

Old but Gold: Prospecting TCP to Engineer and Real-time Monitor DNS Anycast

Giovane C. M. Moura¹ John Heidemann² Wes Hardaker²
Jeroen Bulten³ João Ceron¹ Cristian Hesselman^{1,4}

1: SIDN Labs, 2: USC/ISI, 3: SIDN, 4: University of Twente

OARC 34

Virtual Conference

2021-02-04



Latency is key in DNS (but hard to measure)

- Authoritative OPs will use whatever tools to reduce latency:
 1. multiple NSes
 2. Anycast
 3. Peering/IXPs
 4. ...
- But is **hard to know** client's latency:
 1. Ripe Atlas, Thousand Eyes: good but not complete coverage
 2. Verfploeter [1]: requires ICMP measurements
 - Verfploeter is ran typically daily, as it is expensive
 - Difficult to apply to IPv6 (hitlist)

What if there was a better way ?

- A method that:
 - Comes from *real-clients*
 - Works well with IPv6
 - Requires *no extra* measurements (passive only)
- Well, there is one: **DNS over TCP (DNSTCP)**
 - RTT measured from handshake (or takedown)
 - we've been using for 1.5 years at SIDN (.nl)
 - helped to solve several issues
 - fulfills all the above

What if there was a better way ?

- A method that:
 - Comes from *real-clients*
 - Works well with IPv6
 - Requires *no extra* measurements (passive only)
- Well, there is one: **DNS over TCP (DNSTCP)**
 - RTT measured from handshake (or takedown)
 - we've been using for 1.5 years at SIDN (.nl)
 - helped to solve several issues
 - fulfills all the above

TCP RTT history: old but gold

- TCP RTT estimation has been used since 1996 [2]
- Widely used in passive analysis of HTTP (FB uses it [5])
- It has been applied on DNS multiple times:
 - Roy Arends (2012)
 - Casey Deccio (2018)
 - Maciej Andzinski [3] (2019)
 - Our tech report (2020) [4]

So what's NEW with our work?

- extensive and comprehensive methodology validation
 - Is the TCP data representative?
 - Are the UDP and TCP latency comparable?
- acted upon the data with 4 operators (Anycast A, B, B-Root, and Google)
 - We identify several use cases and issues
 - We manipulated BGP to fix those issues
 - We document it carefully
- use in real-time within .nl to detect anomalies
 - Route leaks

TCP traffic **MUST**:

1. Provide enough **coverage** (spatial and temporal)
 - you know, most DNS traffic is still UDP
2. provide **similar latency** to UDP
 - so we can generalize the results

Is DNS traffic representative?

	Queries		Resolvers		ASes	
	Anycast A	Anycast B	Anycast A	Anycast B	Anycast A	Anycast B
Total	5 237 454 456	5 679 361 857	2 015 915	2 005 855	42 253	42 181
IPv4	4 005 046 701	4 245 504 907	1 815 519	1 806 863	41 957	41 891
UDP	3 813 642 861	4 128 517 823	1 812 741	1 804 405	41 947	41 882
TCP	191 403 840	116 987 084	392 434	364 050	18 784	18 252
<i>ratio TCP</i>	5.02%	2.83%	21.65%	20.18%	44.78%	43.58%
IPv6	1 232 407 755	1 433 856 950	200 396	198 992	7 664	7 479
UDP	1 160 414 491	1 397 068 097	200 069	198 701	7 662	7 478
TCP	71 993 264	36 788 853	47 627	4 6190	3 391	3 354
<i>ratio TCP</i>	6.2%	2.63%	23.81%	23.25%	44.26%	44.85%

Table 1: DNS usage for two authoritative services of .nl (Oct. 15–22, 2019).

- 5% of clients, 20% of resolvers, and 44% of ASes
- You get this for free
- Roots: 1.77–14% of TCP queries (see report [4])

Important ASes use TCP

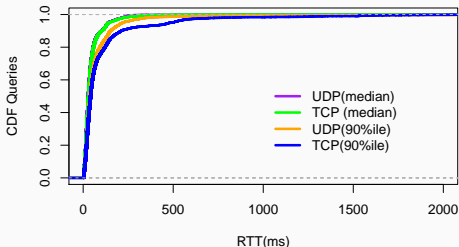
	Anycast A	Anycast B
IPv4	4 005 046 701	4 245 504 907
from TCP ASes	3 926 025 752	4 036 328 314
Ratio (%)	98.02%	95.07%
from TCP resolvers	2 306 027 922	1 246 213 577
Ratio (%)	57.7%	29.35%
IPv6	1 232 407 755	1 433 856 950
from TCP ASes	1 210 649 060	1 386 035 175
Ratio (%)	98.23%	96.66%
from TCP resolvers	533 519 527	518 144 495
Ratio (%)	43.29%	36.13%

Table 2: Queries per Services for ASes and Resolvers that send TCP queries for .nl (Oct. 15–22, 2019).

- ASes that do TCP send most of the traffic

DNS: TCP vs UDP latency are comparable

	K-Root		L-Root	
	UDP	TCP	UDP	TCP
Date	Sept 4–5, 2020		Sept 5–6, 2020	
Freq.	4min	8min	4min	8min
Probes	10520	8676	10586	8989
\cap Probes	8582		8892	
Queries	3749892	1045605	3779763	1062557
\cap Queries	3063836	1034233	3181098	1055888

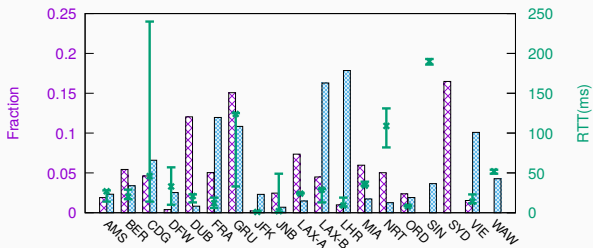
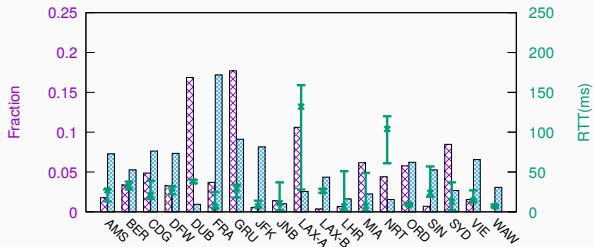


OK, so what can we do with it?

- DNS/TCP provides enough VPs
- Has similar latency than UDP
- Measure real clients
- No costs
- Easily copes with IPv6
- Requires no extra measurements
- Can be run in real time

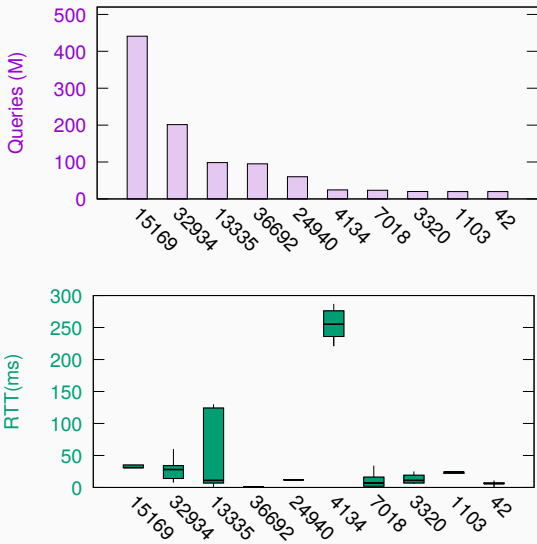
Prioritizing Analysis: by Site

Anycast B: IPv4 and IPv6 RTT per site



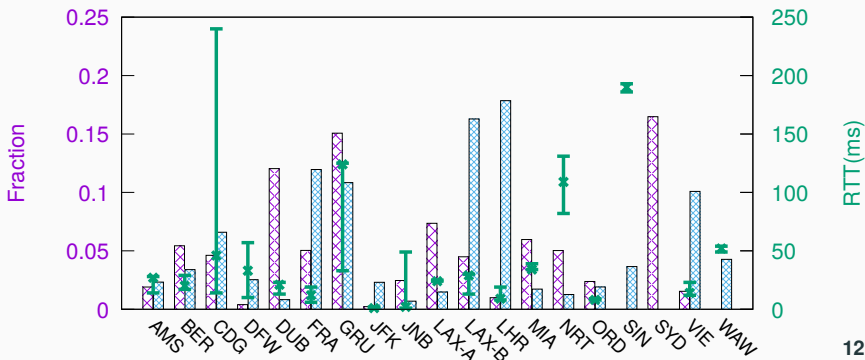
Prioritizing Analysis: by client AS

Anycast B: IPv6 queries and RTT per client AS



Problems: Distant Lands

- A client is mapped by BPG to far distant anycast sites
- Some sites have a large RTT value or spread (CDG, SIN, NRT)
- We can see that using DNS/TCP RTT



Solutions: Distant Lands (NRT)

- Causes: No presence/direct peer with Chinese ISPs
- Chinese int'l connections can exhibit congestion [6]
- Fix: site in China (OPs clients may not be comfortable) or direct peer (\$)

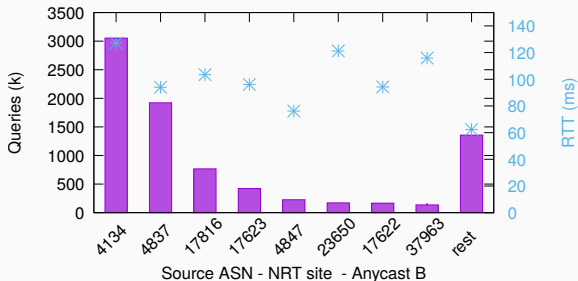


Figure 2: Anycast B, Japan site (NRT): Top 8 querying ASes are Chinese, and responsible for 80% of queries.

Problems: prefer customer to another continent

- Common BGP policy: prefer customer
 - if AS can satisfy route via customer, so be it
- But sometimes it takes clients to another continent
- We found Comcast (US, AS7922) reaching Anycast B via GRU site (Brazil)
- We contacted the Operator; fixed with right BGP community

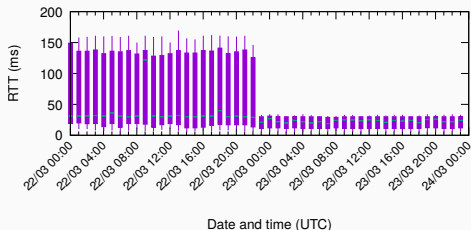
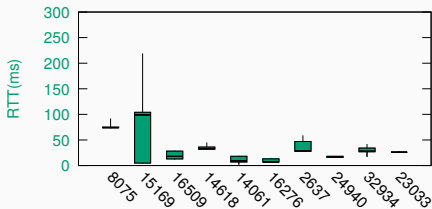
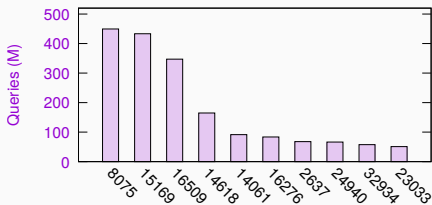


Figure 3: Anycast B and Comcast: RTT before and after resolving IPv6 misconfiguration.

Problem: Anycast Polarization

- We found that MS (8075) and Google (15169) had high latencies to Anycast A
- And they are the top 2 client ASes



Problem: Google Polarized → high latency)

- All Google Traffic was going to AMS site only : RTT 100ms

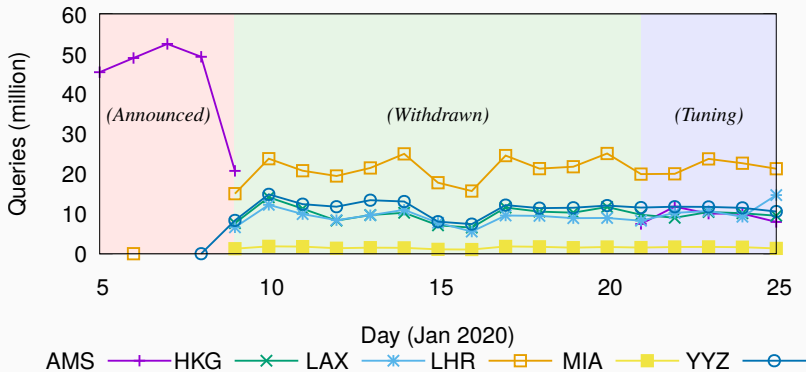


Figure 4: IPv4: Queries and Experiments from Google (AS15169) to Server A

Solution: Depolarizing traffic from Google (BGP)

- We fixed the issue with BGP manipulations
- Median latency: from 100ms to 10ms.

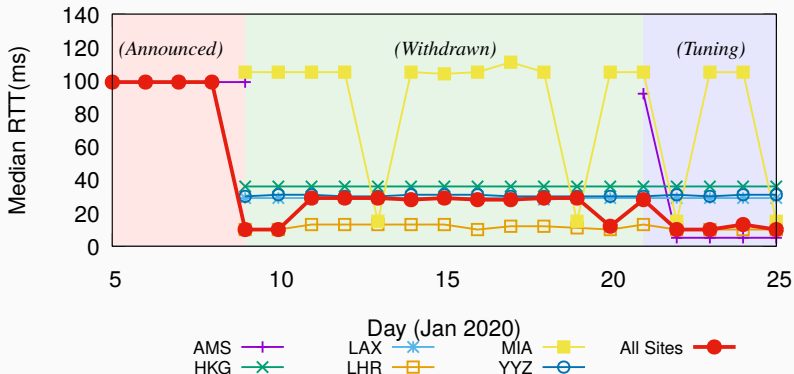


Figure 5: IPv4: Queries and Experiments from Google (AS15169) to Server A

Solution: Depolarizing for Microsoft

- We fixed the issue with route withdraw

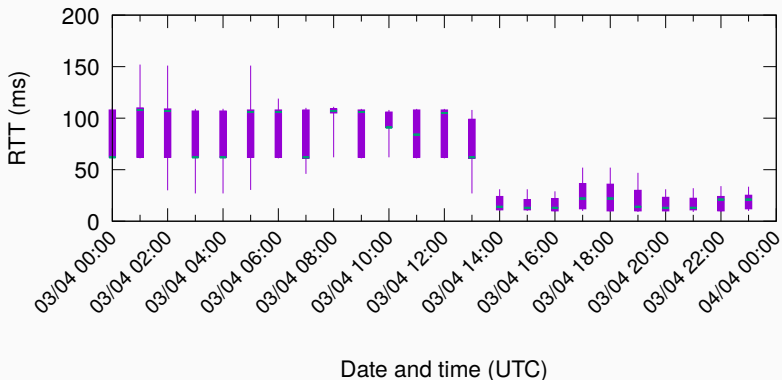


Figure 6: .n1 Anycast A and Microsoft (IPv4): RTT before and after depolarization.

Near-real time Anycast Monitoring: Anteater

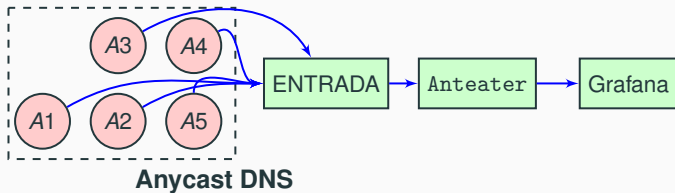
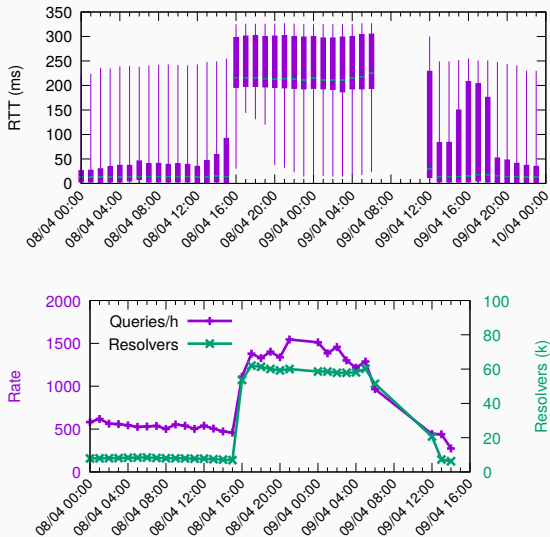


Figure 7: DNS/TCP RTT near real-time monitoring at .nl

Near-real time Anycast Monitoring: Anteater

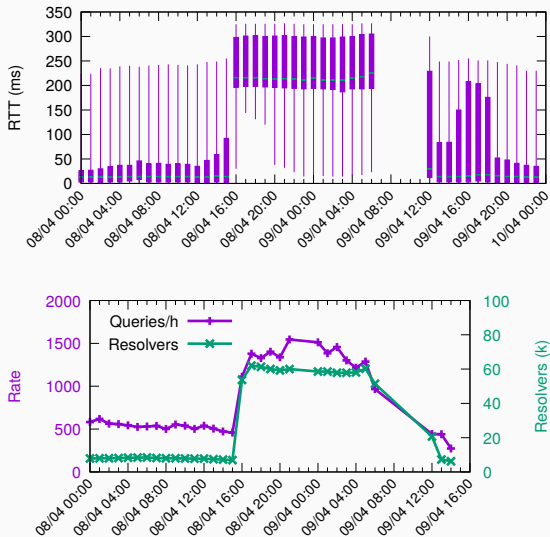


Anteater: detecting routing leaks



EU Traffic went to AUS, tier1 propagated SYD announcements

Anteater: detecting routing leaks



EU Traffic went to AUS, tier1 propagated SYD announcements

Summary

- DNS/RTT are useful for Anycast Engineering
- We show how to prioritize analysis (per site, per client)
- We use our approach in three anycast Services (Services A and B, and B-Root)
- We document Anycast Polarization, and shed latency in 90ms
- Other types of issues covered as well
- ENTRADA, open-source, automatically measures it
- We've been using it for over 1.5 year at SIDN (.nl)
- Tech report:

<https://www.isi.edu/~johnh/PAPERS/Moura20a.html>

- [1] DE VRIES, W. B., DE O. SCHMIDT, R., HARAKER, W., HEIDEMANN, J., DE BOER, P.-T., AND PRAS, A.

Verfploeter: Broad and load-aware anycast mapping.

In *Proceedings of the ACM Internet Measurement Conference* (London, UK, 2017).

- [2] HOE, J. C.

Improving the start-up behavior of a congestion control scheme for tcp.

In *Proceedings of the ACM SIGCOMM Conference* (Stanford, CA, Aug. 1996), ACM, pp. 270–280.

[3] MACIEJ ANDZINSKI .

Passive analysis of DNS server reachability.

https://www.nic.cz/files/nic/IT_19/prezentace/12_andzinski.pdf, 11 2019.

[4] MOURA, G. C. M., HEIDEMANN, J., HARDAKER, W., BULTEN, J., CERON, J., AND HESSELMAN, C.

Old but gold: Prospecting TCP to engineer DNS anycast (extended).

Tech. Rep. ISI-TR-740, USC/Information Sciences Institute, June 2020.

- [5] SCHLINKER, B., CUNHA, I., CHIU, Y.-C., SUNDARESAN, S., AND KATZ-BASSETT, E.

Internet Performance from Facebook's Edge.

In Proceedings of the Internet Measurement Conference (New York, NY, USA, 2019), IMC '19, ACM, pp. 179–194.

- [6] ZHU, P., MAN, K., WANG, Z., QIAN, Z., ENSAFI, R., HALDERMAN, J. A., AND DUAN, H.

Characterizing transnational internet performance and the great bottleneck of china.

Proc. ACM Meas. Anal. Comput. Syst. 4, 1 (May 2020).