45.299s
user    3m45.190s
sys     0m0.012s
root@506079b79c52:/data/docker_share#
```

```
huangweijie@lab-vm ~>
docker exec -it docker_escape_2_cpu_100_io_3 "/bin/
bash"
root@4b1db40ba10f:/# cd /data/docker_share/
root@4b1db40ba10f:/data/docker_share# time ./io

real    3m45.445s
user    3m45.396s
sys     0m0.004s
root@4b1db40ba10f:/data/docker_share#
```

```
huangweijie@lab-vm ~>
docker exec -it docker_escape_2_cpu_100_sync "/bin/
bash"
root@f690532c9360:/# cd /data/docker_share/
root@f690532c9360:/data/docker_share# ./sync
```



```
root@a73d1dbb8174:/data/docker_share# time ./io dat
a1.in data1.out
real    0m48.713s
user    0m17.402s
sys     0m31.269s
root@a73d1dbb8174:/data/docker_share# time ./io dat
a1.in data1.out
real    1m46.446s
user    0m19.548s
sys     0m34.366s
root@a73d1dbb8174:/data/docker_share# [REDACTED]
```

```
root@506079b79c52:/data/docker_share# time ./io da
ta2.in data2.out
real    0m49.450s
user    0m19.183s
sys     0m30.231s
root@506079b79c52:/data/docker_share# time ./io da
ta2.in data2.out
real    1m42.023s
user    0m20.662s
sys     0m33.589s
root@506079b79c52:/data/docker_share# [REDACTED]
```

```
root@4b1db40ba10f:/data/docker_share# time ./io da
ta3.in data3.out
real    0m51.897s
user    0m22.391s
sys     0m29.451s
root@4b1db40ba10f:/data/docker_share# time ./io da
ta3.in data3.out
real    1m49.410s
user    0m22.305s
sys     0m33.443s
root@4b1db40ba10f:/data/docker_share# [REDACTED]
```

```
root@f690532c9360:/data/docker_share# ./sync
root@f690532c9360:/data/docker_share# [REDACTED]
```

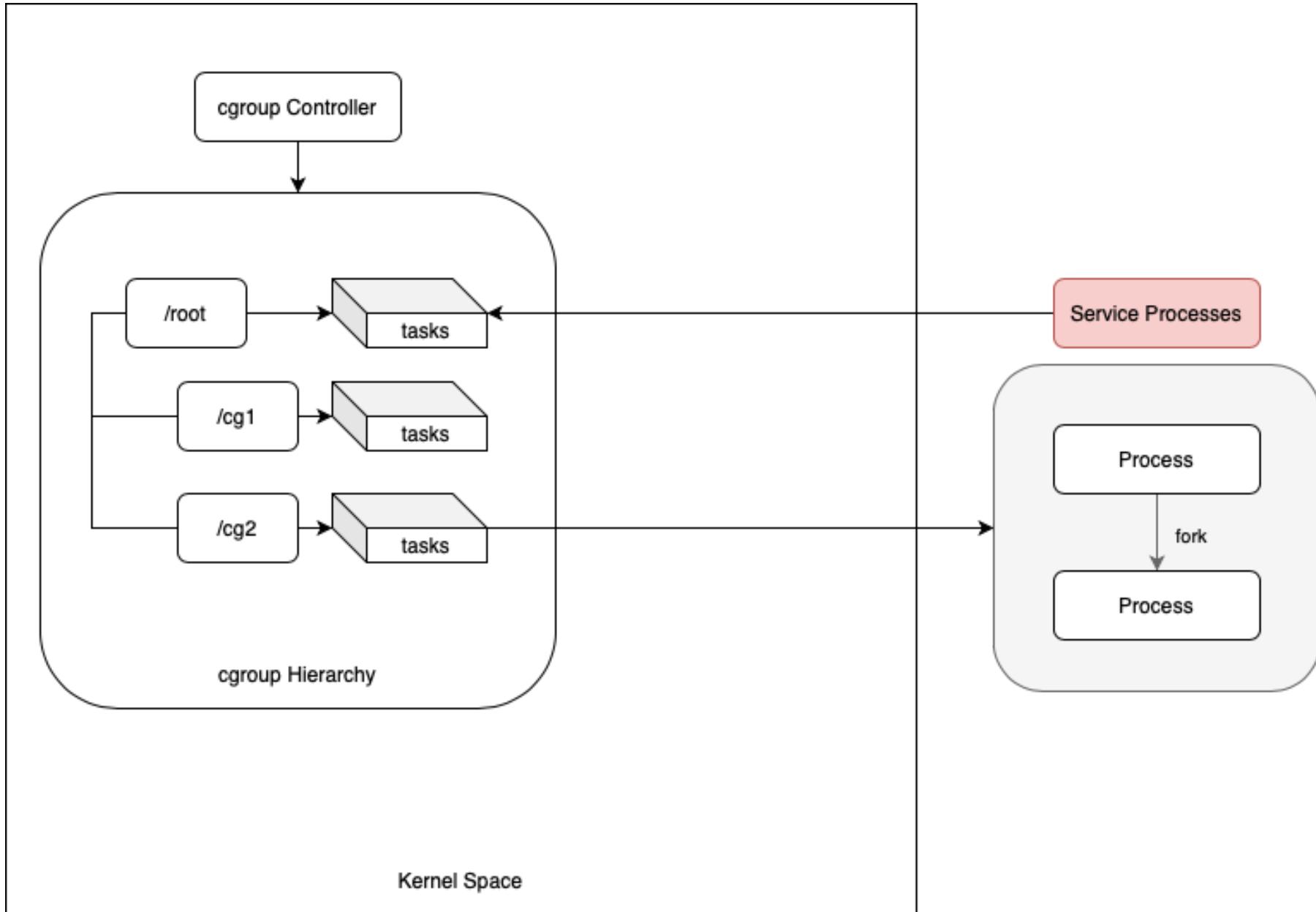
# Result

Time consumption: ~200%

CPU usage: ~50%

## **Part 4: Exploiting Service Processes**

### **System Process Journald**



# Systemd-Journald Service

- Provides a system service to collect system **logging** data, including kernel log messages, system log messages
- Three categories of operations in a container can **force** the journald process to log

# Three categories of operations

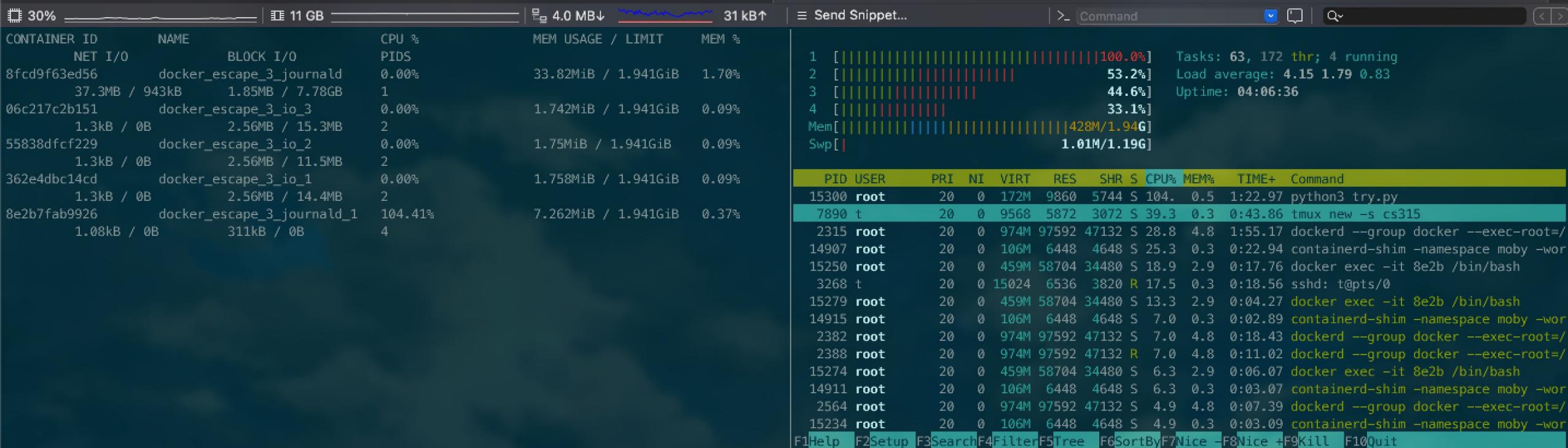
- Switch user (su) command
- Add user/group
- Exception

# Attack Strategy

- System processes of the host are attached to cgroups that are **different** from the processes within containers
- Workloads inside containers can trigger activities for those system processes will **not** be charged to containers' cgroups

# Experiment

- Stop `systemd-journald` service
- Three containers performing file I/O
- Start `systemd-journald` service
- Rerun three I/O containers, and run malicious containers
- Measure time consumotion and CPU usage with/without the malicious system call



```
root@362e4dbc14cd:~# time ./io data1.in data1.out
real    0m47.413s
user    0m21.734s
sys     0m25.563s
root@362e4dbc14cd:~# time ./io data1.in data1.out
real    1m6.999s
user    0m21.727s
sys     0m30.747s
root@362e4dbc14cd:~#
```

```
root@55838dfcf229:~# time ./io data2.in data2.out
real    0m45.638s
user    0m21.559s
sys     0m23.971s
root@55838dfcf229:~# time ./io data2.in data2.out
real    1m11.465s
user    0m22.819s
sys     0m34.295s
root@55838dfcf229:~#
```

```
root@06c217c2b151:~# time ./io data3.in data3.out
real    0m48.033s
user    0m21.958s
sys     0m25.977s
root@06c217c2b151:~# time ./io data3.in data3.out
real    1m23.279s
user    0m23.771s
sys     0m33.087s
root@06c217c2b151:~#
```

```
File "/usr/lib/python3.6/threading.py", line 916
, in _bootstrap_inner
    self.run()
File "try.py", line 14, in run
    run()
File "try.py", line 19, in run
    tot += 1
UnboundLocalError: local variable 'tot' referenced
before assignment

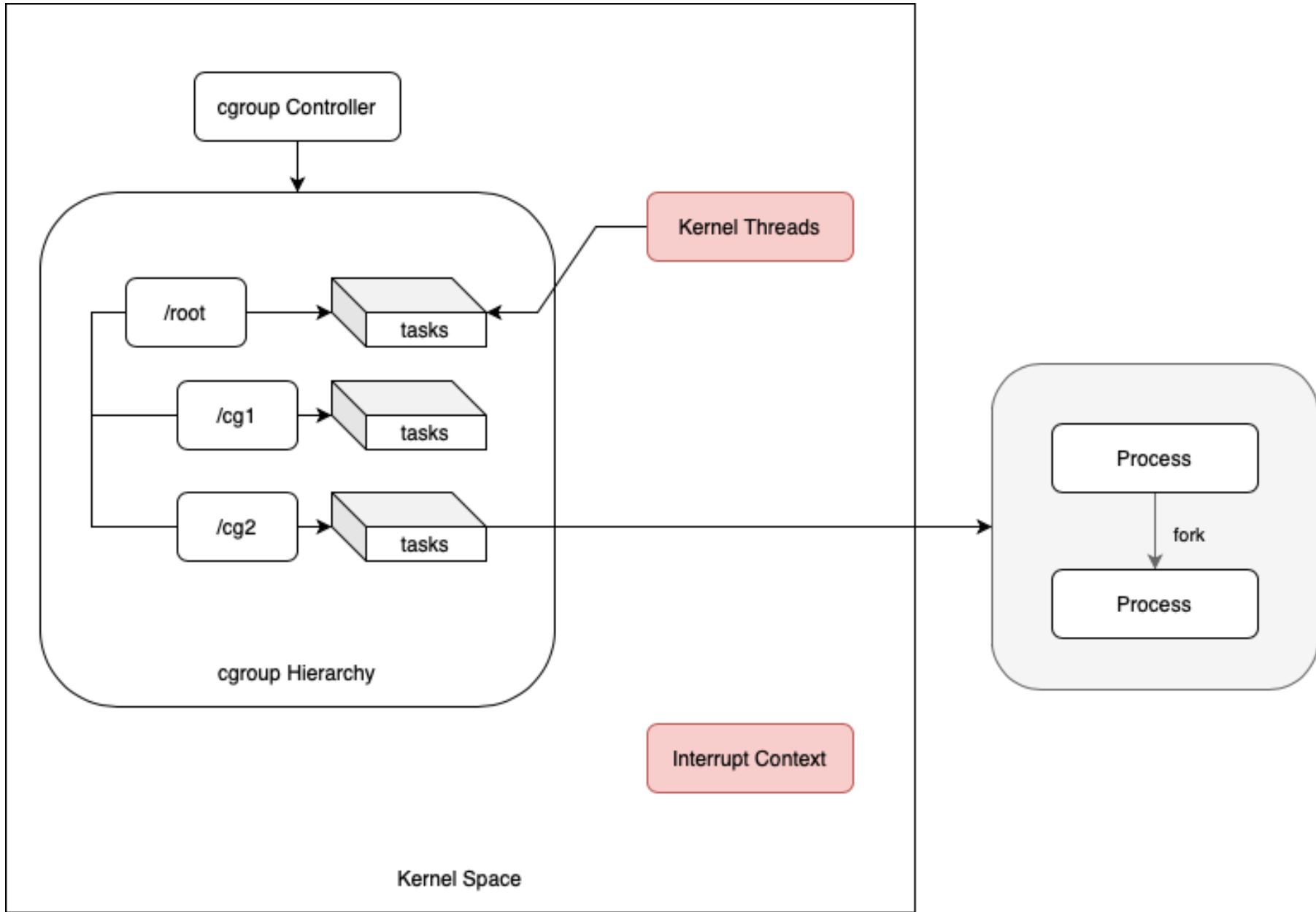
Exception in thread Thread-1:
Traceback (most recent call last):
  File "/usr/lib/python3.6/threading.py", line 916
, in _bootstrap_inner
    self.run()
  File "try.py", line 14, in run
    run()
  File "try.py", line 19, in run
    tot += 1
UnboundLocalError: local variable 'tot' referenced
before assignment
```

# Result

- Time consumption: ~157%
- CPU usage: ~43.6%

# Part 5: Exploiting Kernel threads & Interrupt Context

## Softirq Handling



# Softirq

- Software Interrupt Request
- It is common for hardware interrupt handlers to raise softirqs
- While most hardware interrupts might **not** be directly raised by containers
- Container users are able to manipulate workloads on network interface generating NET softirq
- The handling of those softirqs consumes CPU resources on the process context of kernel thread or interrupt context

## NET softirq

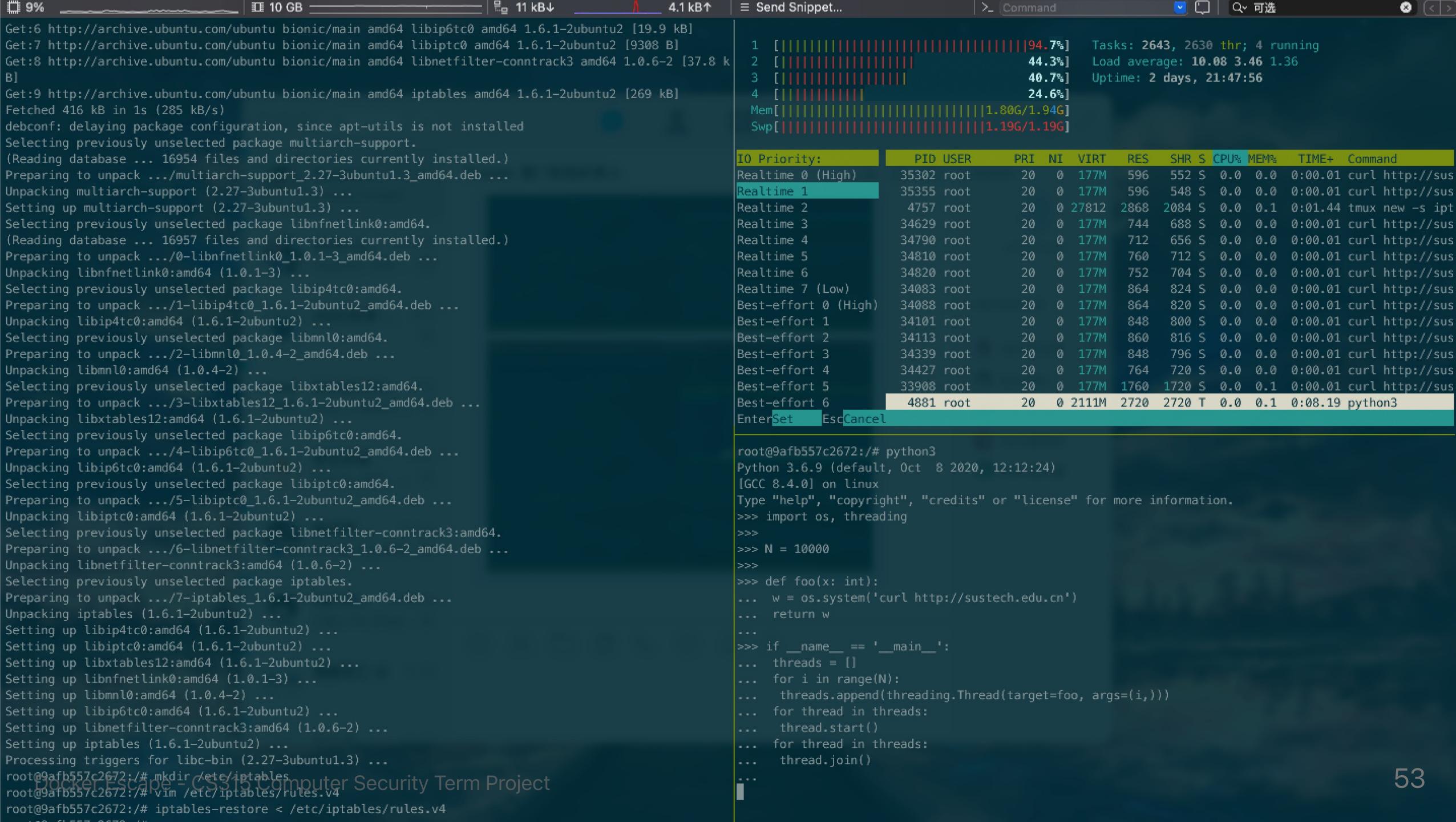
- Interrupt will be raised once the NIC finishes a packet transmission, and softirqs are responsible for moving packets between the driver's buffer and the networking stack.
- **1% overhead for net-working traffic over 1 Gbps**

# Attack Strategy

- The overhead incurred by the networking traffic will be greatly amplified by the firewall system (e.g., iptables) on the server.
- On Linux, Docker relies on configuring iptables rules to provide network isolation for containers.
- Particularly, it might set multiple rules for containers that provide web or networking services.
- The rules **exist** even the container is **stopped**.

# Experiment

- Limit only one cpu core
- Measure CPU usage with the multithreaded script by `htop`



# Result

- CPU usage: ~215.7%

# Docker Escape: Future

- Defence? Think about it.

# Q&A

11712121 胡玉斌/11811902 顾同舟/11812102 黄炜杰/11812103 张文灏/11812104 石文轩