# Dockercon 2017 Networking Workshop

**Mark Church, Technical Account Manager @ Docker**

**Lorenzo Fontana, Docker Captain**

**Nico Kabar, Solutions Architect @ Docker**

# Agenda

1. Fundamentals & Network Drivers
2. Bridge Driver
3. Overlay Driver
4. MACVLAN Driver
5. Network Services: Service Discovery and Load Balancing
6. Network Design
7. Network Troubleshooting
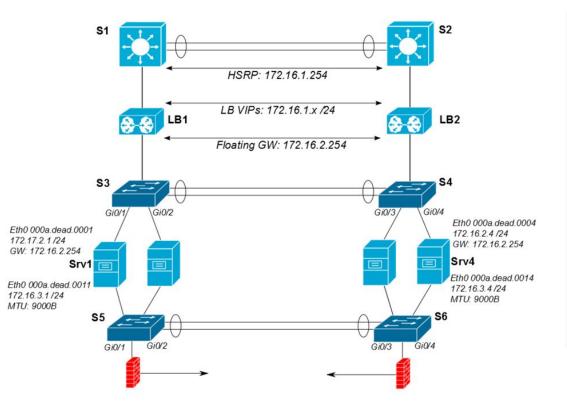8. Deep Dive: Network Namespaces, iptables, and VXLAN

docker

# The Container Network Model (CNM)
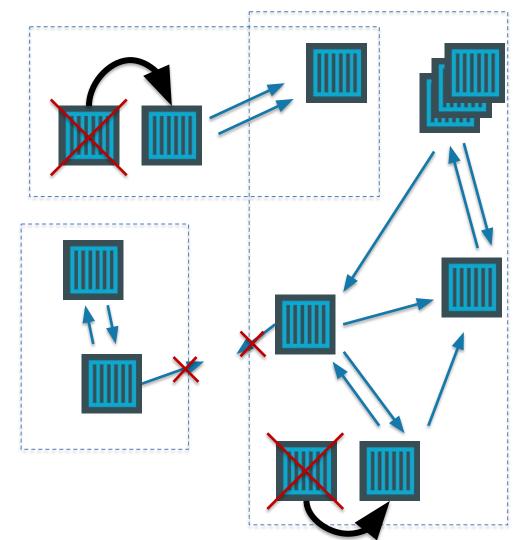
# Networking is hard!

- Distributed in nature

- Many discrete components that are managed and configured differently

- Services that need to be deployed uniformly across all of these discrete components

# Enter containers …

- 100s or 1000s of containers per host

- Containers that exist for minutes or months

- Microservices distributed across many more hosts (>>> E-W traffic)

## … this is worse.

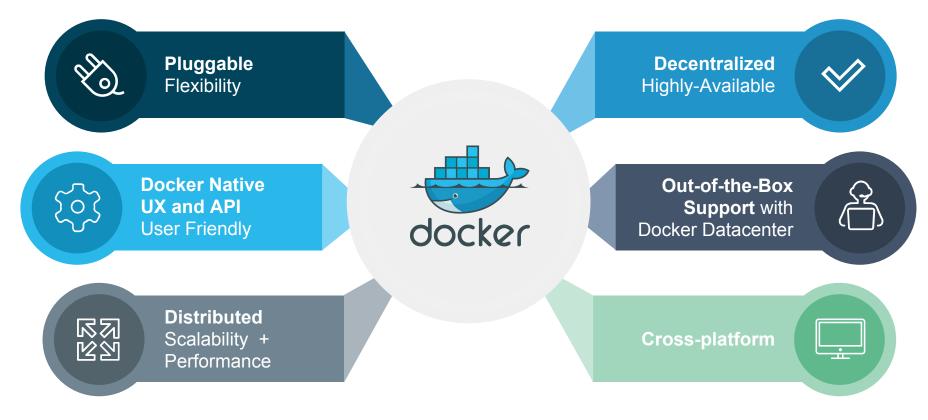# Docker Networking Design Philosophy

**Put Users First**

. . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

Developers and
Operations

**Plugin API Design**

. . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

Batteries included
but removable

# Docker Networking Goals

**Pluggable**
Flexibility

**Docker Native UX and API**
User Friendly

**Distributed**
Scalability +
Performance

**Decentralized**
Highly-Available

**Out-of-the-Box Support** with
Docker Datacenter

**Cross-platform**

# Container Network Model (CNM)

Sandbox

Endpoint
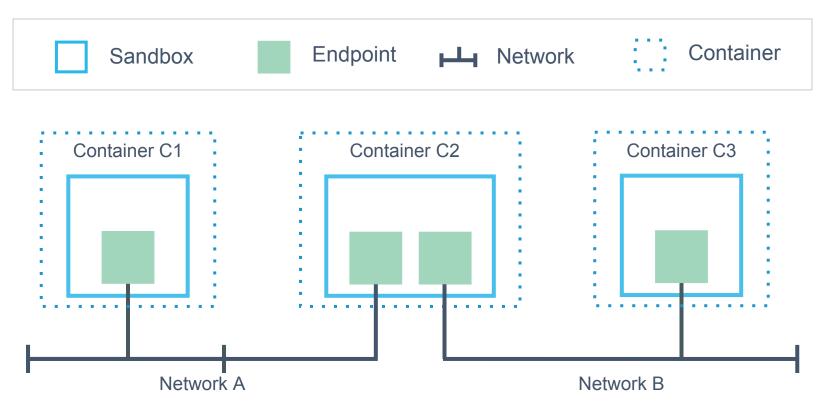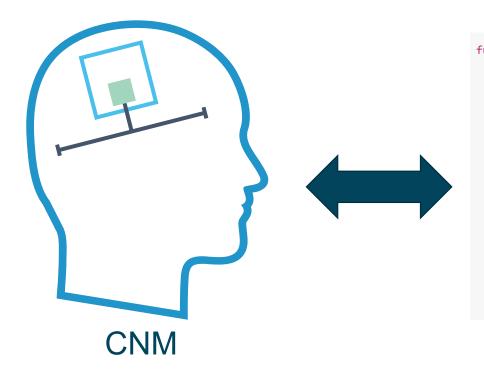
Network

# Containers and the CNM

# What is Libnetwork?

Libnetwork is Docker's native implementation of the CNM



```go
func main() {
    if reexec.Init() {
        return
    }

    // Select and configure the network driver
    networkType := "bridge"

    // Create a new controller instance
    driverOptions := options.Generic{}
    genericOption := make(map[string]interface{})
    genericOption[netlabel.GenericData] = driverOptions
    controller, err := libnetwork.New(config.OptionDriver
    if err != nil {
        log.Fatalf("libnetwork.New: %s", err)
    }
```

CNM

Libnetwork

# What is Libnetwork?

Docker's native implementation of the CNM

Provides built-in service discovery and load balancing

Library containing everything needed to create and manage container networks

Provides a consistent versioned API

Pluggable model (native and remote/3rd party drivers)

Multi-platform, written in Go, open source

docker

# Libnetwork and Drivers

Libnetwork has a
pluggable driver interface

Drivers are used to implement
different networking technologies

**Built-in drivers are called
local drivers, and include:**
bridge, host, overlay, MACVLAN

**3rd party drivers are called
remote drivers, and include:**
Calico, Contiv, Kuryr, Weave…

Libnetwork also supports pluggable IPAM drivers

docker

# Show Registered Drivers

```
$ docker info

Containers: 0
 Running: 0
 Paused: 0
 Stopped: 0
Images: 2
<snip>
Plugins:
 Volume: local
 Network: null bridge host overlay
...
```
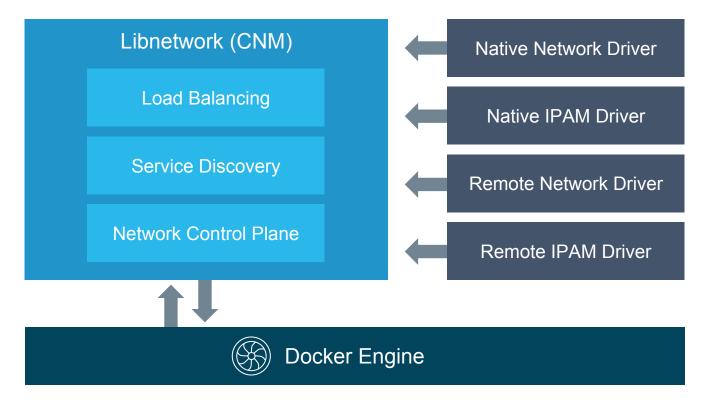
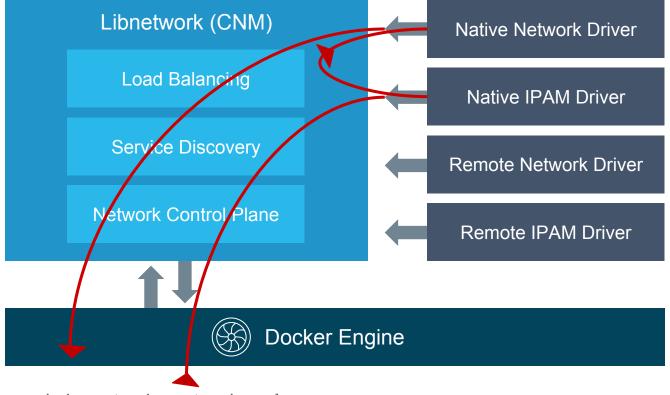# Libnetwork Architecture

# Libnetwork Communication Flow



docker network create -d overlay ov

# Networks and Containers



`docker network create –d <driver> ...`

Defer to Driver

`docker run --network ...`

Libnetwork   Driver   Driver   Engine

# Detailed Overview: Summary

- The CNM is an open-source container networking specification contributed to the community by Docker, Inc.

- The CNM defines sandboxes, endpoints, and networks

- Libnetwork is Docker's implementation of the CNM

- Libnetwork is extensible via pluggable drivers

- Drivers allow Libnetwork to support many network technologies

- Libnetwork is cross-platform and open-source

The CNM and Libnetwork **simplify** container networking and improve **application portability**

docker

# Docker Networking Fundamentals

docker

Libnetwork
(CNM)

- Multihost Networking
- Plugins
- IPAM
- Network UX/API

- Aliases
- DNS Round
  Robin LB

- HRM
- Host-Mode

| 1.7 | 1.8 | 1.9 | 1.10 | 1.11 | 1.12 | 1.13 17.03 |

Service
Discovery

Distributed
DNS

- Secure out-of-the-box
- Distributed KV store
- Load balancing
- Swarm integration
- Built-in routing mesh
- Native multi-host
  overlay
- …

docker

# Docker Networking on Linux

- The Linux kernel has extensive networking capabilities (TCP/IP stack, VXLAN, DNS…)
- Docker networking utilizes many Linux kernel networking features (network namespaces, bridges, iptables, veth pairs…)
- Linux bridges: L2 virtual switches implemented in the kernel
- Network namespaces: Used for isolating container network stacks
- veth pairs: Connect containers to container networks
- iptables: Used for port mapping, load balancing, network isolation…

# Docker Networking *is* Linux (and Windows) Networking

Host

User Space

Kernel

Devices

Docker Engine

Linux Bridge    VXLAN    net namespaces

IPVS    iptables    veth    TCP/IP

eth0    eth1

# Docker Networking on Linux and Windows

## Linux

- Network Namespace

- Linux Bridge

- Virtual Ethernet Devices

- IP Tables

## Windows

- Network Compartments

- VSwitch

- Virtual nics

- Firewall & VFP Rules

# Docker Windows Networking

# Container Network Model (CNM)

Sandbox

Endpoint

Network

# Linux Networking with Containers

- Namespaces are used extensively for container isolation

- Host network namespace is the default namespace

- Additional network namespaces are created to isolate containers from each other

Docker host

Host Network Namespace

Cntnr1    Cntnr2    Cntnr3

eth0    eth1

# Host Mode Data Flow



Docker host 1

Docker host 2

eth0

172.31.1.5

eth0

192.168.1.25

# Demo: Docker Networking Fundamentals

docker

# Lab Section 1

docker

# Bridge Driver

docker

# What is Docker Bridge Networking?

**Single-host networking!**

- Simple to configure and troubleshoot
- Useful for basic test and dev

Docker host

Cntnr1    Cntnr2    Cntnr1

Bridge

# What is Docker Bridge Networking?

- Each container is placed in its own network namespace
- The bridge driver creates a bridge (virtual switch) on a single Docker host

- All containers on this bridge can communicate

- The bridge is a private network restricted to a single Docker host

Docker host

Cntnr1    Cntnr2    Cntnr1

Bridge

docker

# What is Docker Bridge Networking?



Docker host 1

CntnrA    CntnrB

Bridge

Docker host 2

CntnrC    CntnrD

Bridge

Docker host 3

CntnrE    CntnrF

Bridge 1    Bridge 2

Containers on different **bridge** networks cannot communicate

# Bridge Networking in a Bit More Detail

- The bridge created by the bridge driver for the pre-built bridge network is called docker0
- Each container is connected to a bridge network via a veth pair which connects between network namespaces
- Provides single-host networking
- External access requires port mapping



Docker host

| Cntnr1 | Cntnr2 | Cntnr1 |

veth | veth | veth

Bridge

eth0

# Docker Bridge Networking and Port Mapping

Docker host 1

Cntnr1

10.0.0.8  **:80**

Bridge

172.14.3.55  **:8080**

L2/L3 physical network

Host port                   Container port

```
$ docker run -p 8080:80 ...
```

# Bridge Mode Data Flow

Docker host 1

172.17.10.5

172.17.10.6

eth0

eth0

veth

veth

Bridge

eth0

192.168.2.17

Docker host 2

eth0

172.17.8.3

veth

Bridge

eth0

192.168.1.25

# Demo

BRIDGE

docker

# Lab Section 2

docker

# Overlay Driver

# What is Docker Overlay Networking?

The **overlay** driver enables simple and secure **multi-host** networking

Docker host 1

CntnrA          CntnrB

Docker host 2

CntnrC          CntnrD

Docker host 3

CntnrE          CntnrF

Overlay Network

All containers on the **overlay** network can communicate!

# Building an Overlay Network (High level)

Docker host 1

Docker host 2

10.0.0.3

10.0.0.4

Overlay 10.0.0.0/24

172.31.1.5

192.168.1.25

# Docker Overlay Networks and VXLAN

- The **overlay** driver uses VXLAN technology to build the network
- A **VXLAN tunnel** is created through the **underlay network(s)**
- At each end of the tunnel is a VXLAN tunnel end point (**VTEP**)
- The **VTEP** performs encapsulation and de-encapsulation
- The **VTEP** exists in the Docker Host's network namespace

Docker host 1      Docker host 2

VTEP — VXLAN Tunnel — VTEP

172.31.1.5      192.168.1.25

Layer 3 transport
(underlay networks)

# Building an Overlay Network (more detailed)

Docker host 1

Docker host 2

veth

C1: 10.0.0.3

Br0

Network namespace

VTEP
:4789/udp

VXLAN Tunnel

172.31.1.5

veth

C2: 10.0.0.4

B 0

Network namespace

VTEP
:4789/udp

192.168.1.25

Layer 3 transport
(underlay networks)

# Overlay Networking Ports

Docker host 1

Docker host 2

10.0.0.3

10.0.0.4

Management Plane (TCP 2377) - Cluster control

Data Plane (UDP 4789) - Application traffic (VXLAN)

Control Plane (TCP/UDP 7946) - Network control

172.31.1.5

192.168.1.25

# Overlay Network Encryption with IPSec

Docker host 1

verh

C1: 10.0.0.3

B-0

Network
namespace

VTEP
:4789/udp

VXLAN Tunnel

172.31.1.5

Docker host 2

verh

C2: 10.0.0.4

H-0

Network
namespace

VTEP
:4789/udp

192.168.1.25

**IPsec Tunnel**

Layer 3 transport
(underlay networks)

# Overlay Networking Under the Hood

- Virtual eXtensible LAN **(VXLAN)** is the **data transport** (RFC7348)
- Creates a new L2 network over an L3 transport network
- Point-to-Multi-Point tunnels
- VXLAN Network ID (**VNID**) is used to map frames to VLANs
- Uses Proxy ARP
- Invisible to the container
- The **docker_gwbridge** virtual switch per host for default route
- Leverages the distributed KV store created by Swarm
- Control plane is encrypted by default
- Data plane can be encrypted if desired

# Demo

OVERLAY

docker

# MACVLAN Driver

docker

# What is MACVLAN?

- A way to attach containers to existing networks and VLANs
- Ideal for apps that are not ready to be fully containerized
- Uses the well known MACVLAN Linux network type

Docker host 1

Cntnr1          Cntnr2

10.0.0.8        10.0.0.9

eth0:
10.0.0.40

V               P

10.0.0.68       10.0.0.25

L2/L3 physical underlay (10.0.0.0/24)

# What is MACVLAN?

A way to connect containers to virtual and physical machines on existing networks and VLANs

**Parent interface** has to be connected to physical underlay

**Sub-interfaces** used to trunk 802.1Q VLANs

Each container gets its own **MAC** and **IP** on the underlay network

Each container is visible on the physical underlay network

Gives containers direct access to the underlay network **without port mapping** and without a **Linux bridge**

Requires **promiscuous mode**

# What is MACVLAN?

Docker host 1

Cntnr1

10.0.0.18

Cntnr2

10.0.0.19

eth0:
10.0.0.30

Docker host 2

Cntnr3

10.0.0.10

Cntnr4

10.0.0.11

eth0:
10.0.0.40

Docker host 3

Cntnr5

10.0.0.91

Cntnr6

10.0.0.92

eth0:
10.0.0.50

V

10.0.0.68

P

10.0.0.25

L2/L3 physical underlay (10.0.0.0/24)

Promiscuous mode

docker

# What is MACVLAN?

Cntnr2
10.0.0.18

Cntnr2
10.0.0.19

Cntnr3
10.0.0.10

Cntnr4
10.0.0.11

Cntnr5
10.0.0.91

Cntnr6
10.0.0.92

V
10.0.0.68

P
10.0.0.25

L2/L3 physical underlay (10.0.0.0/24)

docker

# MACVLAN and Sub-interfaces

- MACVLAN uses **sub-interfaces** to process 802.1Q VLAN tags.

- In this example, two sub-interfaces are used to enable two separate VLANs

- Yellow lines represent VLAN 10

- Blue lines represent VLAN 20



Docker host

Cntnr2

Cntnr2

eth0: 10.0.0.8

eth0: 10.0.20.3

MACVLAN 10

MACVLAN 20

eth0.**10**

eth0.**20**

eth0

802.1q trunk

L2/L3 physical underlay

VLAN 10
10.0.10.1/24

VLAN 20
10.0.20.1/24

# MACVLAN Summary

- Allow containers to be plumbed into existing VLANs

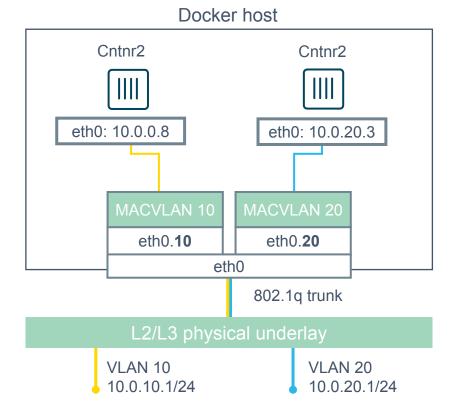- Ideal for integrating containers with existing networks and apps

- High performance (no NAT or Linux bridge…)

- Every container gets its own **MAC** and **routable IP** on the physical underlay

- Uses **sub-interfaces** for 802.1q VLAN tagging

- Requires **promiscuous** mode!

# Demo

MACVLAN

docker

# Use Cases Summary

- The bridge driver provides simple single-host networking
  - Recommended to use another more specific driver such as overlay, **MACVLAN** etc…
- The overlay driver provides native out-of-the-box multi-host networking
- The MACVLAN driver allows containers to participate directly in existing networks and VLANs
  - Requires promiscuous mode
- Docker networking will continue to evolve and add more drivers and networking use-cases

# Docker Network Services

SERVICE REGISTRATION, SERVICE DISCOVERY, AND LOAD BALANCING

docker

# What is Service Discovery?

The ability to discover services within a Swarm

Every **service** registers its name with the Swarm

Every **task** registers its name with the Swarm

Clients can lookup service **names**

Service discovery uses the DNS resolver embedded inside each container and the DNS server inside of each Docker Engine

# Service Discovery in a Bit More Detail

Docker host 1

Docker host 2

task1.myservice          task2.myservice

task3.myservice

"mynet" network (overlay, MACVLAN, user-defined bridge)

```
task1.myservice        10.0.1.19
task2.myservice        10.0.1.20
task3.myservice        10.0.1.21
myservice              10.0.1.18
```

Swarm DNS (service discovery)

# Service Discovery in a Bit More Detail

Docker host 1

task1.myservice          task2.myservice

DNS resolver          DNS resolver
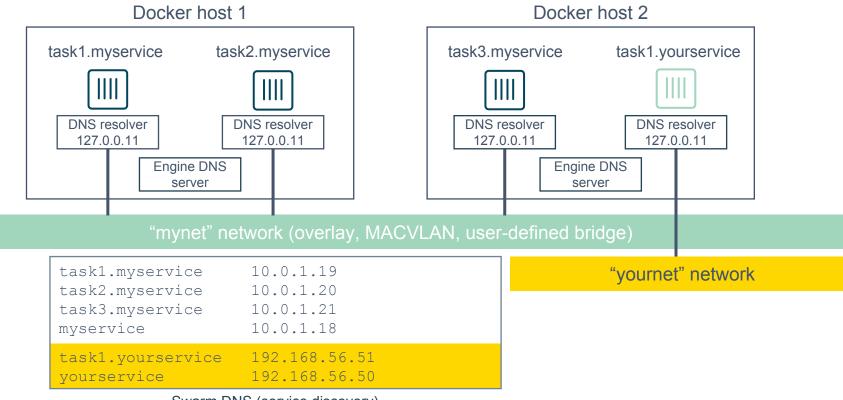127.0.0.11            127.0.0.11

Engine DNS
server

Docker host 2

task3.myservice          task1.yourservice

DNS resolver          DNS resolver
127.0.0.11            127.0.0.11

Engine DNS
server

"mynet" network (overlay, MACVLAN, user-defined bridge)

```
task1.myservice      10.0.1.19
task2.myservice      10.0.1.20
task3.myservice      10.0.1.21
myservice            10.0.1.18
task1.yourservice    192.168.56.51
yourservice          192.168.56.50
```

Swarm DNS (service discovery)

"yournet" network
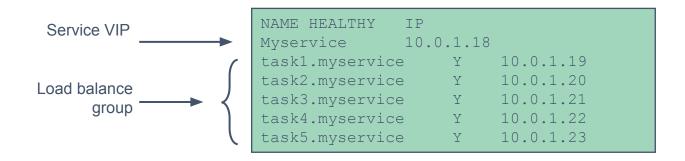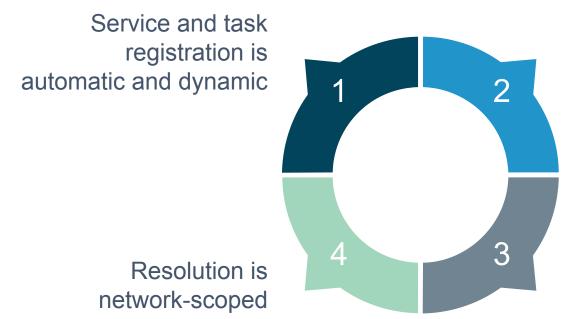
# Service Virtual IP (VIP) Load Balancing

- Every **service** gets a **VIP** when it's created
  - This stays with the service for its entire life
- Lookups against the VIP get load-balanced across all **healthy tasks** in the service
- Behind the scenes it uses Linux kernel **IPVS** to perform transport layer load balancing
- `docker inspect <service>` (shows the service VIP)

Service VIP →

Load balance group →

```
NAME HEALTHY    IP
Myservice         10.0.1.18
task1.myservice      Y     10.0.1.19
task2.myservice      Y     10.0.1.20
task3.myservice      Y     10.0.1.21
task4.myservice      Y     10.0.1.22
task5.myservice      Y     10.0.1.23
```

# Service Discovery Details

**Service and task registration is automatic and dynamic**

**Name-IP-mappings stored in the Swarm KV store**

1

2

4

3

**Resolution is network-scoped**

**Container DNS and Docker Engine DNS used to resolve names**

- Every container runs a local DNS resolver (127.0.0.1:53)
- Every Docker Engine runs a DNS service

# Q & A

docker

# Demo

SERVICE DISCOVERY

docker

# Load Balancing External Requests

ROUTING MESH

# What is the Routing Mesh?

Native load balancing of requests coming from an external source
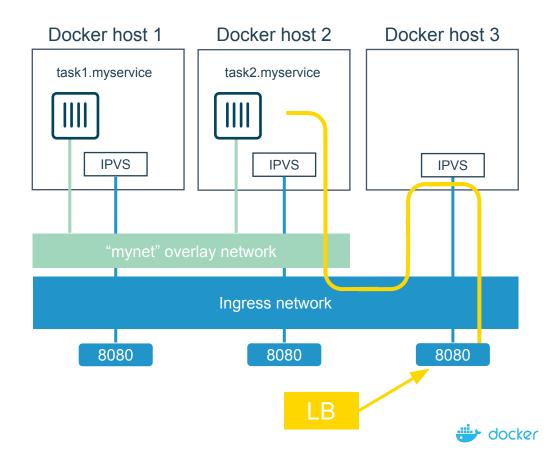
Services get published on a single port across the entire Swarm

A special overlay network called "**Ingress**" is used to forward the requests to a task in the service

Incoming traffic to the published port can be handled by all Swarm nodes

Traffic is internally load balanced as per normal service VIP load balancing

# Routing Mesh Example

1. Three Docker hosts

2. New service with 2 tasks

3. Connected to the **mynet** overlay network

4. Service published on port 8080 swarm-wide

5. External LB sends request to Docker host 3 on port 8080

6. Routing mesh forwards the request to a healthy task using the ingress network

Docker host 1

task1.myservice

IPVS

Docker host 2

task2.myservice

IPVS

Docker host 3

IPVS

"mynet" overlay network

Ingress network

8080

8080

8080

LB

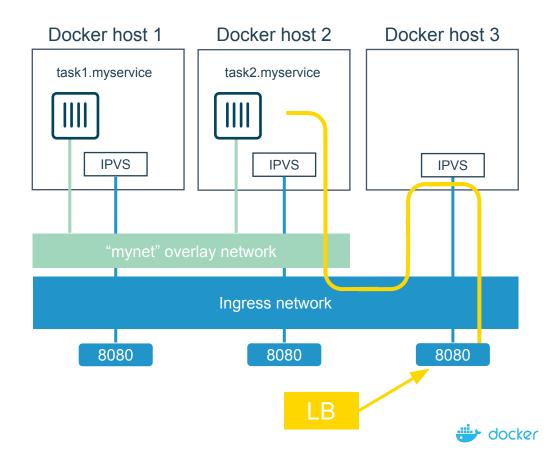# Routing Mesh Example

1. Three Docker hosts
2. New service with 2 tasks
3. Connected to the **mynet** overlay network
4. Service published on port 8080 swarm-wide
5. External LB sends request to Docker host 3 on port 8080
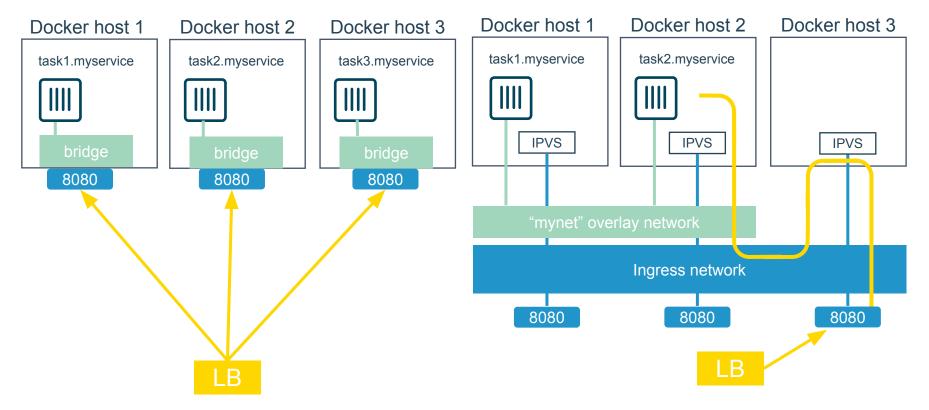6. Routing mesh forwards the request to a healthy task using the ingress network

### Docker host 1

task1.myservice

IPVS

### Docker host 2

task2.myservice

IPVS

### Docker host 3

IPVS

"mynet" overlay network

Ingress network

8080

8080

8080

LB

# Host Mode vs Routing Mesh

# Demo

ROUTING MESH

docker

# HTTP Routing Mesh (HRM) with Docker Datacenter

APPLICATION LAYER LOAD BALANCING (L7)

docker

# What is the HTTP Routing Mesh (HRM)?

Native **application layer (L7)** load balancing of requests coming from an external source

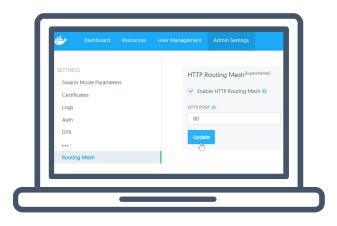Load balances traffic based on hostnames
from HTTP headers

Allows multiple services to be accessed
via the same published port

Requires Docker Enterprise Edition

Builds on top of transport layer routing mesh

docker

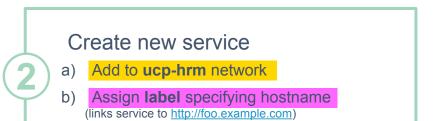# Enabling and Using the HTTP Routing Mesh



```
docker service create -p 8080
\
--network ucp-hrm \
--label
com.docker.ucp.mesh.http=8080=
http://foo.exsample.org
...
```
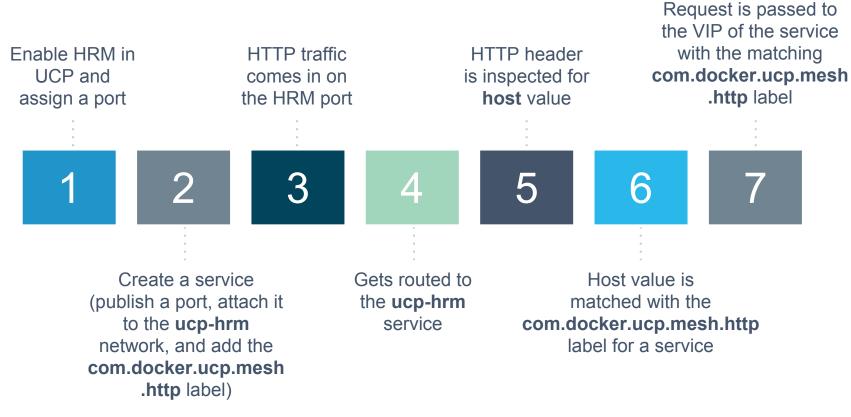
**Enable HTTP routing mesh in UCP**

**1**

a) Creates **ucp-hrm** *network*

b) Creates **ucp-hrm** *service* and exposes it on a port (80 by default)

**Create new service**

**2**

a) Add to **ucp-hrm** network

b) Assign **label** specifying hostname
(links service to http://foo.example.com)

# HTTP Routing Mesh (HRM) Flow

Enable HRM in UCP and assign a port

HTTP traffic comes in on the HRM port

HTTP header is inspected for **host** value

Request is passed to the VIP of the service with the matching **com.docker.ucp.mesh .http** label

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |

Create a service (publish a port, attach it to the **ucp-hrm** network, and add the **com.docker.ucp.mesh .http** label)

Gets routed to the **ucp-hrm** service

Host value is matched with the **com.docker.ucp.mesh.http** label for a service

docker

# HTTP Routing Mesh Example

Docker host 1

user-svc.1
com.docker.ucp.mesh.http=
8080=http://foo.example.com

ucp-hrm.1 :80

Docker host 2

user-svc.2
com.docker.ucp.mesh.http=
8080=http://foo.example.com

ucp-hrm.2 :80

Docker host 3

ucp-hrm.3 :80

**ucp-hrm**
http://foo.example.com
Service: user-svc
VIP: 10.0.1.4

"ucp-hrm" overlay network

Ingress network

ucp-hrm:80

foo.example.com:80

**LB**

docker

# Demo

HRM

docker

# Q & A

docker

# Docker Network Troubleshooting

docker

# Common Network Issues

**Blocked ports, ports required to be open for network mgmt, control, and data plane**

## Iptables issues

Used extensively by Docker Networking, must not be turned off

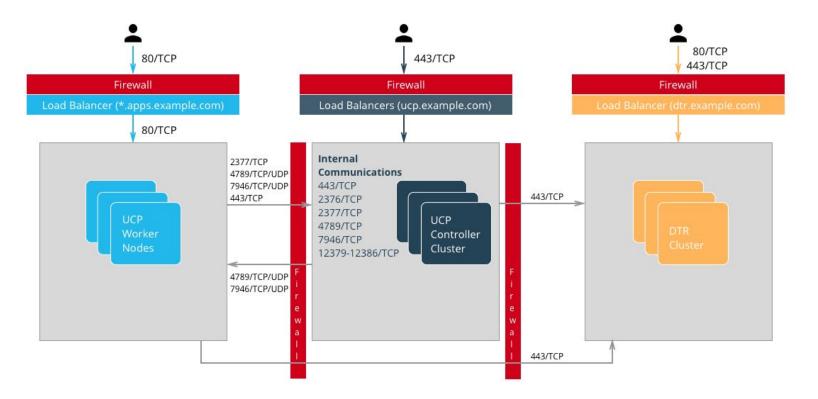List rules with $ iptables -S, $ iptables -S -t nat

**Network state information stale or not being propagated**

Destroy and create networks again with same name

**General connectivity problems**

# Required Ports

# General Connectivity Issues

## Network always gets blamed first :(

Eliminate or prove connectivity first, connectivity can be broken at service discovery or network level

## Service Discovery

Test service name resolution or container name resolution

```
drill <service name> (returns
the service VIP DNS record)

drill tasks.<service name>
(returns all task DNS records)
```

## Network Layer

Test reachability using VIP or container IP

```
task1$ nc -l 5000, task2$
nc <service ip> 5000

ping <container ip>
```

docker

# Netshoot Tool

Has most of the tools you need **in a container** to troubleshoot common networking problems

```
iperf, tcpdump, netstat, iftop, drill, netcat-openbsd, iproute2,
util-linux(nsenter), bridge-utils, iputils, curl, ipvsadmin, ethtool…
```

## Two Uses

Connect it to a specific **network namespace** (such as a container's) to view the network from that container's perspective

Connect it to a **docker network** to test connectivity on that network

# Netshoot Tool

## Connect to a container namespace

```
docker run -it --net container:<container_name> nicolaka/netshoot
```

## Connect to a network

```
docker run -it --net host nicolaka/netshoot
```

Once inside the **netshoot** container, you can use any of the network troubleshooting tools that come with it

# Network Troubleshooting Tools

## Capture all traffic to/from port 999 on eth0 on a myservice container

```
docker run -it --net
container:myservice.1.0qlf1kaka0cq38gojf7wcatoa  nicolaka/netshoot
tcpdump -i eth0 port 9999 -c 1 -Xvv
```

## See all network connections to a specific task in myservice

```
docker run -it --net
container:myservice.1.0qlf1kaka0cq38gojf7wcatoa  nicolaka/netshoot
netstat -taupn
```

# Network Troubleshooting Tools

## Test DNS service discovery from one service to another

```
docker run -it --net
container:myservice.1.bil2mo8inj3r9nyrss1g15qav  nicolaka/netshoot drill
yourservice
```
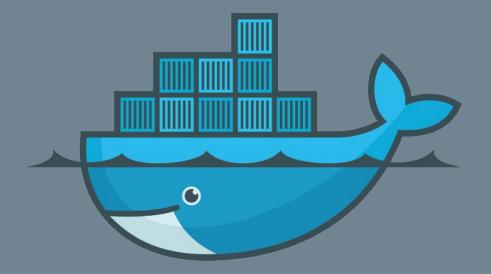
## Show host routing table from inside the netshoot container

```
docker run -it --net host nicolaka/netshoot ip route show
```

# Lab Section 3

THANK YOU