

Kubernetes Security

Kubernetes Security

```

0x004013d0  0x3290/bin/initproc -pxw 0xentry90
0x004013e0  0x8949ed31 0x89485ed1 0x848348e2 0x495
0x004013f0  0x4430c0c7 0x57b00400 0x040430c1 0xc7c
0x00401400  0x004013f0 0x4b6f15ff 0x0ff40020 0x000
0x00401410  0x607218f8 0x2d485988 0x00000000 0x0ef
0x00401420  0x76e58948 0x00000b1b 0x85480000 0x5d10
0x00401430  0x607218f8 0x660e0ff00 0x00841f0f 0x000
0x00401440  0x1f0f335d 0x2e600840 0x00841f0f 0x000x
0x00401450  0x607218b0 0x8b148550 0x607218e1 0xfecp
0x00401460  0xe5894083 0x48f08948 0x483f8ec1 0xd14
0x00401470  0xb815747e 0x00000000 0x74c08548 0x181b
0x00401480  0xf00607f2 0x001f0fe0 0x0f66c35d 0x000
0x00401490  0x5dc13d80 0x75200620 0x89485511 0xff6
0x004014a0  0xc65df0ff 0x205da0e5 0xc3f31001 0x000
0x004014b0  0x0060618f 0x3f834800 0xeb005700 0x001
0x004014c0  0x00000008 0x00854800 0x4855f154 0xd0f
0x004014d0  0xbffa9e5d 0x2e66ff0f 0x00841f0f 0x000
0x004014e0  0xb5554441 0x00000005 0x4548be53 0xfb8
0x004014f0  0x8348f13f 0x8b48800c 0x205d732d 0xb84x4
0x00401500  0x00282504 0x89480000 0x31827440 0xfbf

```



Agenda

- Basic concepts
- Architecture
- Networking
- Direct access & security
- CVE-2018-1002105 explained





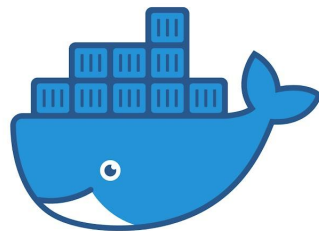
Why me

- I am guilty of setting-up a Kubernetes cluster on-premise in 2015/16
- (Un)fortunately it is still in-use today
- CoreOS + manually installed systemd services
- Bootstrapped through kubelet (migrated from fleet)
- Currently running k8s 1.14
- N+1 redundancy (one node at a time automatically reboots&upgrades with CLUO)
- **Don't try this at home work** (or if you need to, follow Kelsey Hightower: Kubernetes The Hard Way)
- Simple way: **kops** (for cloud) and **kubeadm** (for on-prem)



Docker

- Layer around Linux namespaces and cgroups
- Virtualisation vs. containerization
- Docker image is a standard
- Pro-tip: **don't just blindly run someone else's code from DockerHub on your machine**
- Running an image with `--privileged` or mounting `/dev`
= giving direct access to your host





DockerHub

Authored by: [Kent Lamb](#)

Unauthorized access to Docker Hub database

Article ID: KB000968

HUB

SECURITY

Issue

On Thursday, April 25th, 2019, we discovered unauthorized access to a single Docker Hub database storing a subset of non-financial user data. Upon discovery, we acted quickly to intervene and secure the site.

We want to update you on what we've learned from our ongoing investigation, including which Hub accounts are impacted, and what actions users should take.

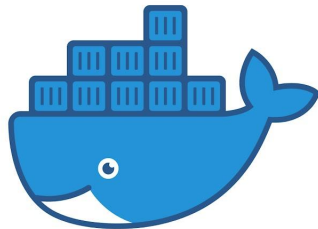
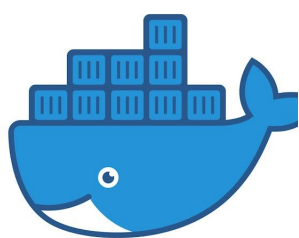


Container manager

How to manage workload (multiple dockers on multiple server).

Diverse workload: batch processing, realtime

- **Kubernetes**
- Docker Swarm
- ...
- CoreOS fleet -> Kubernetes
- Rancher -> wrapper (Kubernetes)





Kubernetes

- Means helmsman / pilot
- Open-sourced from Google 2014 (inspired by Borg), written in Go
- De-facto standard API for IaaS
- Helm is the package manager for k8s





Core concepts

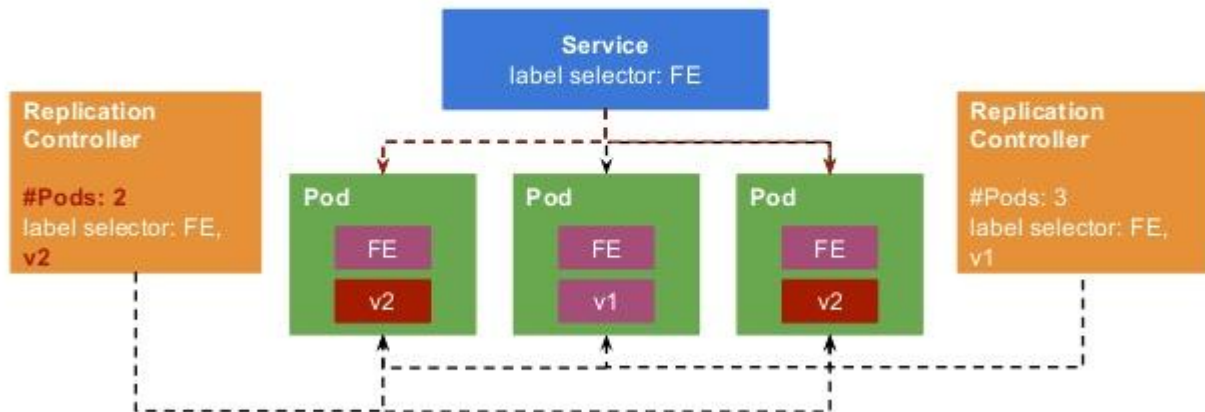
- Label (everything can be tagged with a key/value pair)
- Pod = group of containers running on the same node (~ VM)
- Service = a way to expose network services from a pod
- Replication-controller = makes sure there are N pods
 - deprecated in favor of **ReplicaSets**
 - or declarative configuration - you usually have a **Deployment** object
 - **Deployment** manages **ReplicaSets** which in turn manage **pods**



Tying it all together

Rolling update

 **docker**
TREPTIK.
THE FUTURE OF DEV OPS





Architecture (1/2)

- Administrator
 - **kubectl** (available even as a Windows binary)
- Worker (minion)
 - **kubelet** (docker orchestrator, can be standalone / bootstrap everything else from static config files)
 - **kube-proxy** (enables simple service discovery, more about it in networking section)
- Master
 - **kube-apiserver** (stateless, receives HTTP REST requests)
 - **kube-scheduler** (schedules stuff around the cluster)
 - **kube-controller-manager** (implements the feedback loop - desired state, spec -> actual state, status)



Architecture (2/2)

- Database
 - well actually there is none :), you need **etcd** (distributed reliable key-value store)
 - Raft algorithm and quorum -> having an even amount of servers is stupid
- Only api server will write and read
 - Authentication
 - Authorization
 - Admission controller



Networking (1/2)

- Docker containers are usually connected to a bridge
... and then you expose ports
- Of course there is also host networking
- But when you have lot of containers (on different hosts) this becomes messy
(you want to abstract where a container is running)
- Pods (multiple containers with same IP) are in a flat (virtual) network
- So they can talk among themselves without problems
- However a random IP is assigned to each pod



Networking (2/2)

- You need to use services (which is just a destination NAT trick through kube-proxy for interpod communication - proxy, iptables rules, lvs)
- NodePort service is basically similar to docker & expose port
- In the cloud “service” (with type LoadBalancer) can automatically configure e.g., AWS ELB
- Poor-man’s solution - **MetalLB**





Networking - how does the illusion work?

- Flat net is created with an overlay network (**beware of hybrid cloud set-up and plaintext traffic**)
- We use Flannel (but you could have Calico, Weave, or sth else)
- Called CNI (container network interface)
- VXLAN (L2 encapsulated in L4 UDP port 4789) vs. VLAN (802.1q)



Flannel

Implement IPSEC mode #6

[New issue](#)**Open**

eyakubovich opened this issue on 13 Aug 2014 · 26 comments



eyakubovich commented on 13 Aug 2014

Contributor



For users that require network level encryption, flannel should be able to use IPSEC as the encapsulation protocol.



22



1

eyakubovich added the **enhancement** label on 13 Aug 2014

Assignees

No one assigned

Labels

kind/enhancement

Projects

None yet



Kubelet

- kube-apiserver talks to kubelet to make it start containers
- “By default kubelet allows anonymous authentication”

```
--anonymous-auth
```

```
Enables anonymous requests to the Kubelet server. Requests that are not rejected by another authentication method are treated as anonymous requests. Anonymous requests have a username of system:anonymous, and a group name of system:unauthenticated. (default true)
```

- “The default authorization mode is AlwaysAllow, which allows all requests”
- When you see 10250/TCP open there is a big chance you can **execute code**



Kubelet - verification

```
$ curl -k https://localhost:10250/pods/ | jq . | more
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
           Dload  Upload  Total   Spent    Left     Speed
100 145k    0 145k    0    0 5397k    0 --:--:-- --:--:-- --:--:-- 5397k
{
  "kind": "PodList",
  "apiVersion": "v1",
  "metadata": {},
  "items": [
    {
      "metadata": {
        "name": "kube-proxy-lju1-mts-coreos4",
        "namespace": "kube-system",
        "selfLink": "/api/v1/namespaces/kube-system/pods/kube-proxy-lju1-mts-coreos4",
        "uid": "fbec2ae2a98e09a0c95bf2c1bcb1f7",
        "creationTimestamp": null,
        "annotations": {
          "kubernetes.io/config.hash": "fbec2ae2a98e09a0c95bf2c1bcb1f7",
          "kubernetes.io/config.seen": "2019-04-27T00:31:37.07508847Z",
          "kubernetes.io/config.source": "file"
        }
      },
      "spec": {
        "volumes": [
          {
            "name": "etc-kubernetes",
```



Exposed database (1/2)

- State is stored in etcd (scan for TCP ports 2379/2380 or legacy 4001/7001)
- Kubernetes heavily relies on JWT, so leaking a service token is (often) enough
- RCE -> update deployment and let controller-manager start rogue pods
- How to protect?
 - etcd3 supports roles, but easiest is just to authenticate with cert
 - encrypt data at rest - from k8s side (EncryptionConfig yaml)



Exposed database (2/2)

```
$ ETCCTL_API=3 etcdctl get --endpoints=http://127.0.0.1:2379 /registry/secrets/fiction/default-token-n6tsp | tail -n +2 | jq . | head -15
{
  "kind": "Secret",
  "apiVersion": "v1",
  "metadata": {
    "name": "default-token-n6tsp",
    "namespace": "fiction",
    "uid": "10c32037-ba4a-11e7-ab84-aa0000eecd2",
    "creationTimestamp": "2017-10-26T12:34:52Z",
    "annotations": {
      "kubernetes.io/service-account.name": "default",
      "kubernetes.io/service-account.uid": "10b2c4dc-ba4a-11e7-ab84-aa0000eecd2"
    }
  },
  "data": {
    "ca.crt": "LS0tLS1CRUdJTiBDRVJUSUZJQ0FURSB0tLS0tCk1JSURNakNDQWhxZ0F3SUJBZ0lKQUs2VVhVc0FpQlpaTUeWR0NTcUdTswIzRFFfQkN3VUFNQLV4RXpBUkJnTlYkQkFNTUNtdDFZbVZ5Ym1WMFpYTXdIaGN0TVRZd0lqRTFNVEV4TURNd1doY05Na1l3TWpFeU1URXhNRE13V2pBVgpNUk13RVFZRFRUUREQXByZFdkbGNtNWxkR1Z6TU1JQklqQU5CZ2txaGtpRzl3MEJBUEUvGQUFPQ0FROEUFNSU1CCKNnS0NBUEVBBeTNjbVlac1Bvc3Iya1NUZERsdldqV3BuWlBVNUd0dEJESlJGR0pmOXVZSWlpd3NEOWLucDRQR1AKZlBUZW9hc2x1MWsvQVixOXR0RHJXZnY5U0lGWE94dThwZHLBc
```



“Quasi” SSRF

- The normal SSRF is still a vector (try `https://first_ip_in_cluster_ip_range`, `https://kubernetes.cluster.local` or sth)
- Additionally: each pod by default contains `/var/run/secrets/kubernetes.io/serviceaccount/`
- So rogue code just read the “token” file
- Contact kube-api with that JWT
- Depending on privileges (in kube-system pod -> game-over)
- You could at least see all other pods, if not execute custom code and/or DoS existing stuff in the cluster



Critical security issue

“In all Kubernetes versions prior to v1.10.11, v1.11.5, and v1.12.3, incorrect handling of error responses to proxied upgrade requests in the kube-apiserver allowed specially crafted requests to establish a connection through the Kubernetes API server to backend servers, then send arbitrary requests over the same connection directly to the backend, authenticated with the Kubernetes API server's TLS credentials used to establish the backend connection.” --CVE-2018-1002105

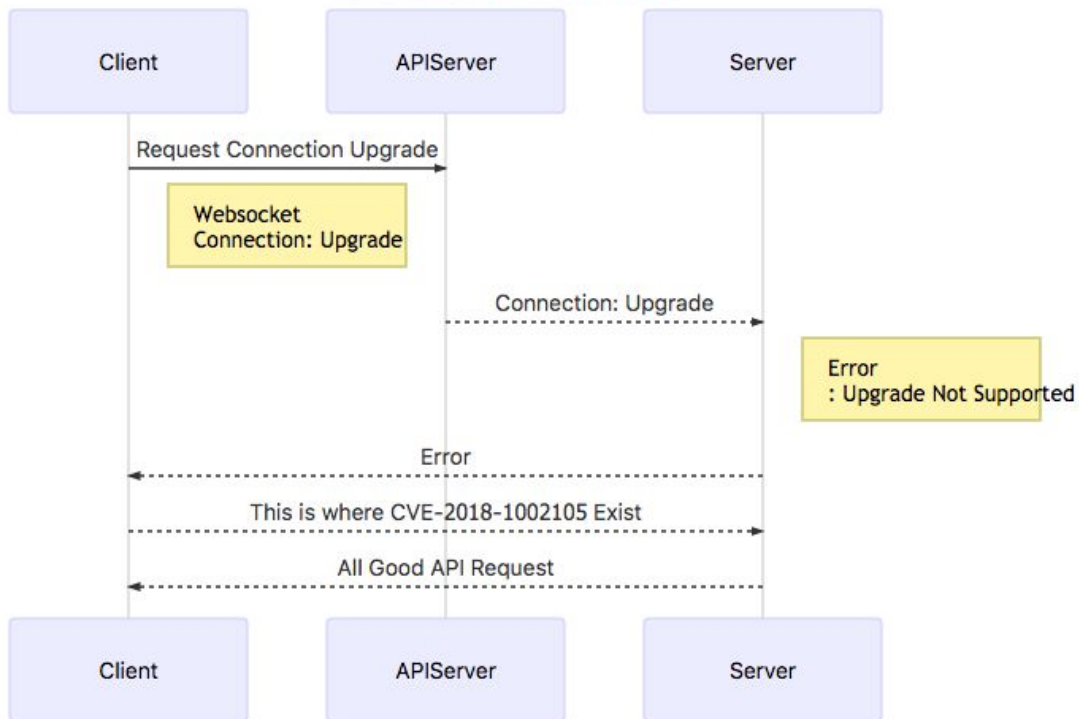
Exec (or portforward) permission to cluster-admin

https://github.com/evict/poc_CVE-2018-1002105



Prettier picture

CVE-2018-1002105





Checklist

- Enable RBAC
- Use TLS between all components (also etcd!) and use IPsec for hybrid set-ups
- Patch your servers (make sure you have redundancy to be able to do it at any time)
- Try to not execute privileged containers
- Beware of DockerHub
- For untrusted code (like CTF challenges) make sure to “overwrite” /run/secrets that gets mounted inside of pod



Questions

?