### Move Fast, Unbreak Things! Network debugging at scale

Petr Lapukhov Network Engineer



## People who made this possible

Aijay Adams Lance Dryden Angelo Failla Zaid Hammoudi James Paussa James Zeng

### Basics of fault detection How people fix broken networks



# Data-center network (3)

- Multi-stage Clos Topologies
- Lots of devices and links
- BGP Only
- IPv6 >> IPv4
- Large ECMP fan-out
- L2 and L3 ECMP

#### Spine switches



#### Rack switches



## Backbone network (3)

- MPLS core
  - BB = Backbone Router (LSR)
- Data-center attachment
  - DR = Datacenter Router (LER)
- Auto-bandwidth
- ECMP over MPLS tunnels



## Detecting packet loss (4)

#### Standard counters



Non-Standard counters fsw001.p001.f01.atn1# show platform trident counters Debug counters Description T2Fabric19/0/1 RX - Non congestion discards T2Fabric19/0/1 TX - IPV4 L3 unicast aged and dropped pkts T2Fabric19/0/1 RX - Receive policy discard T2Fabric1)/0/1 TX - 12 Julticast crop T2Fabric1)/0/1 RK - Tuniel eiror paskets T2Fabric19/0/1 TX - Invalid VLAN T2Fabric19/0/1 RX - Receive VLAN drop T2Fabric19/0/1 RX - Receive multicast drop T2Fabric19/0/1 TX - Dropped because TTL counter T2Fabric19/0/1 RX - Receive uRPF drop T2Fabric19/0/1 TX - Packet dropped due to any condition T2Fabric19/0/1 RX - IBP discard and CPB full T2Fabric19/0/1 TX - Miss in VXLT table counter

# How human debugs it? (4)

- Ping/hping/nping (TCP/ICMP/UDP probing)
- Change src port to try all ECMP paths
- Find a broken path, then run traceroute over it
- ping && traceroute are still important

### NetNORAD The network fault detector



## Massive pinging FTW

- Run "pingers" on <u>some</u> machines
- Run responders on lots of machines
  - Targets count ~= 100x pingers count
- Collect packet loss and RTT...
- Analyze and report!

#### Pingers



#### Responders

# NetNORAD evolution (4)

- 1<sup>st</sup> Run`ping` from python agent
- 2<sup>nd</sup> Raw sockets, fast TCP probes
- 3<sup>rd</sup> Raw sockets, fast ICMP probes
- Now: UDP probing + responder agent

# Pinger and Responders (5)

#### Pingers

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

### Open sourced (C++), https://github.com/facebook/UdpPinger

#### Responders

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



# Allocating pingers and targets (2)

### Pingers

1+ cluster per DC 10+ racks per cluster Two pingers <u>per rack</u>

#### Targets

2+ targets per <u>each rack</u>10's of thousands targetsConsult host alarms

![](_page_11_Picture_5.jpeg)

# Probe timestamping

- Path changes / congestion
- Kernel time-stamps
- Application timeout tuning

![](_page_12_Picture_4.jpeg)

# Why UDP probing?

- No TCP RST packets
- Efficient ECMP
- RSS friendly
- Extensible

![](_page_13_Picture_5.jpeg)

#### Probe Format

#### Signature

#### SentTime

#### RcvdTime

#### ResponseTime

#### Traffic Class

## Deployment caveats (4)

### Caveat

Polarization with ICMP Slow IPv6 FIB lookups High-CPU boxes Checksum offloading

#### Solution

### Use UDP 4.X kernels Multi-threaded responder/RSS Disable offloading ☺

### NetNORAD How to ping and process data?

![](_page_15_Figure_1.jpeg)

![](_page_15_Picture_2.jpeg)

# Challenges (4)

- Nx 100Gbps of ping traffic
- Tens thousands of targets
- <u>Hundreds</u> of pingers
- Lots of data to process
- We really do not care about each host...
- ... The unit of interest is "cluster" health

each host... ster" health

![](_page_17_Figure_1.jpeg)

# Pinging inside clusters (4)

- Detect issues with rack switches
- Dedicated pingers per cluster
- Probe ALL machines in cluster
- Store time-series per host/rack
  - Think HBase for storage
- Lags real-time by ~2-3 minutes

![](_page_18_Figure_7.jpeg)

CSW – cluster switch RSW – rack switch

![](_page_18_Figure_9.jpeg)

![](_page_19_Figure_1.jpeg)

Pinger 2: Same Region Pinger 3: Outside of region

## Proximity tagging (3)

![](_page_20_Picture_1.jpeg)

![](_page_20_Picture_2.jpeg)

![](_page_20_Picture_3.jpeg)

![](_page_20_Picture_4.jpeg)

# Pinging hierarchy

Proximity	Scope	Goal
Outside of region	Across backbone network	End-to-end issues

![](_page_21_Picture_2.jpeg)

# Pinging hierarchy

Proximity	Scope	Goal
Outside of region	Across backbone network	End-to-end issues
Same region	Between data-centers in region	Issues inside/between DCs

![](_page_22_Picture_2.jpeg)

# Pinging hierarchy

Proximity	Scope	Goal
Outside of region	Across backbone network	WAN issues
Same region	Between data-centers in region	Issues between DCs
Same DC	Inside one data-center	Issues in cluster switches

![](_page_23_Picture_2.jpeg)

## Processing the data

#### Processing pipeline: Scribe (4) scribe - Scribe: distributed logging system - Similar OSS project: Kafka Data-set - Pingers write results Shard Shard Shard Shard Processors consume them - Propagation delay ~1-20 seconds pingers Processors (write) (read)

![](_page_25_Picture_6.jpeg)

# Alarming on packet loss (4)

- Build packet-loss time-series
- Track percentiles
- Alarm on rising threshold
- Clear on falling threshold
- Time to detect loss: 20 seconds

![](_page_26_Figure_6.jpeg)

## Visual analysis: Scuba

- In-memory row-oriented storage
- "Scuba: Diving into Data at Facebook"
- Similar OSS project: InfluxDB

![](_page_27_Picture_5.jpeg)

Detecting false-positives

# "Bad" target detection (3)

- Baseline loss
- Packet loss spike
- Filter outliers
- Done in pinger

Rack <sup>-</sup>
target
target

![](_page_29_Figure_6.jpeg)

## "Bad" Pinger problem (3)

- Bad cluster switch!
- Pingers see loss everywhere
- Population size is small
- Harder to weed outliers

![](_page_30_Figure_5.jpeg)

# "Bad" Pinger detection (2)

![](_page_31_Figure_1.jpeg)

- Need more data...
- Monitor pinger cluster
- Use DC/Region pingers
- Mark "bad" clusters
- Done in processor

## Conclusions

- Pinger/responder asymmetry
- Real-time is key
- Pinging hierarchy
- False positive elimination

### Isolating network faults Detecting is not everything

![](_page_33_Figure_1.jpeg)

![](_page_33_Picture_2.jpeg)

![](_page_34_Figure_1.jpeg)

Pinger 2: | see loss! Pinger 3: I see loss too!

### Downstream suppression (3) Data-Center Single alarm Data-Center LOSS Cluster 4 Cluster 1 Cluster 2 Cluster 3 OSS LOSS

![](_page_35_Figure_2.jpeg)

Multiple alarms

## Next steps to isolate (4)

- Approximate location
- Still lots of devices/links
- Check device counters
- if that does not help...
- Remember traceroute?

![](_page_36_Picture_6.jpeg)

### Fbtracert: fast and wide traceroute (6)

![](_page_37_Figure_1.jpeg)

## **Fbtracert: fast and wide traceroute**

#### Port 32701

Path	Sent	Rcvd
1	20	20
2	20	20
4	20	20
8	20	20
10	20	14
TGT	20	15

#### Port 32702

Path	Sent	Rcvd
1	20	20
3	20	20
6	20	20
9	20	20
10	20	20
TGT	20	20

#### Port 32703

Path	Sent	Rcvd
1	20	20
2	20	20
5	20	20
8	20	20
10	20	16
TGT	20	17

#### Port 32704

Path	Sent	Rcvd
1	20	20
3	20	20
7	20	20
9	20	20
10	20	20
TGT	20	20

![](_page_38_Picture_9.jpeg)

## Fbtracert limitations (5)

- CoPP drops ICMP responses
- Paths may flap (MPLS LSP)
- ICMP gets tunneled with MPLS TE
- ICMP responses from wrong interfaces

Open sourced (Golang), https://github.com/facebook/fbtracert

![](_page_39_Picture_8.jpeg)

## Conclusions

- Fault isolation is actively evolving
- Traceroute approach looks generic
- Limited by current hardware
- Backbone path churn is a serious challenge

### Evolving fault detection & isolation Near and far future

![](_page_41_Figure_1.jpeg)

![](_page_41_Picture_2.jpeg)

# Support for on-box agents (4)

- Run same code on routers
- POSIX API
- Other SDK is welcome
- Some vendors already do that
- Be like FBOSS ③

![](_page_42_Picture_7.jpeg)

# Streaming telemetry (3)

- Publishing device counters
- Faster detection
- Protobuf/Thrift for encoding
- Limited amount of counters
- Platform-specific

![](_page_43_Figure_6.jpeg)

![](_page_43_Picture_7.jpeg)

## In-band telemetry (4)

- Next generation of silicon emerging
- Embed device stats in packets
  - E.g. device ID, or queue depth
- Use extra space in UDP probes
- Allow for real-time path tracing

IP/UDP hdr Switch Queue depth **Device ID** IP/UDP hdr

![](_page_45_Picture_1.jpeg)