# Untrusting the network
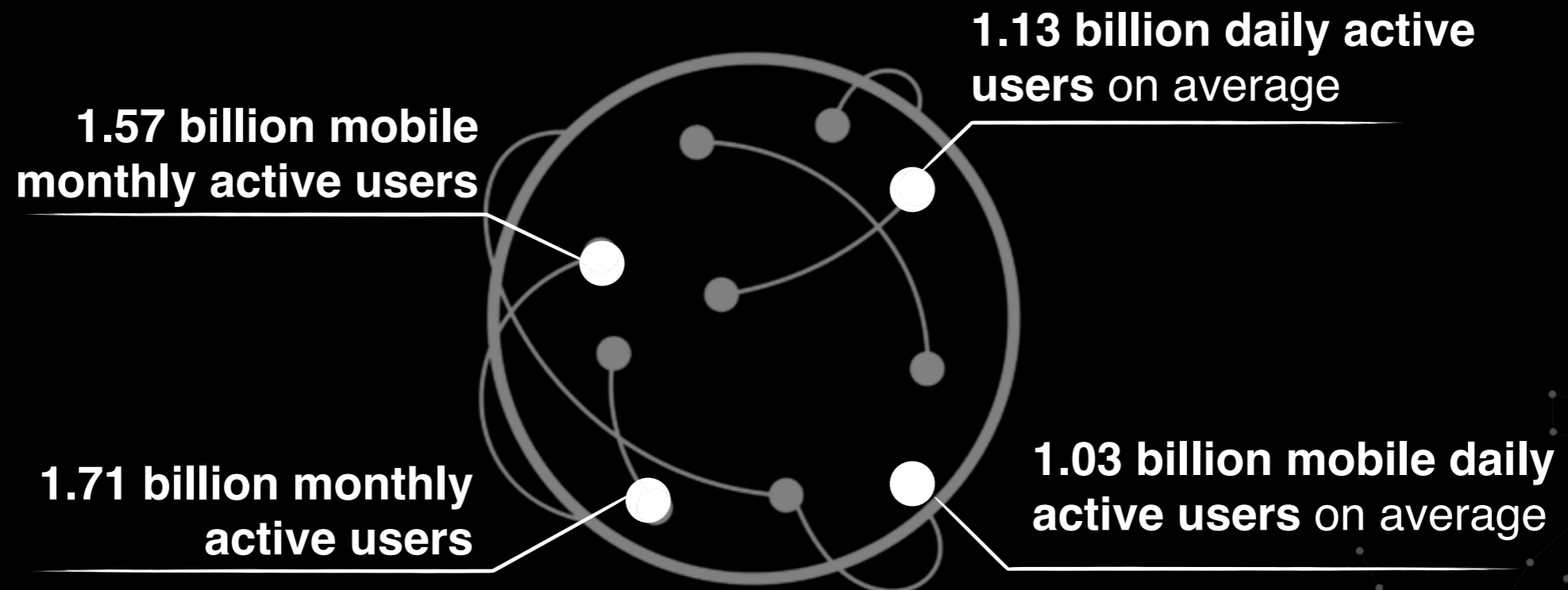
Aijay Adams
Jose Leitao
Production Network Engineers

facebook
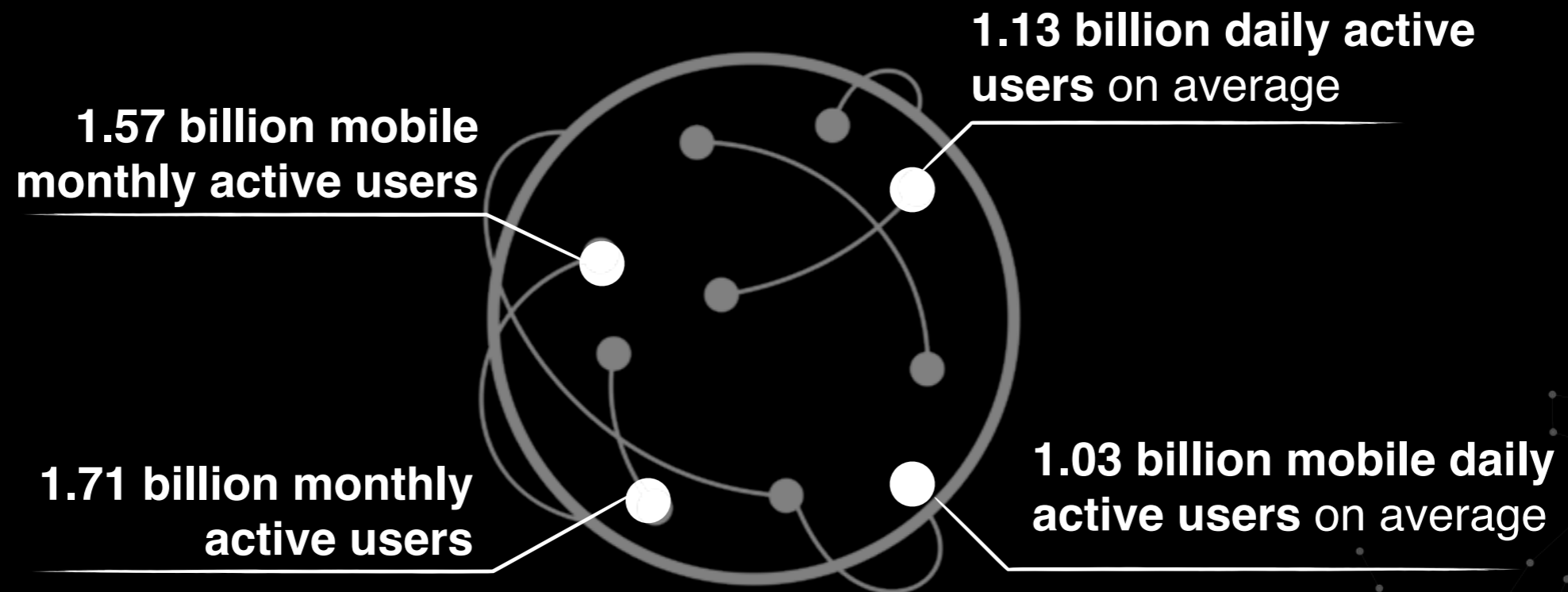
# facebook

as of June 2016
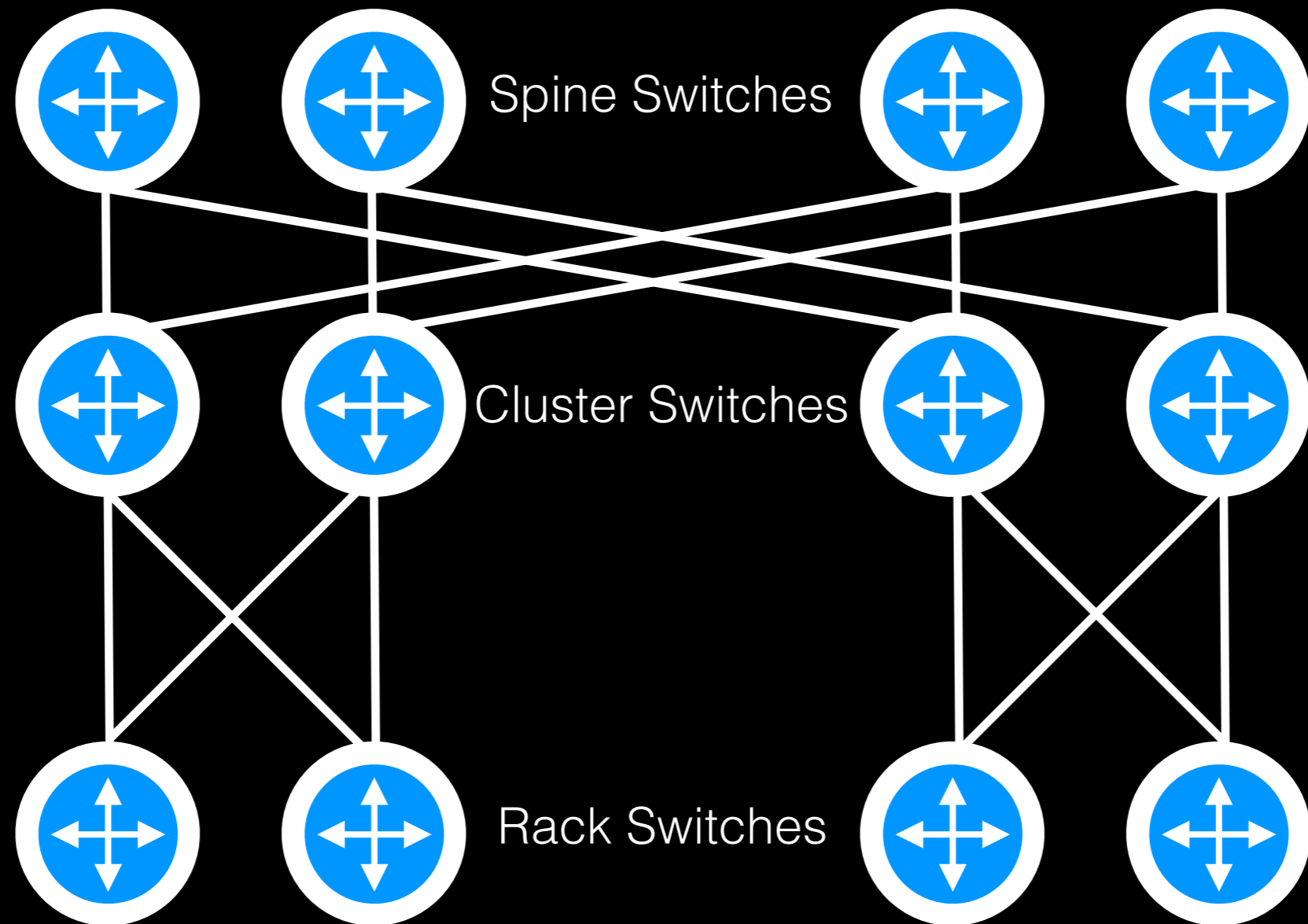
1.13 billion daily active users on average

1.57 billion mobile monthly active users

1.71 billion monthly active users

1.03 billion mobile daily active users on average

# facebook

as of June 2016

**1.13 billion daily active users** on average

**1.57 billion mobile monthly active users**

**1.71 billion monthly active users**

**1.03 billion mobile daily active users** on average

**Approximately 84.5% of our daily active users are outside the US and Canada.**

# Data Center Network



Spine Switches

Cluster Switches

Rack Switches

Wide ECMP, many paths!

# Backbone Network

MPLS Backbone

Data center



Data center

Auto Bandwidth

ECMP over MPLS Tunnels

# Monitoring the Network

?? ?? ?? ?
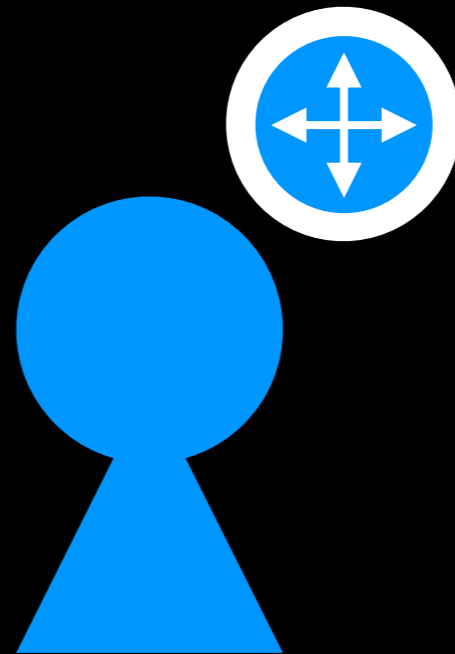
# Monitoring the Network

Counters and Logs

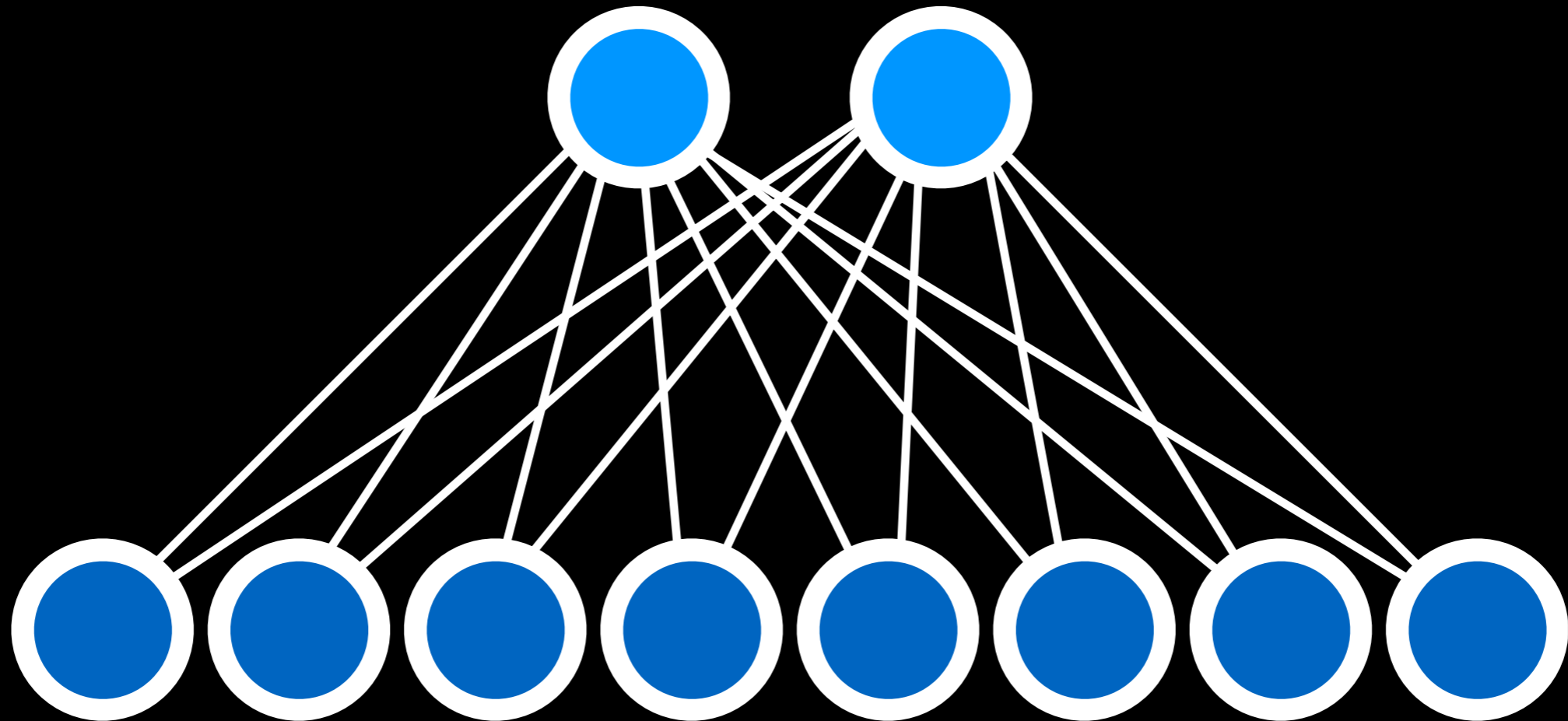# Monitoring the Network



Coworkers

# NetNORAD

The network fault detector

github.com/facebook/UdpPinger

# Ping all the things!



Run pingers on some machines
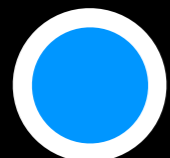Run responders on all machines
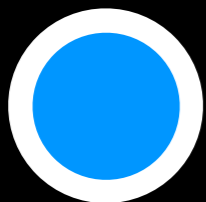Collect and analyse data
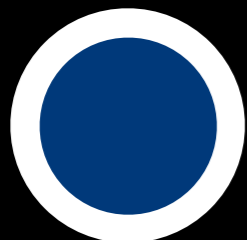
# Evolution

● Run /bin/ping from a python agent

● Raw Sockets, Fast TCP Probes

● Raw Sockets, Fast ICMP Probes
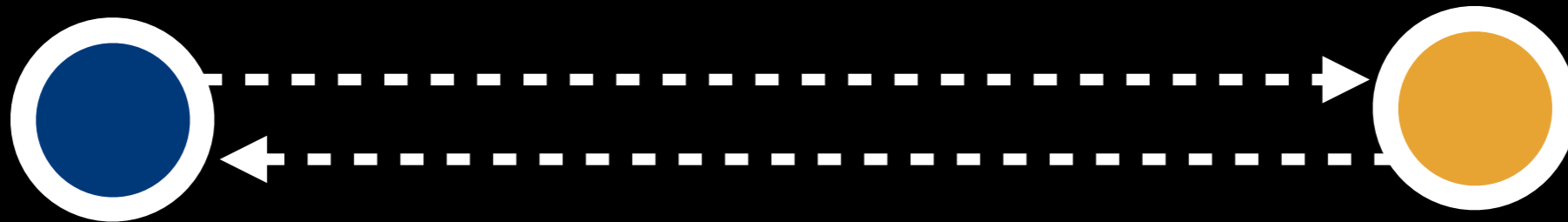
● UDP Probes and Responder

● UDP Probes and Responder + Fast ICMP Probes

# Ping Pong

**Pingers**

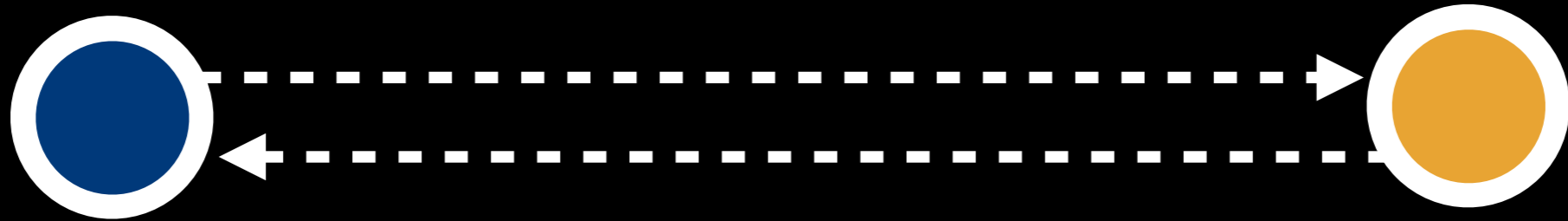**Responders**

- Send UDP and ICMP probes
      to target list
- Timestamp & Log results
- High ping-rate (up to 1Mpps)
- Set DSCP marking

- Receive/Reply to probe
- Timestamp
- Low load: thousands of pps
- Reflect DSCP value back

# Why UDP?

- No TCP RST packets
- Efficient ECMP coverage
- Extensible

**Probe Structure**

Signature

Send Time

Receive Time

Response Time

Traffic Class

# NetNORAD

Ping and Process Data

github.com/facebook/UdpPinger
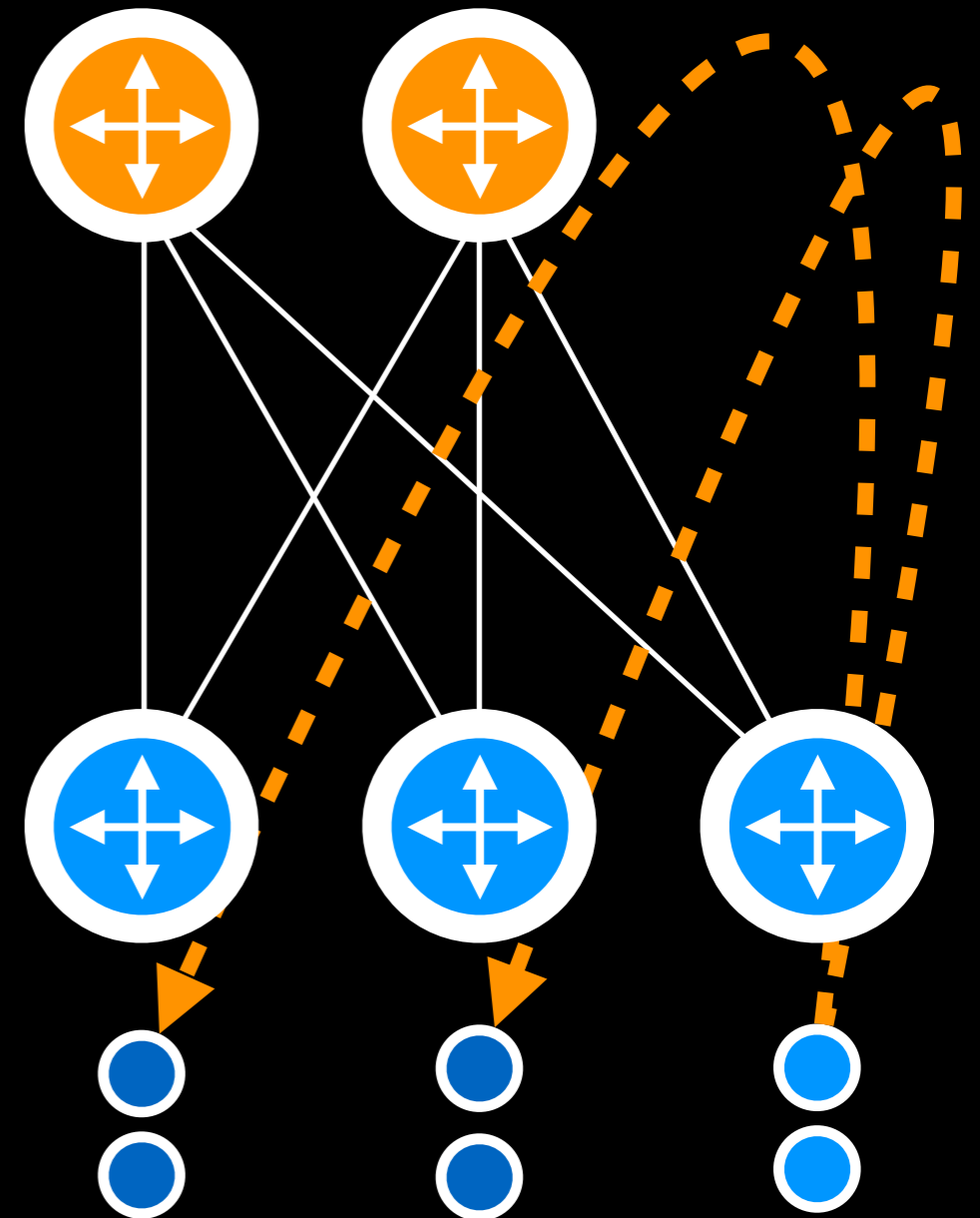
# Challenges

Tens of thousands of targets
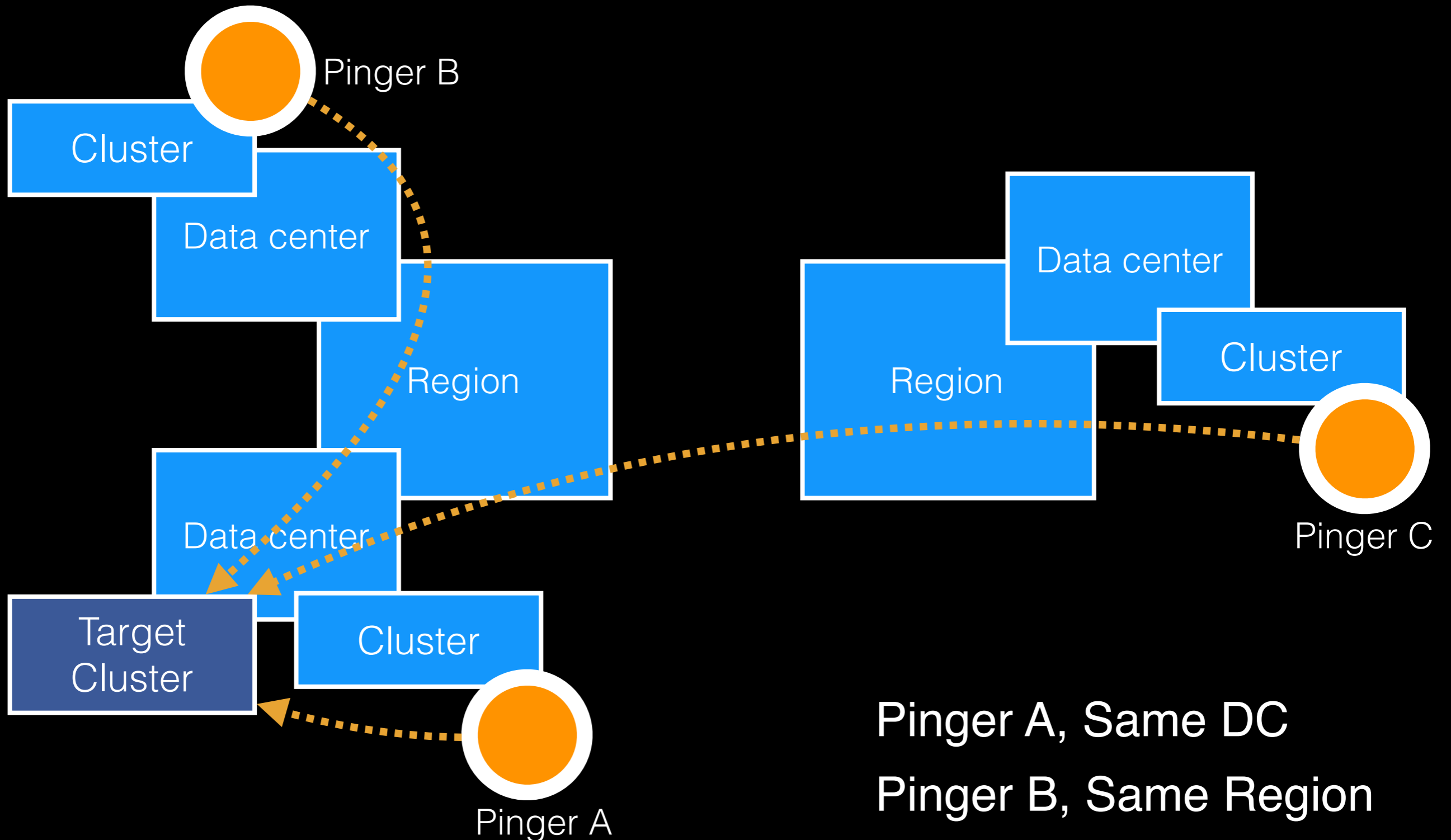
Hundreds of pingers

Lots of data to process

We really do not care about each host…
…The unit of interest is cluster health

# Pinging inside clusters

- Detect issues with rack switches

- Dedicated pingers per cluster

- Probe ALL machines in cluster

- Store time-series per host/rack
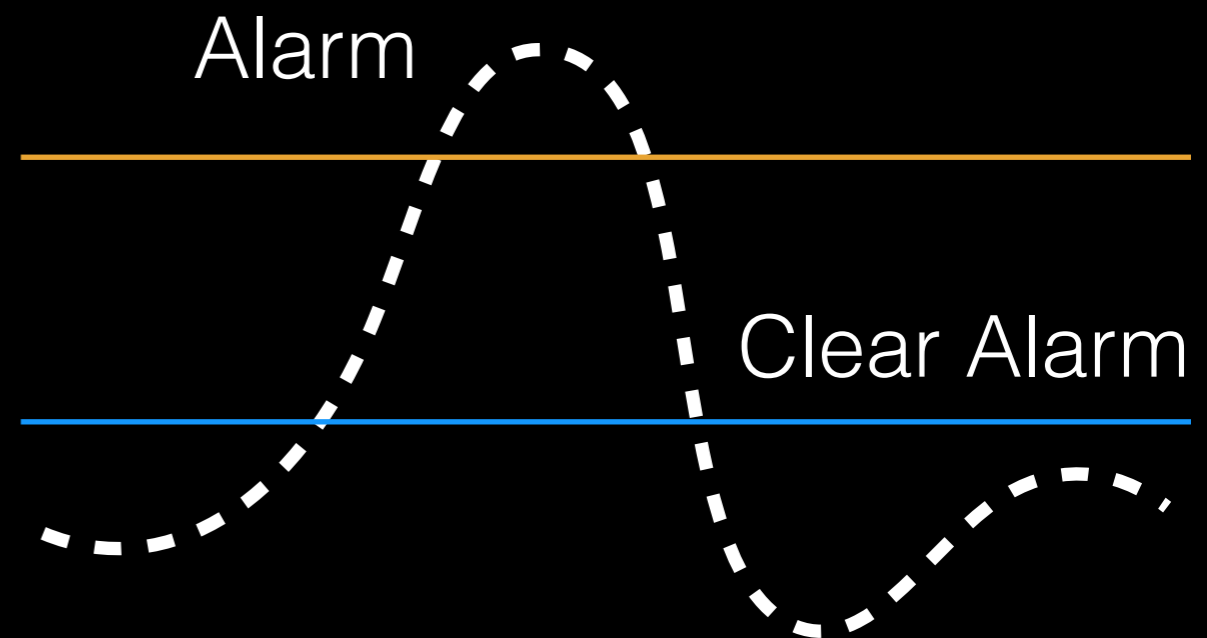
- Lags real-time by 2 minutes

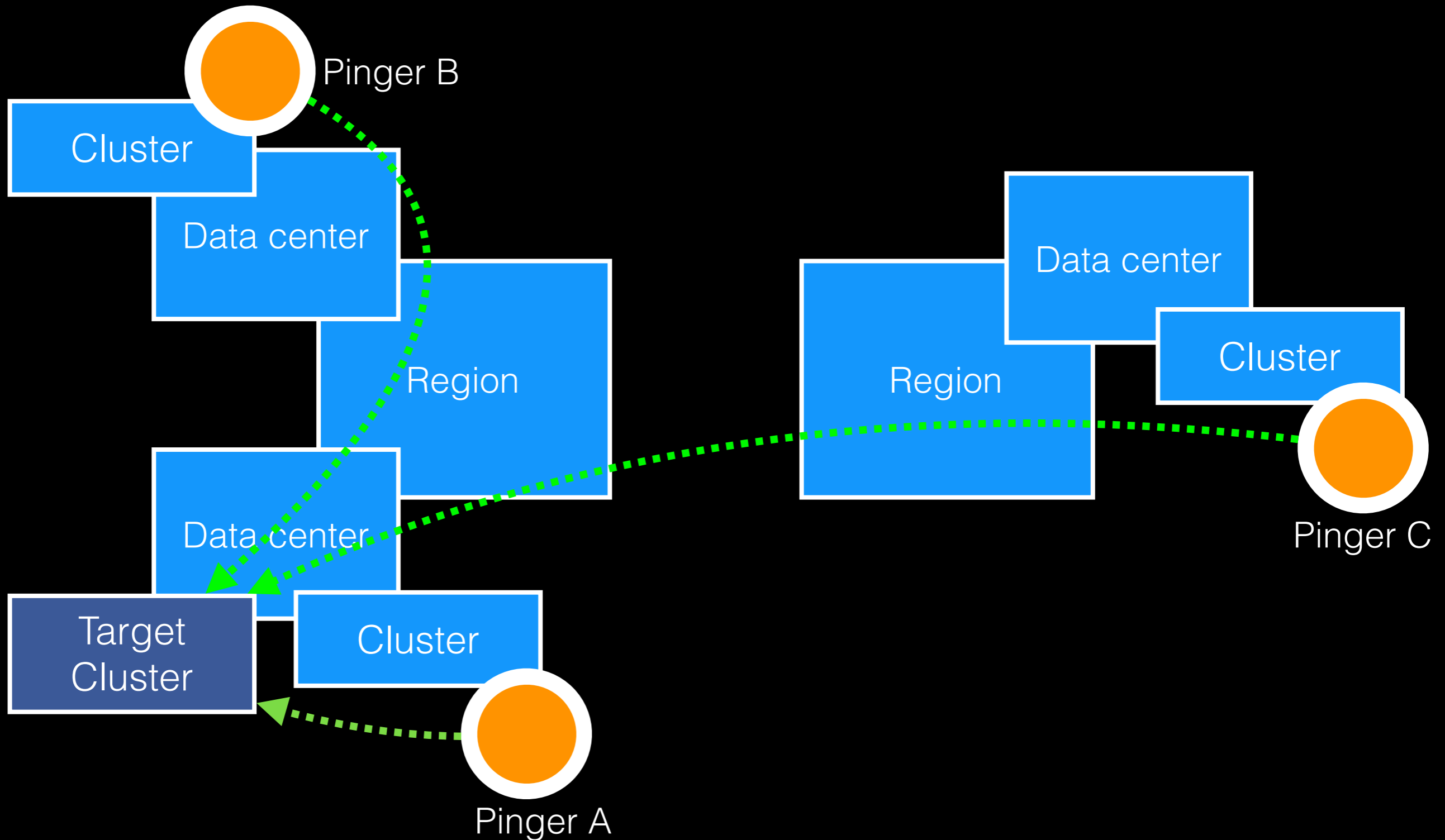# Pinging the clusters



Pinger A, Same DC

Pinger B, Same Region

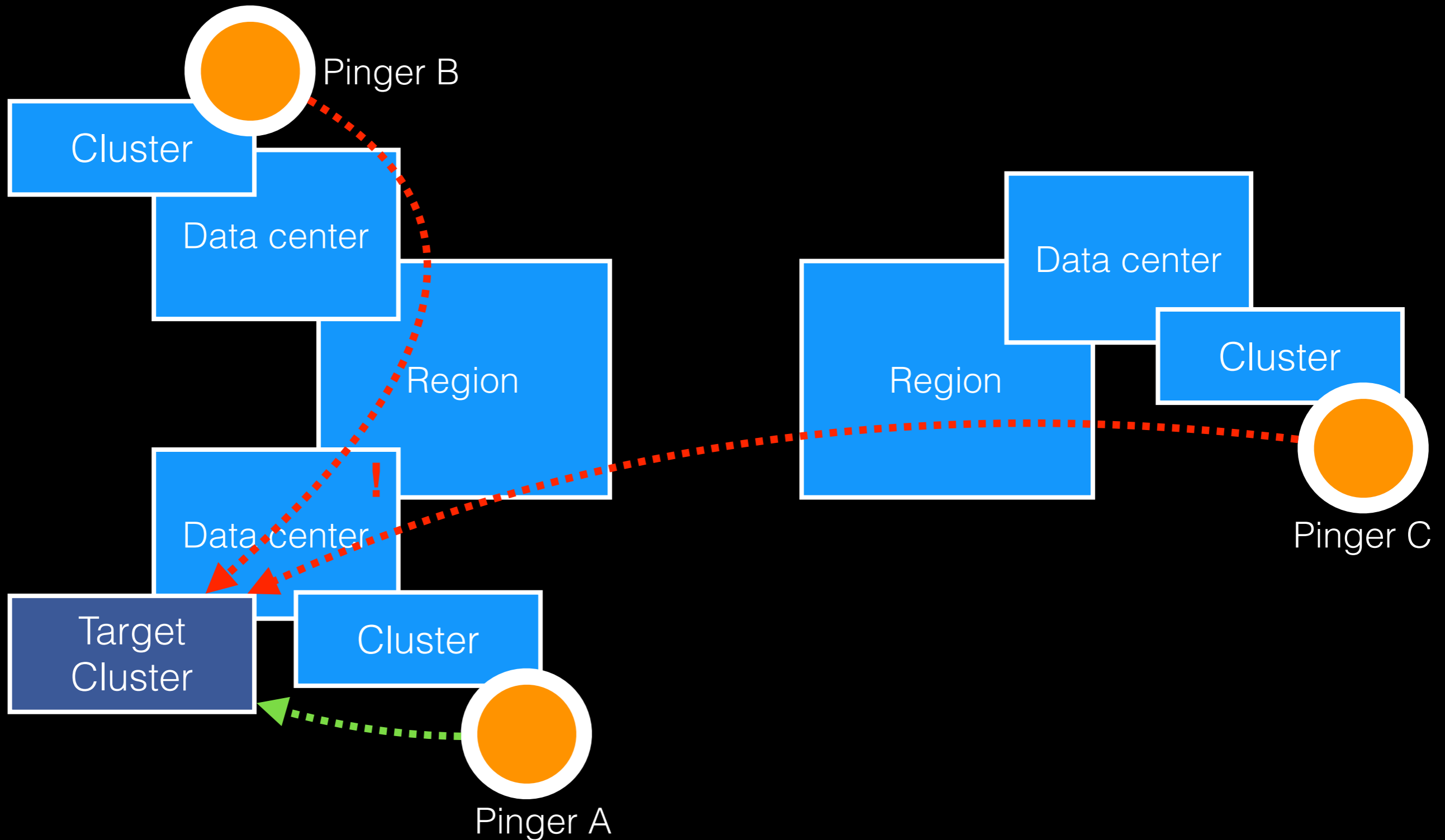Pinger C, Outside of region

# Alarming on Loss

- Build packet loss time-series
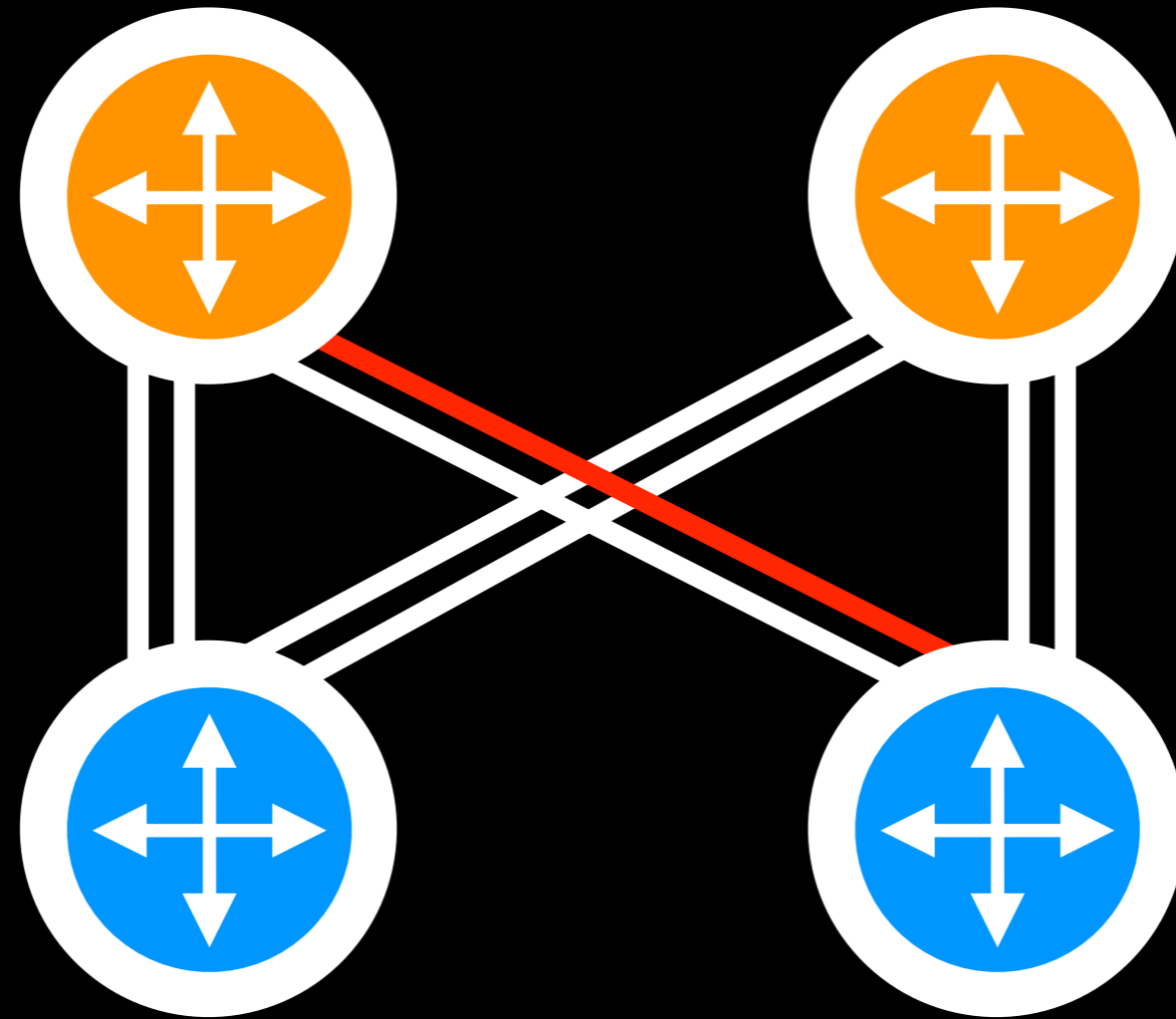
- Track percentiles

- Time to detect loss? 20 Seconds

Alarm

Clear Alarm

# Finding the problem

# Finding the problem

# Finding the problem



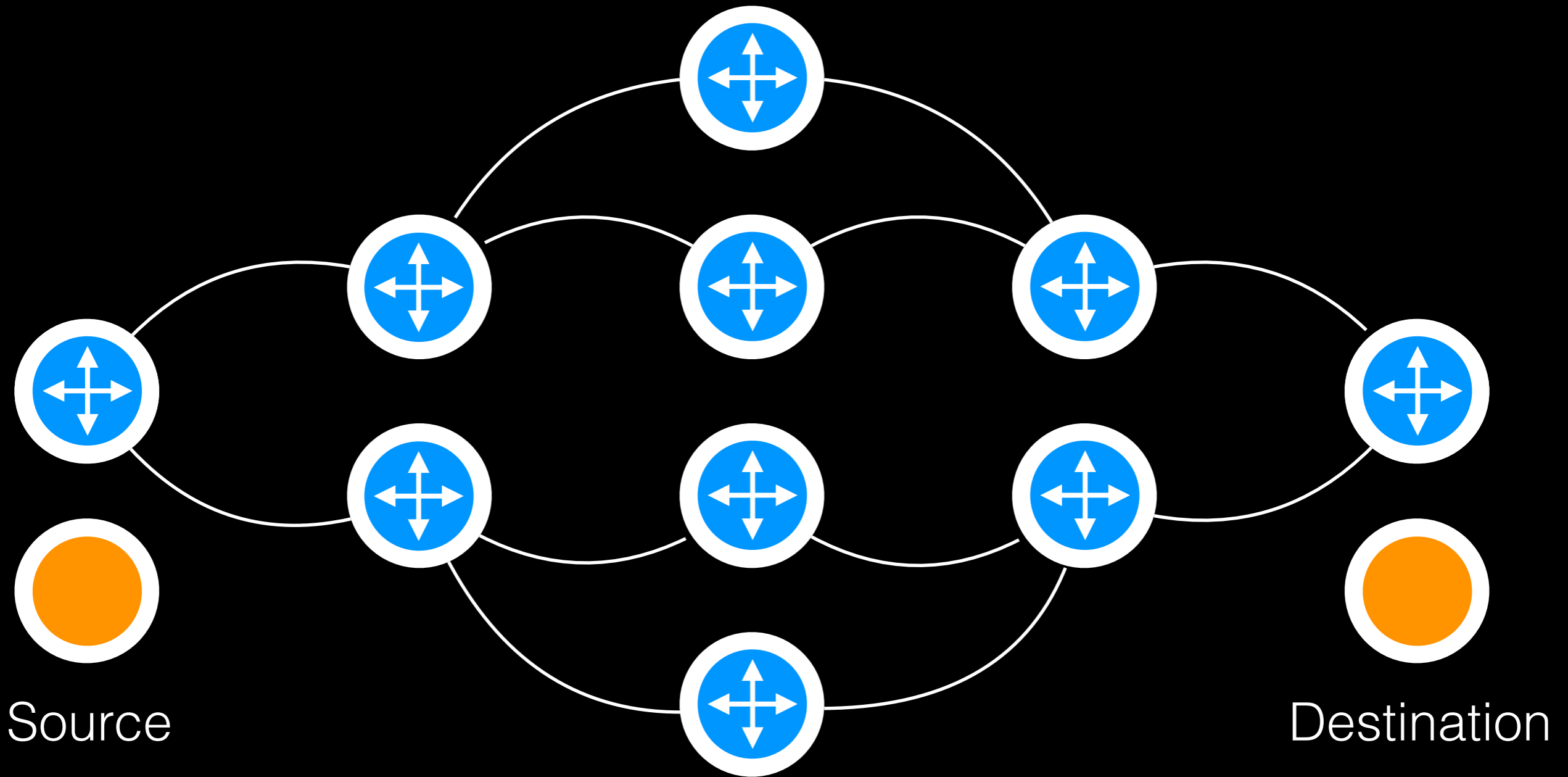Each layer of the network contains many devices
...and many more links

# fbtracert

Isolating Network Faults
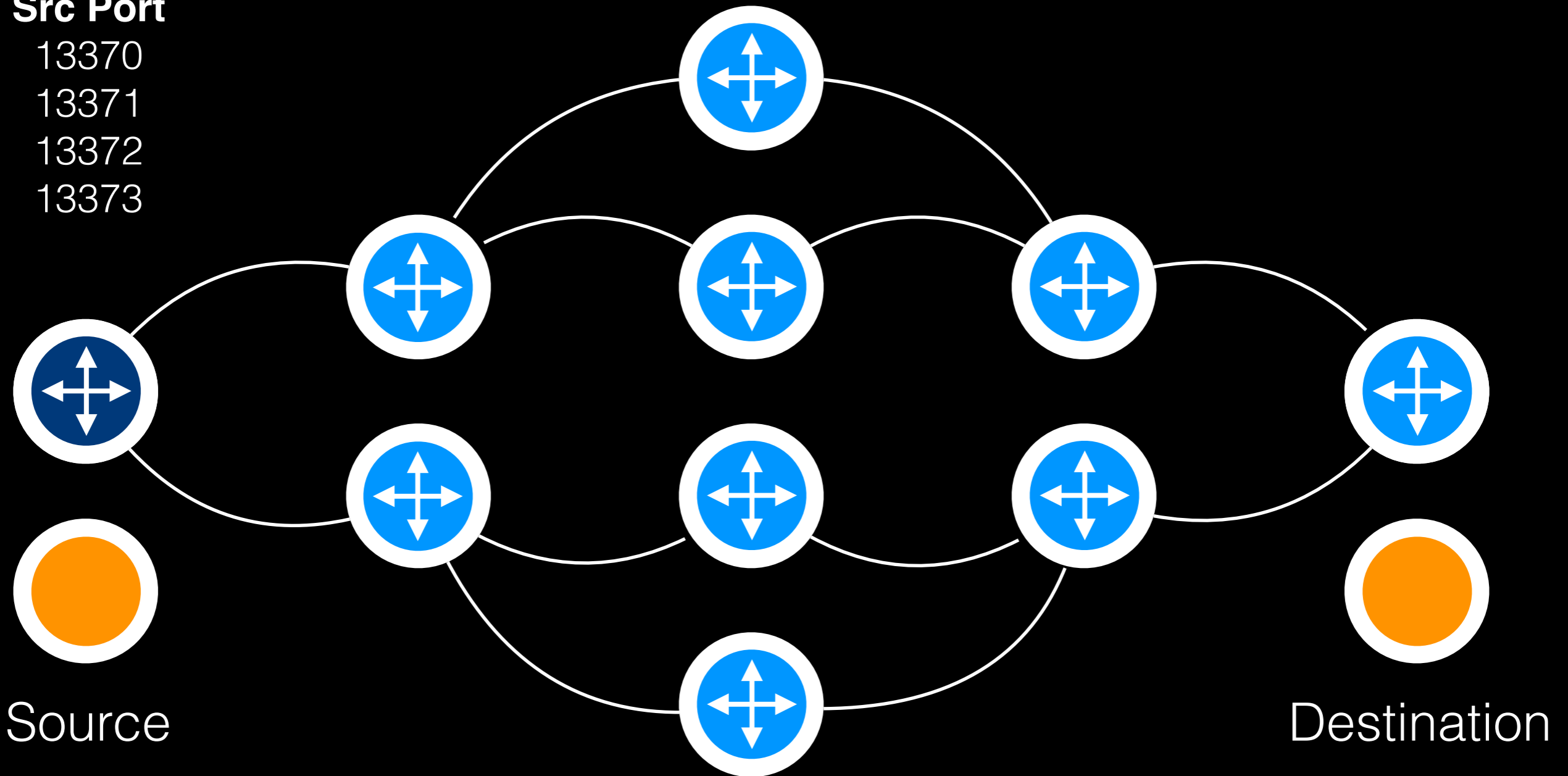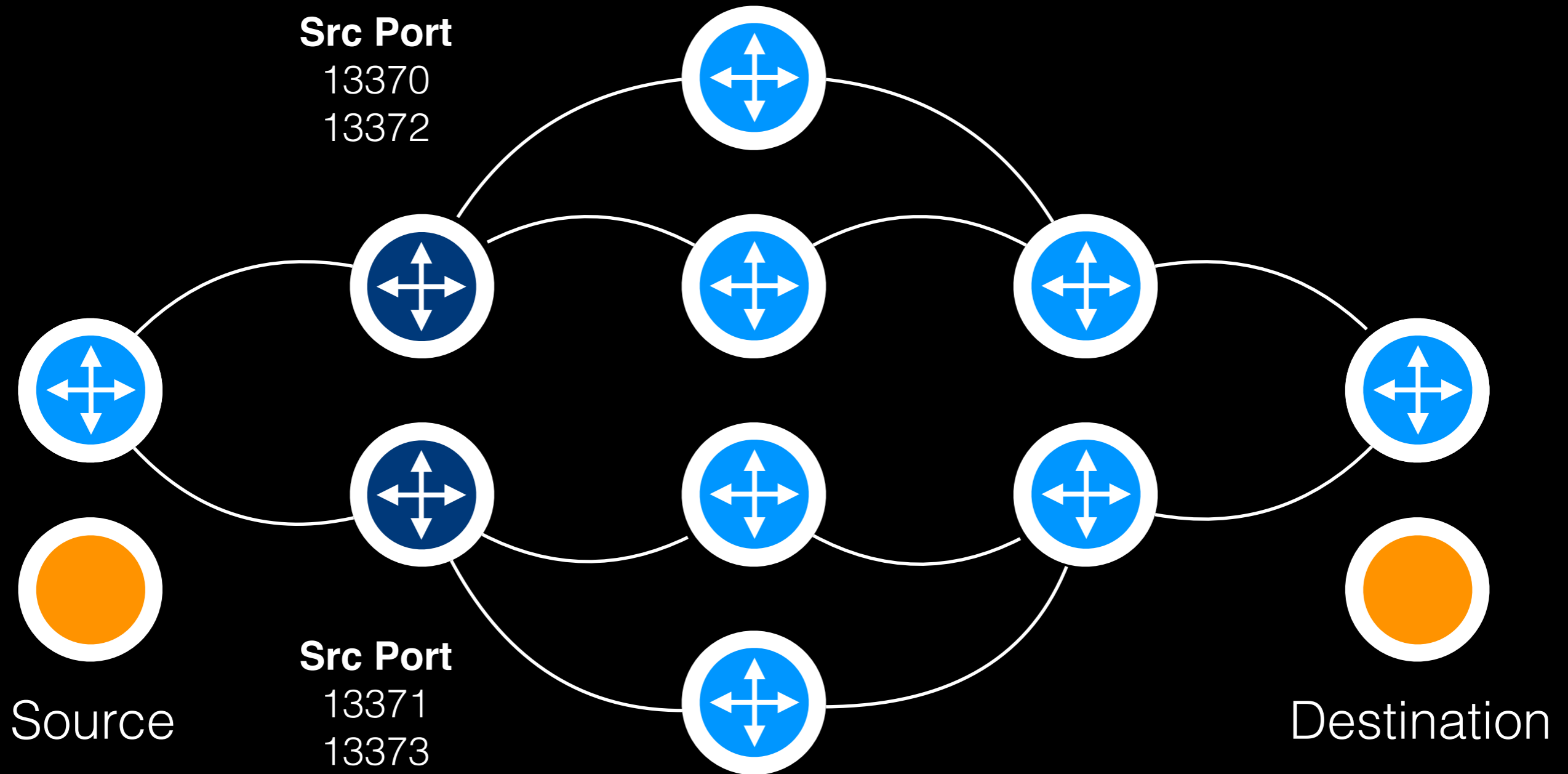
github.com/facebook/fbtracert

# fbtracert



Source
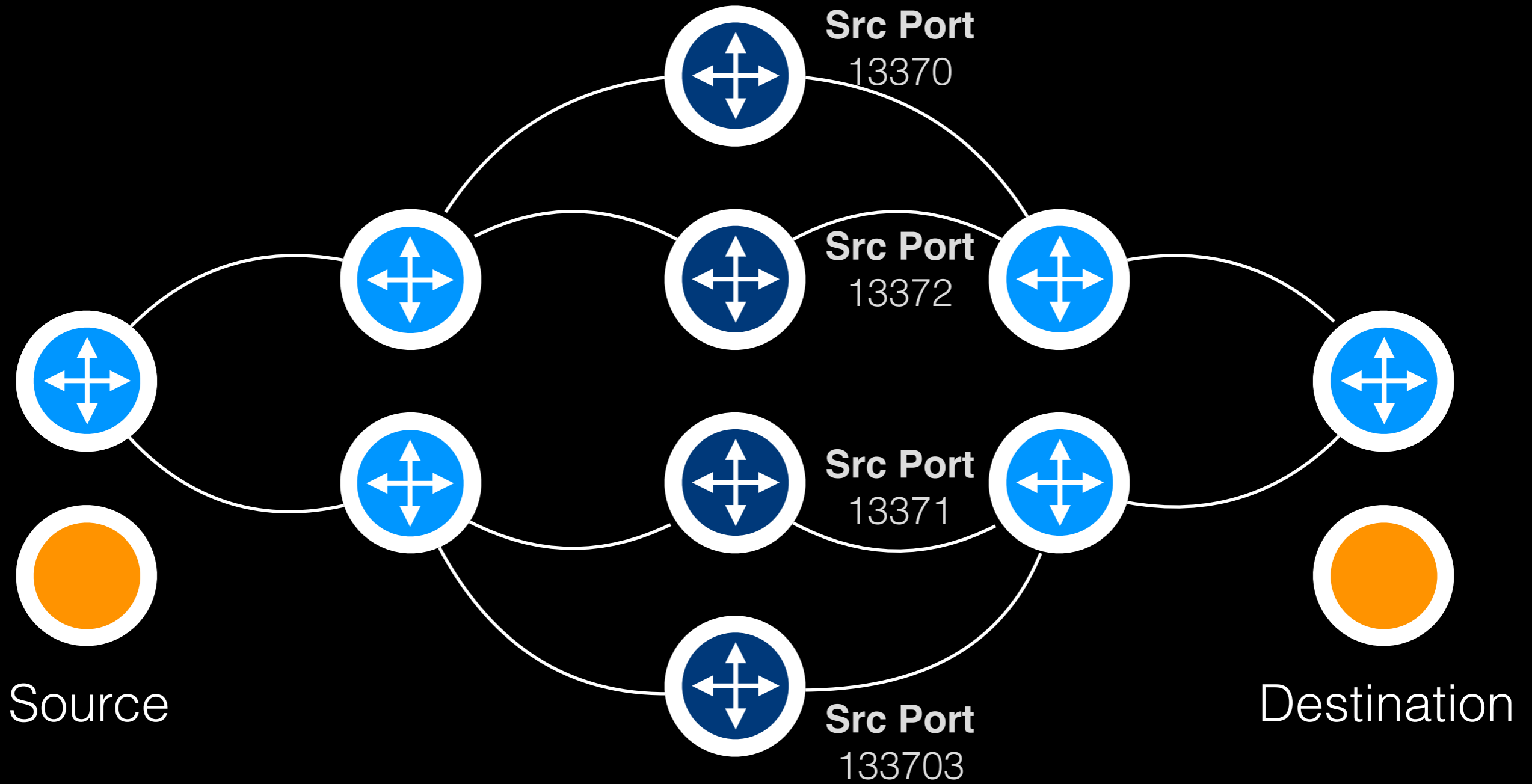
Destination

# fbtracert



**Src Port**
13370
13372

**Src Port**
13371
13373

Source

Destination

# fbtracert

# fbtracert

# fbtracert



**Src Port**
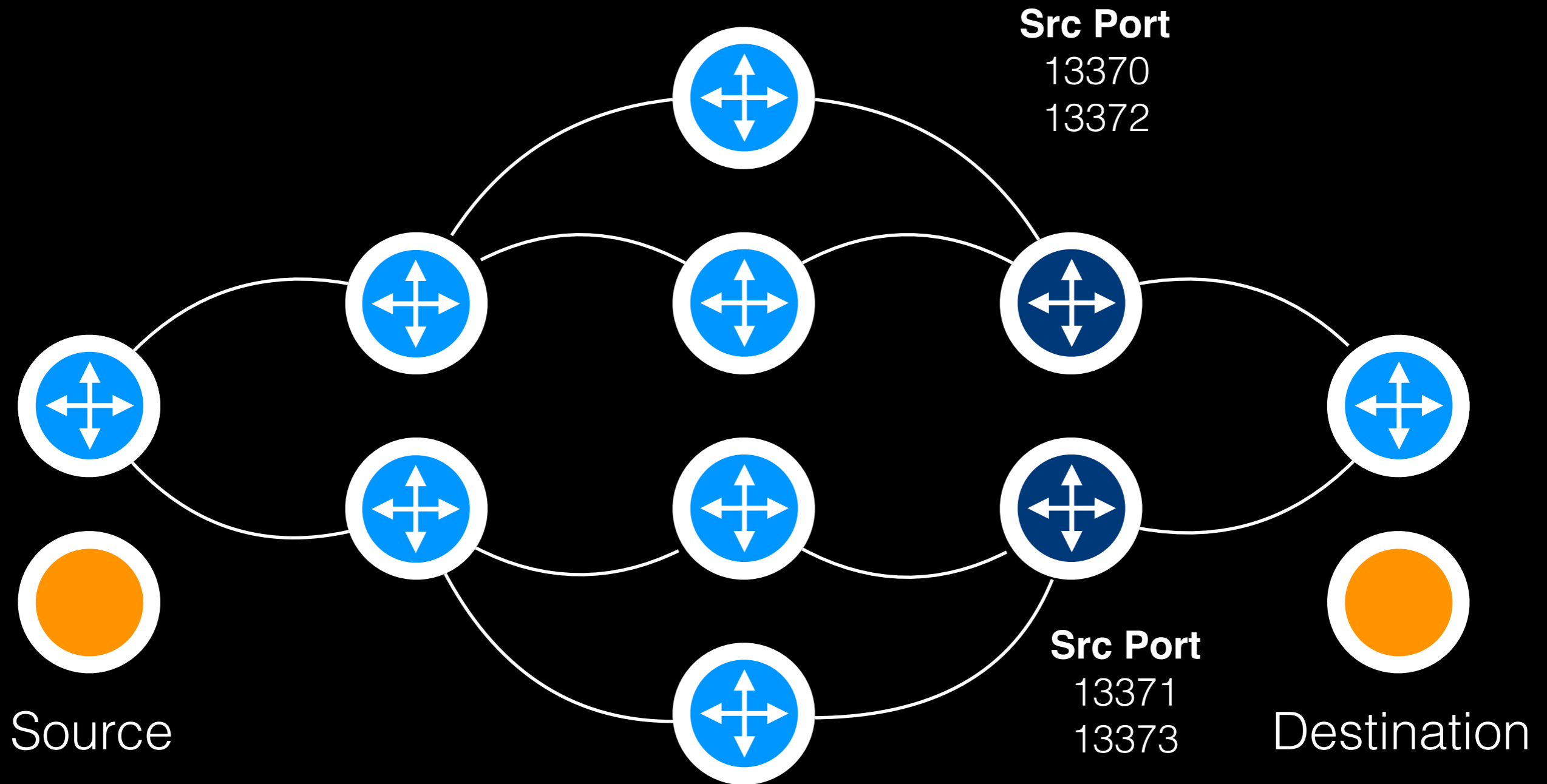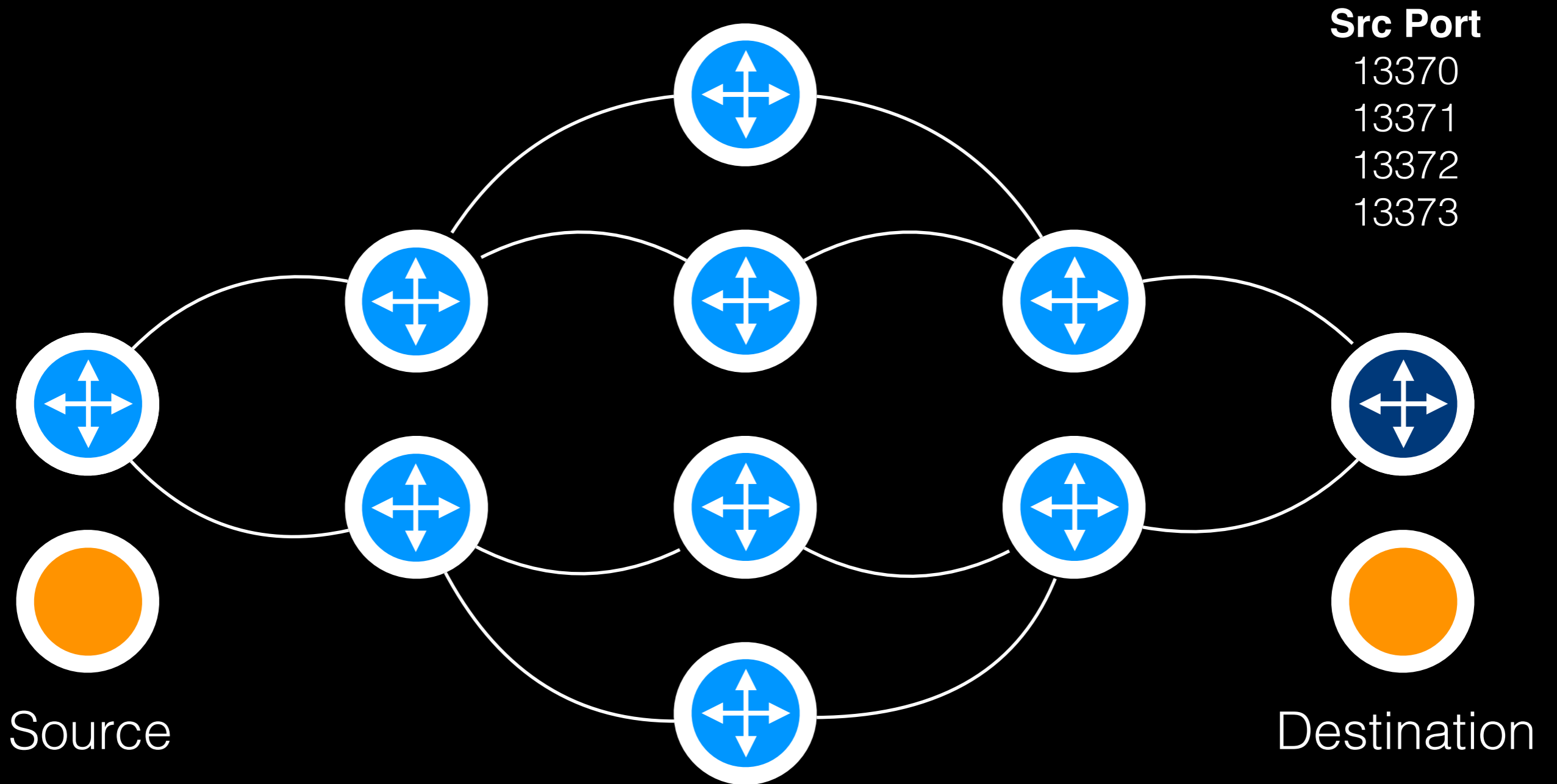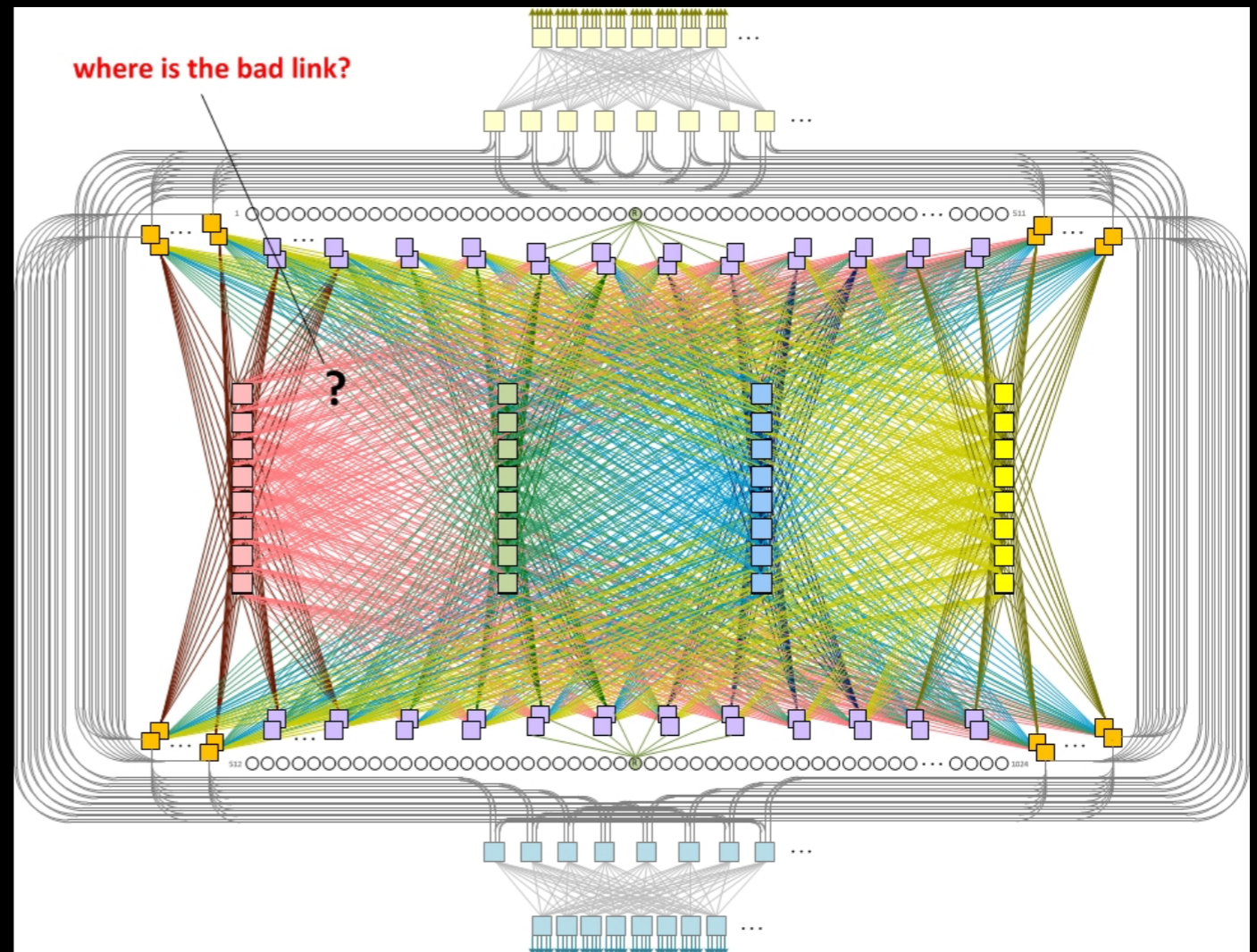13370
13371
13372
13373

Source

Destination
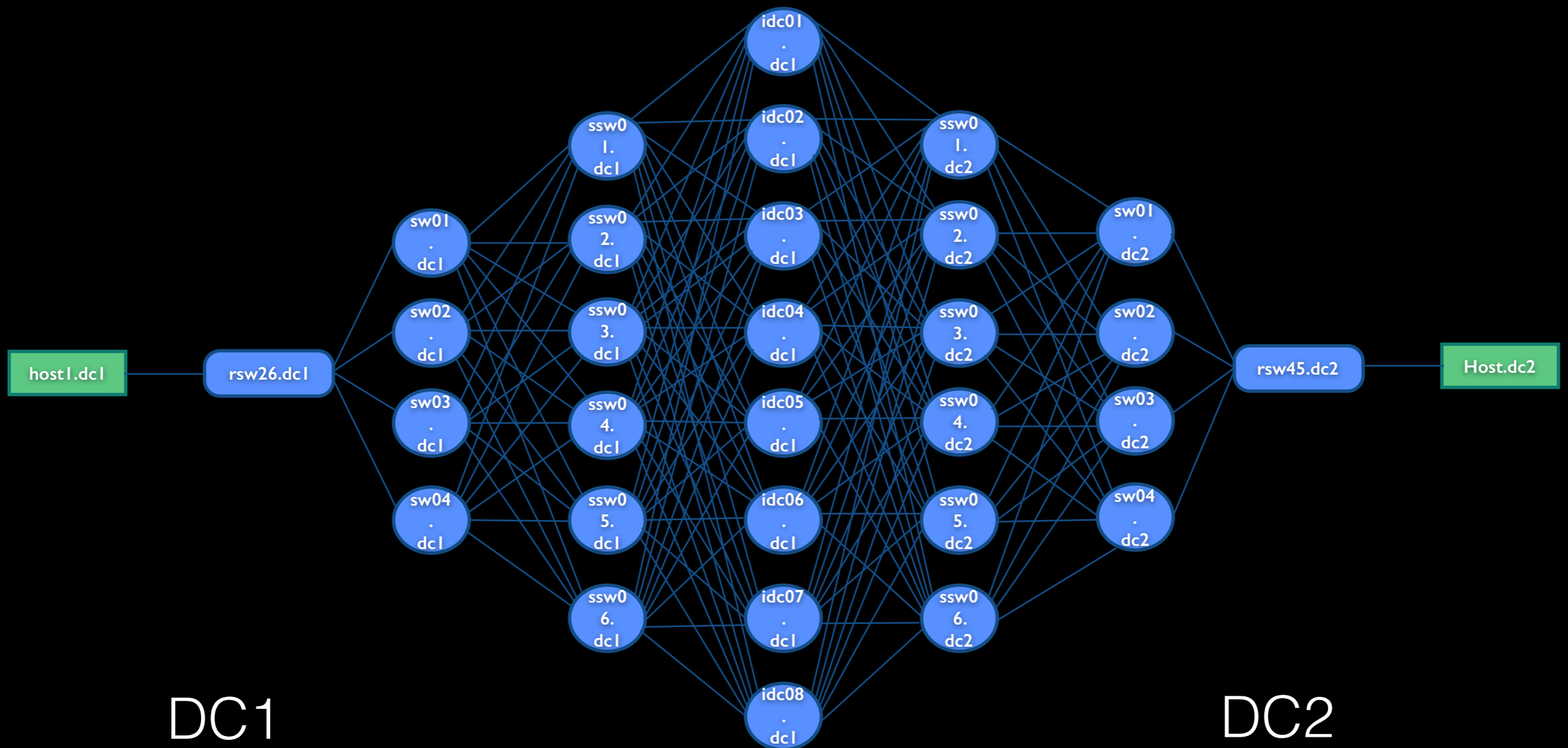
# DC Network Fault Isolation

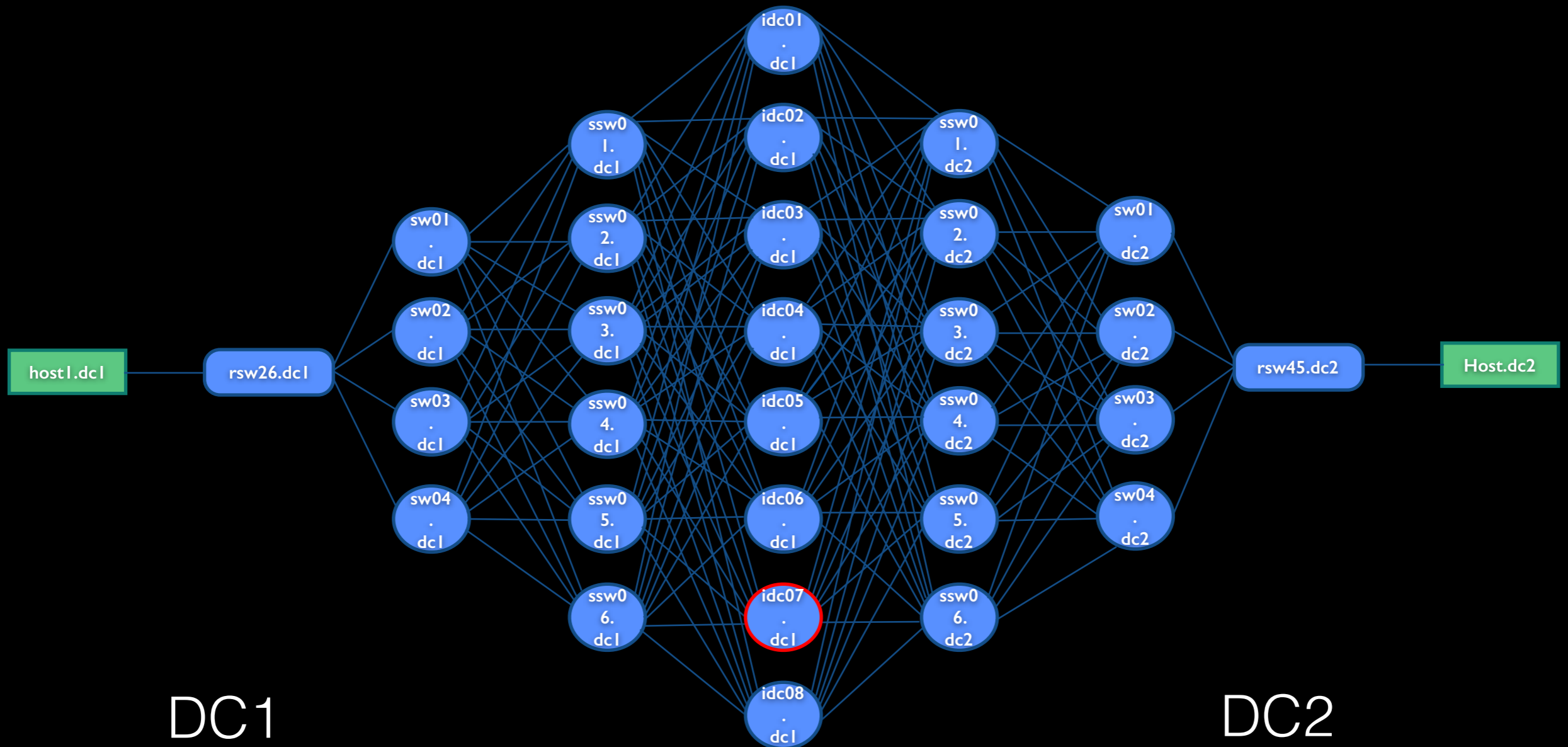## Isolating Network Faults

# Big Fat Fabrics

- Over 1,000 L3 Links between devices in different DCs in the SAME Region

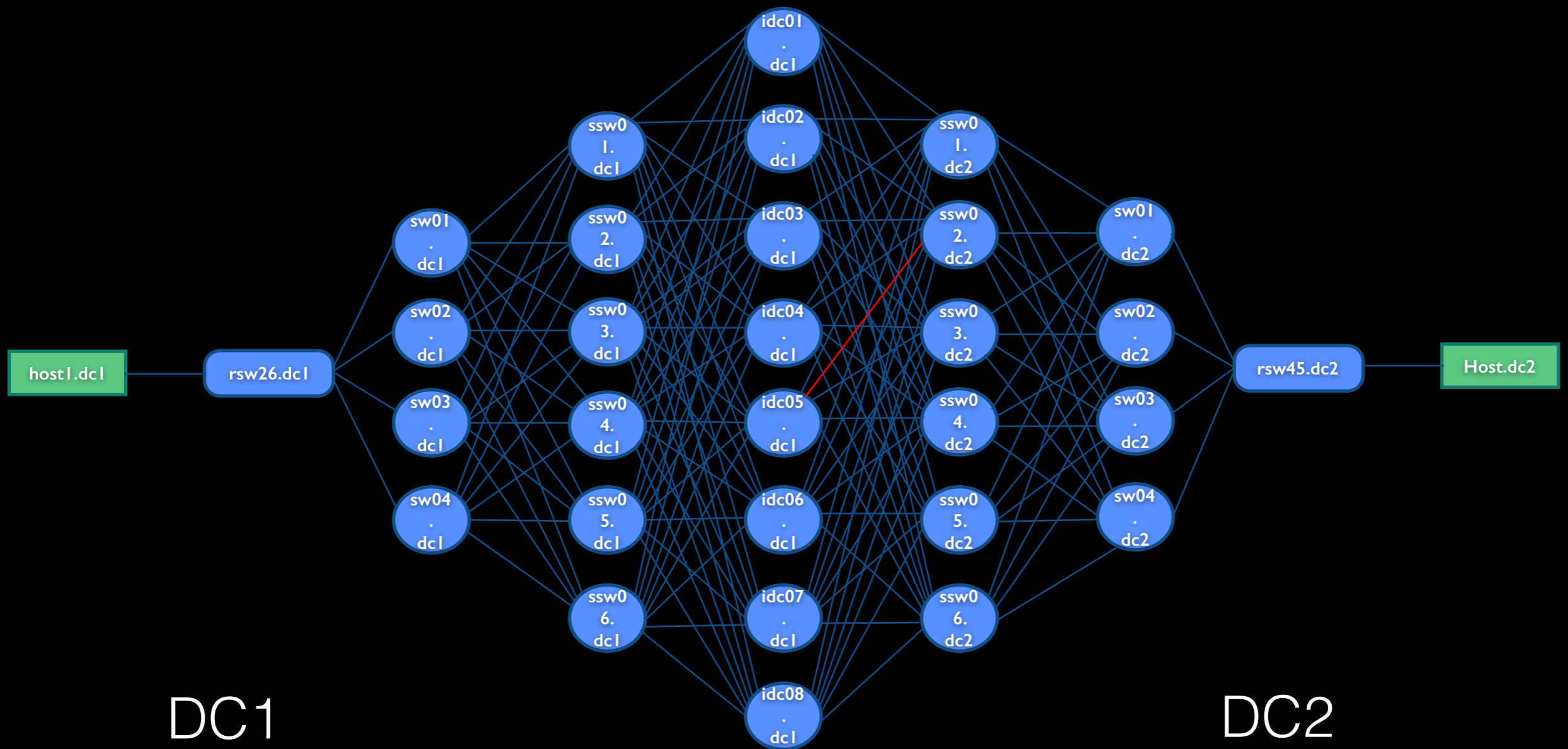- We know there is loss between hosts, but where?
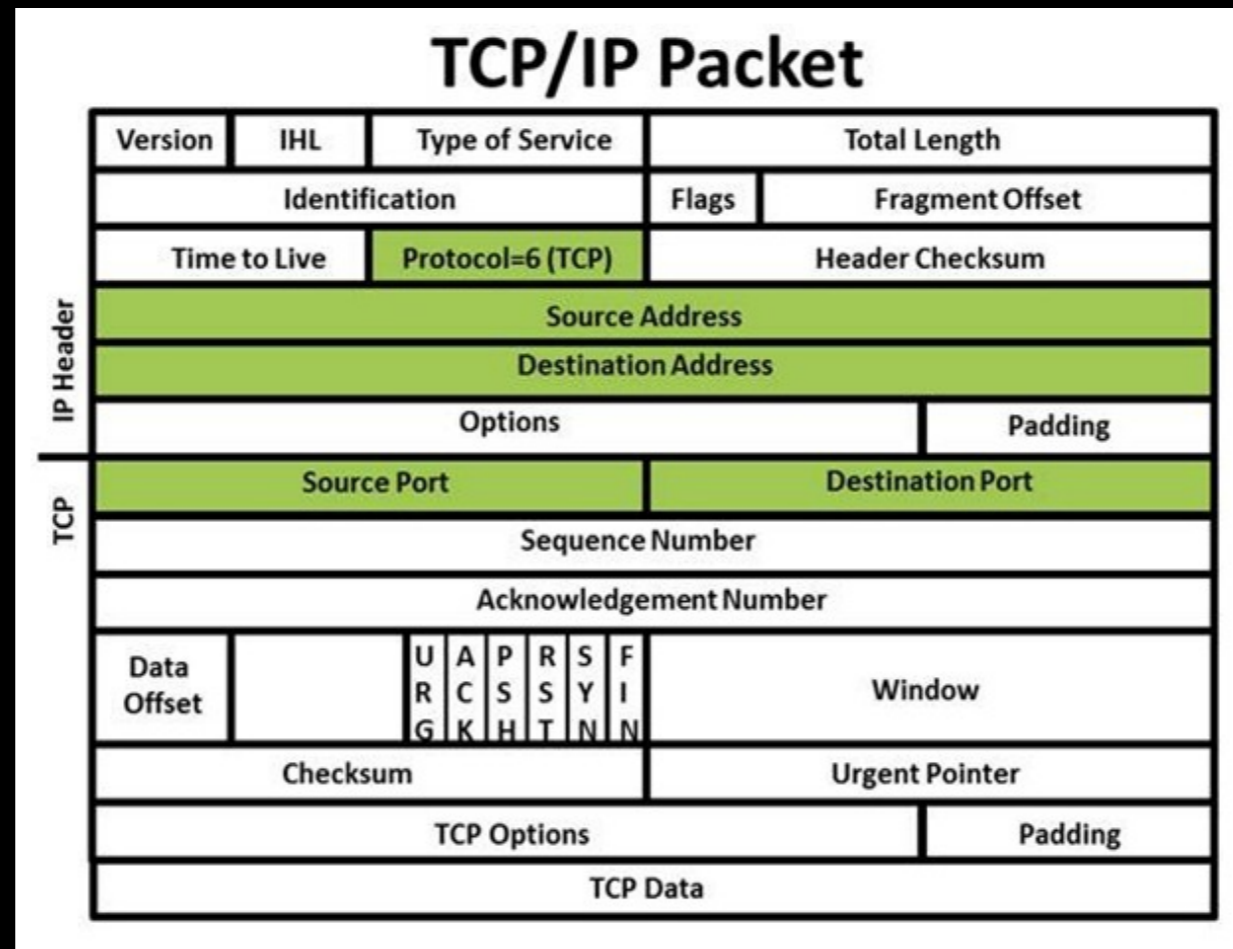
# DC Network Fault Isolation

# Bad Fabric Card


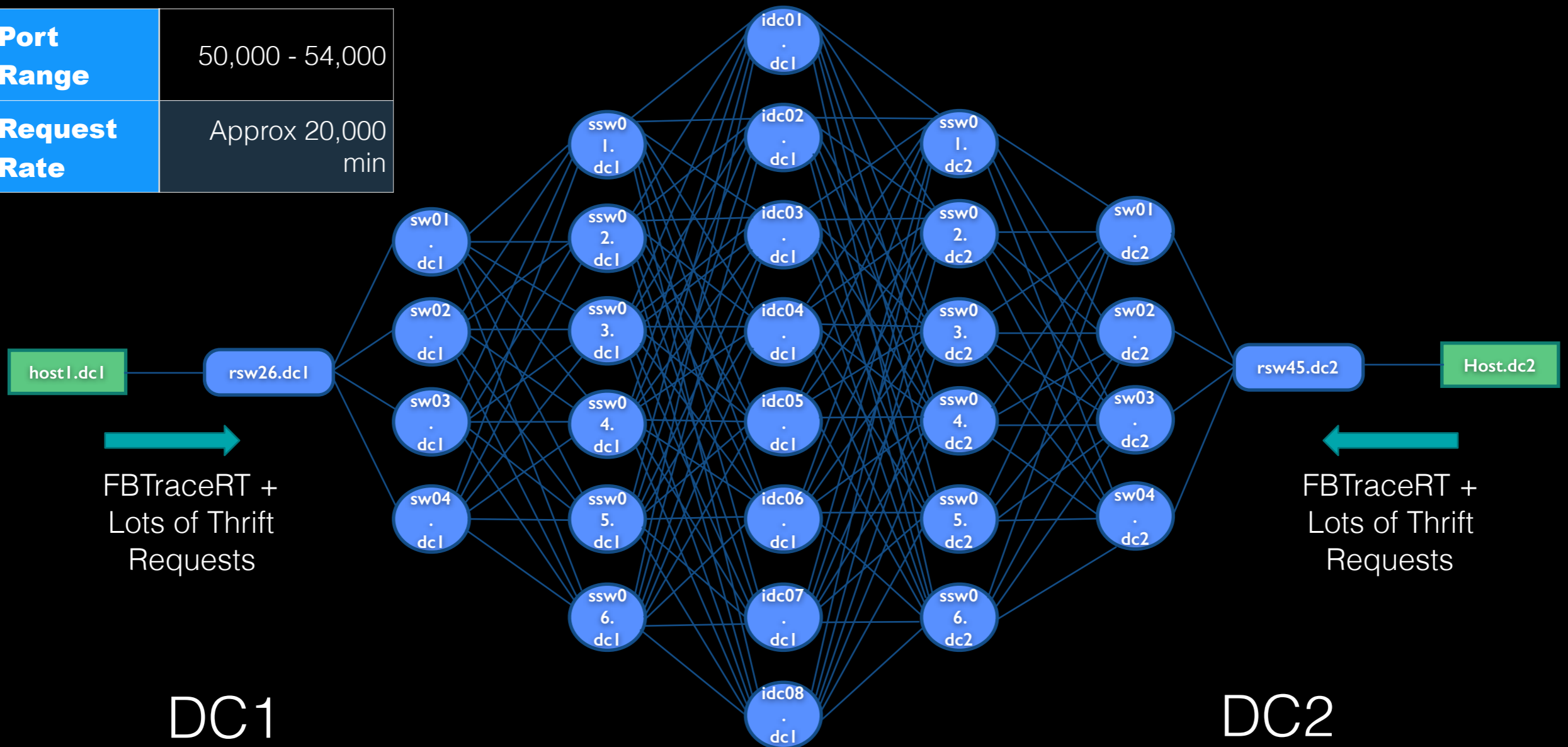
DC1

DC2

# ECMP Packet Hashing



TCP/IP Packet

# DC Network Fault Isolation

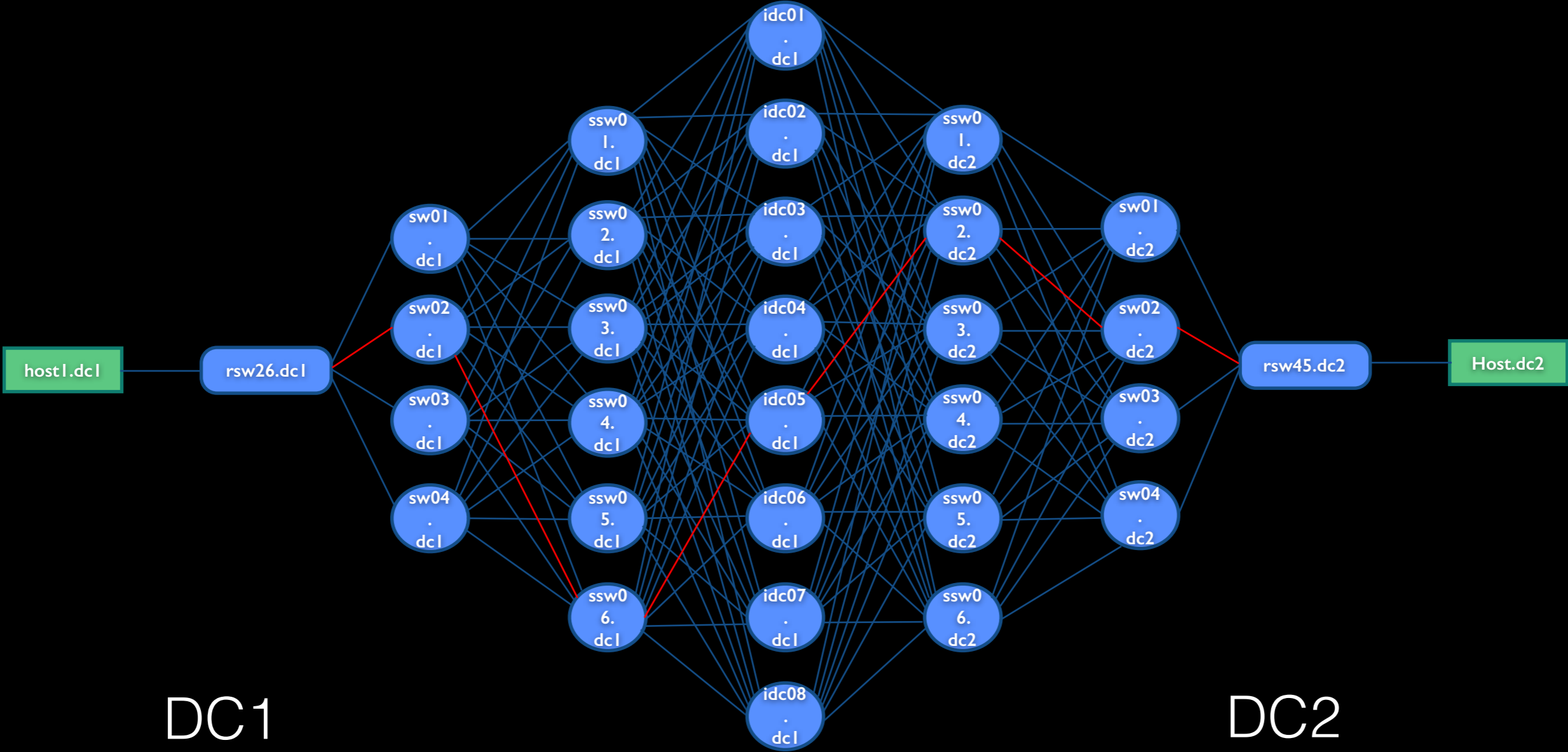Map The Loss for Each Request

# Map The Loss for Each Request

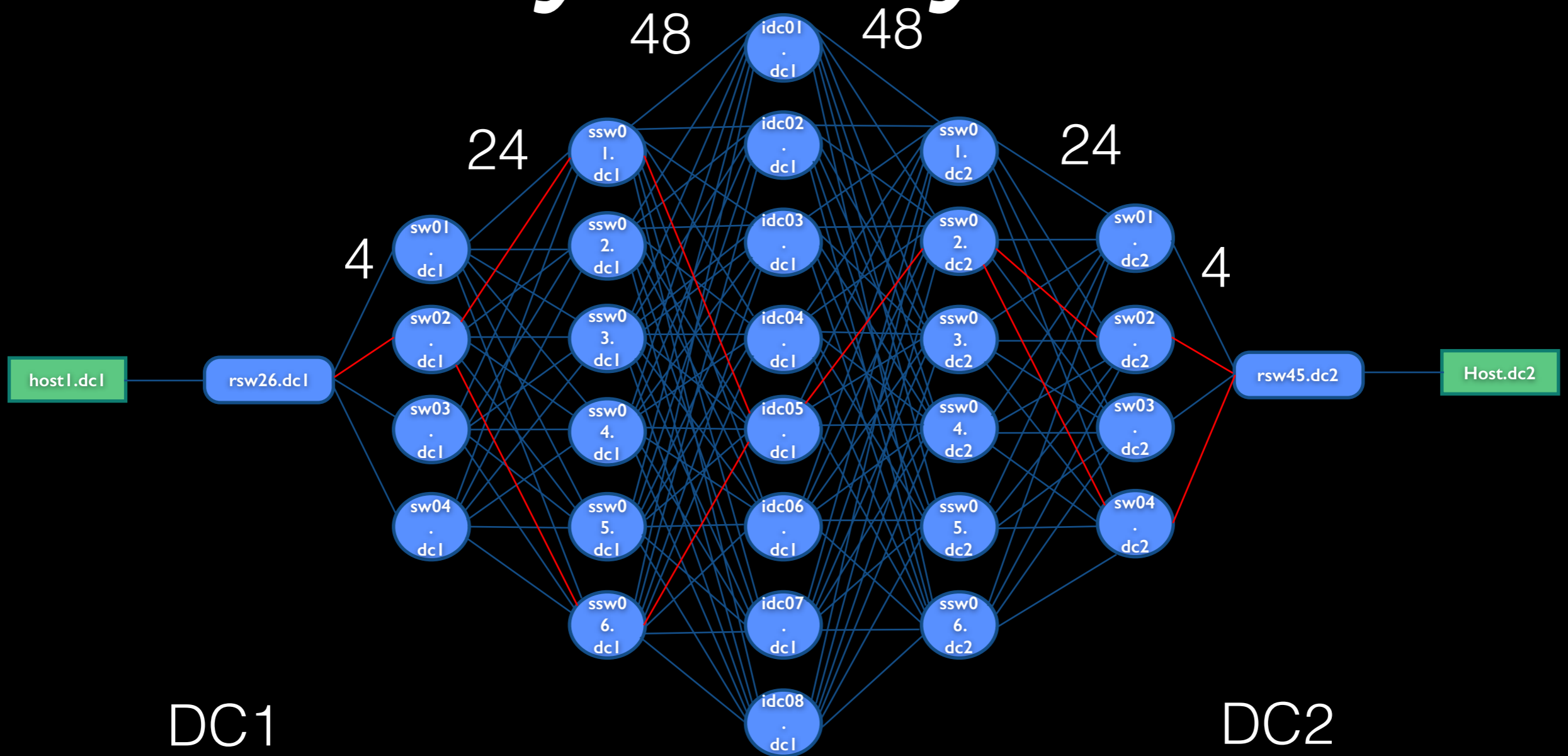# Analyze the Loss

# Clearer Signal

# Fabric Grey Failure Detection



TimeSeries (4 minutes) ordered by ECount

# Conclusions

- Fault isolation is <span style="color:yellow">actively evolving</span>

- Traceroute + probing approach is quite generic

- Limited by current hardware