# Agenda

- Background
- NodeLocal Caching
- Local Benchmarks
- Cluster Benchmarks
- e2e Application Benchmarks
- Outlook
- Questions

KubeCon | CloudNativeCon
North America 2018

# Background

# Kube-dns History

- Bundle of dnsmasq + [SkyDNS](#)

- SkyDNS written in Go by [Miek Gieben](#) and others*

- Partly maintained

*Erik st. Martin, Brian Ketelsen, Michael Crosby

# CoreDNS History

- Authored by [Miek Gieben](#)

- Based on [Caddy](#) (Golang webserver)

- Plugin-based architecture for extensibility++

- GA in 1.11, default in 1.13

# Protocol support

- UDP (RFCs: 1034, 1035, etc.)

- TCP (RFC 7766)

- TLS (DoT) (RFCs: 7858, 8310)

- HTTPS (DoH) (RFC 8484)

- GRPC (DoG?) (RFCs: none)

# Internet Considered Harmful

"Aaaand it's gone"
- *South Park*

- Standards from 70s & 80s

- Assumption of most DNS records being ~static

- Congestion + availability > consistency + reliability

- Old decisions can't keep up w/ new usage patterns

# DNS 1.0 (RFCs: 1034, 1035)

- Requests generally occur over **UDP**, except under special circumstances. Per RFC1035 4.2:

  - "The DNS assumes that messages will be transmitted as datagrams (UDP) or in a byte stream carried by a virtual circuit (TCP). While virtual circuits can be used for any DNS activity, **datagrams are preferred for queries due to their lower overhead and better performance.**"

  - "Depending on how well connected the client is to its expected servers, the minimum retransmission interval should be **2-5 seconds**."

**timeout**:n

> Sets the amount of time the resolver will wait for a response from a remote name server before retrying the query via a different name server ... Measured in **seconds**, the default is RES_TIMEOUT (**currently 5**, see <resolv.h>) ...

**attempts**:n

> Sets the number of times the resolver will send a query to its name servers before giving up ... The default is RES_DFLRETRY (**currently 2**, see <resolv.h>).

# conntrack limits & races

- Cluster DNS is a k8s **Service**

- DNAT rules used to translate ClusterIP to Pod IP

- conntrack table usually limited to **65536 entries**
  => **dropped packets**

- Multiple conntrack table entries per 'connection' (including UDP)

  - No UDP 'close' → entries persist long after they're useful

# conntrack limits & races

- multiple UDP reqs from the same ip:port can results in **race conditions**
      => **dropped packets**

- Races aggravated by **parallelized** reqs for

  different record types (e.g. A, AAAA)

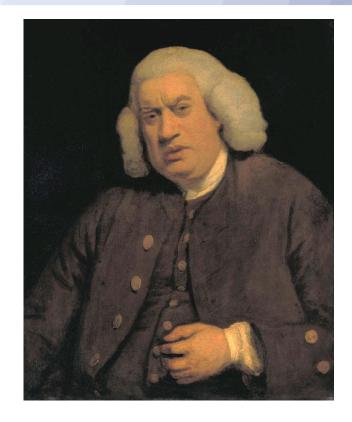- N search paths mean **N times** more requests for

  failed queries

# We've All Been There

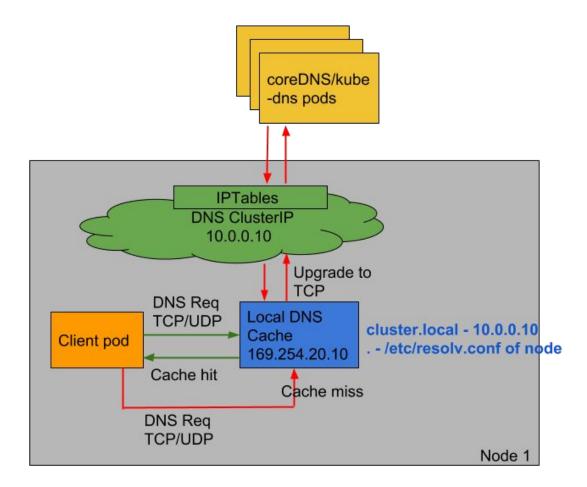nf_conntrack: table full, dropping packet

# Enter: NodeLocal DNS Cache

Coming soon (optionally) to a 1.13 cluster near you!

Runs on every node, serves DNS for pods that are using cluster DNS.

- Improve latency by reducing communication over the network
- **Skips conntrack** for pod-cache connection
  - Less dropped packets!
- Proxy queries over TCP (and preserves the connection)
  - DARPA-grade reliability & consistency!
  - Even less pressure on the cluster DNS's conntrack tables (see above)
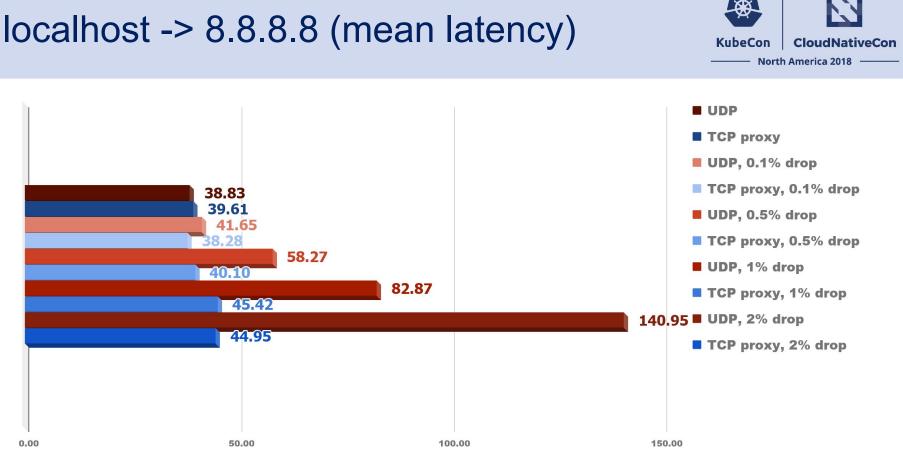- Node-level DNS metrics

Special thanks: Pavithra Ramesh, GKE
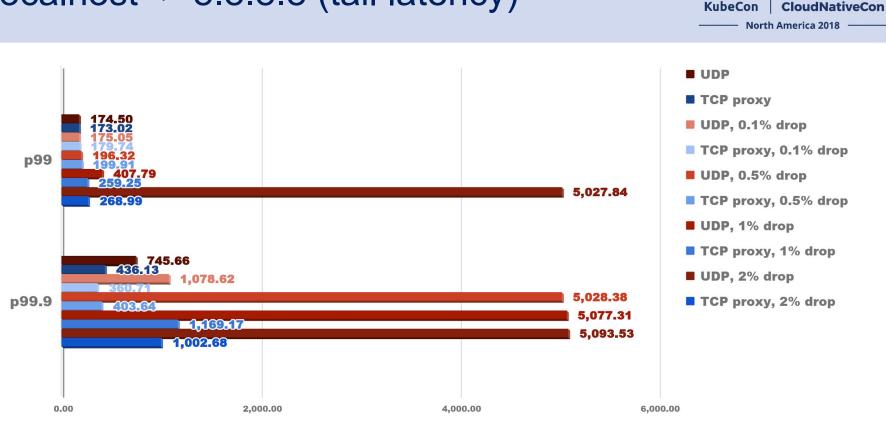
Local Benchmarks

# localhost -> 8.8.8.8 (mean latency)



Code @ https://github.com/dekkagaijin/coredns-demo/

# localhost -> 8.8.8.8 (tail latency)



Legend:
- UDP
- TCP proxy
- UDP, 0.1% drop
- TCP proxy, 0.1% drop
- UDP, 0.5% drop
- TCP proxy, 0.5% drop
- UDP, 1% drop
- TCP proxy, 1% drop
- UDP, 2% drop
- TCP proxy, 2% drop

**p99**
- 174.50
- 173.02
- 175.05
- 179.74
- 196.32
- 199.91
- 407.79
- 259.25
- 5,027.84
- 268.99

**p99.9**
- 745.66
- 436.13
- 1,078.62
- 360.71
- 5,028.38
- 403.64
- 5,077.31
- 1,169.17
- 5,093.53
- 1,002.68

x-axis: latency (ms)

0.00    2,000.00    4,000.00    6,000.00

# CoreDNS vs Kube-DNS: Memory



CoreDNS vs Kube-DNS Est Memory at Scale

coredns    8.58E-04*x + 54.1    kube-dns    9.07E-04*x + 95.9

Credit: Chris O'Haver, Infoblox

# CoreDNS vs Kube-DNS: Queries

| DNS Server | Query Type | QPS | Avg Latency (ms) |
|---|---|---|---|
| CoreDNS | external | 6733 | 12.02 |
| CoreDNS | internal | 33669 | 2.608 |
| Kube-dns | external | 2227 | 41.585 |
| Kube-dns | internal | 36648 | 2.639 |

Credit: Chris O'Haver, Infoblox

Application Benchmarks

# Our application: TXTDirect

- DNS [TXT record](#)-based redirects

- Control over your entrypoint and data

- Open Source based on Caddy

- Does a lot of DNS requests

# TXTDirect: Request flow

*"GET"* kubernetes.opensourcesoftware.rocks

    *"A/AAAA/CNAME"* for **kubernetes.opensourcesoftware.rocks**

    *TXT* for **_redirect.kubernetes.opensourcesoftware.rocks**

    "v=txtv0;type=host;to=https://kubernetes.io"



**kubernetes.io**

Learn more on txtdirect.org

# Setup: Standard Kube-dns



Pod, UDP → Kube-dns

# Setup: Standard CoreDNS



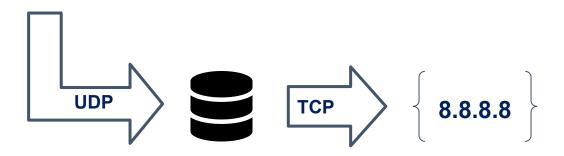Pod, UDP → CoreDNS, TCP → 8.8.8.8

# Setup: NodeLocal cluster



Pod, UDP → NL, TCP → CoreDNS, TCP → 8.8.8.8

# Setup: NodeLocal direct



Pod, UDP → NL, TCP → 8.8.8.8

# Show me the numbers!

# Probability Distribution (Latency)

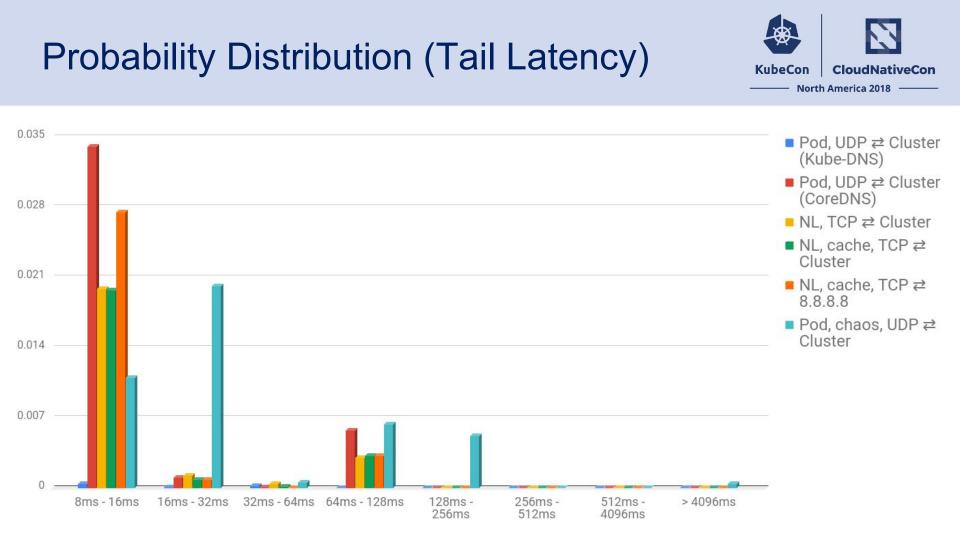# Probability Distribution (Tail Latency)

# What we learned

- TCP does what it's supposed to

- TCP forwarding improves the performance of traditional (UDP) clients, with less variance, and without incurring a ton of overhead

- CoreDNS's plugins make it a good fit for special use-cases

Outlook

# Future Work

- Native DNS over GRPC

- Watch based DNS records

- Performance and reliability improvements

- Ideas? Let us know after the talk!

# Drinks…