

A Perspective Study of DNS Queries for Non-Existent Top-Level Domains

A Technical Report from the Name Collision Analysis Project (NCAP)

26 January 2022

Preface

This is a report to the ICANN Board, the ICANN organization (ICANN org), the ICANN community from the Name Collision Analysis Project.

Table of Contents

Executive Summary	4
Terminology	5
Background	5
Studies	6
Study 1: Root Server Identifier Comparison	7
Data	7
Notable Limitations of the Data	7
Measurements	7
Query Volume per RSI	8
Top Talkers	8
Geographic Relevance	13
Non-Existent TLDs with Highest Query Count	16
Study 1 Key Observations:	19
Study 2: Public Recursive Resolver and Root Comparison	20
Data	20
Notable Limitations of the Data	20
Measurements	21
Total Query Volume per TLD Distribution	21
A and J Root Servers Compared to a PRR Using Total Query Volume per TLD Ranking as a Function	22
A, J, and L Root Servers Compared To Public Recursive Using Distinct Source IPs per TLD Ranking Function	25
Study 2 Key Observations:	27
Key Findings	27
Annex 1: Statistical Methods	29
Jaccard Index	29
Gini Coefficient	29
Annex 2: Non-Top-Talkers	30

Executive Summary

As part of the goals and objectives of the Name Collision Analysis Project Study Two, a study was commissioned by the NCAP Discussion Group to better understand within the context of name collisions how representative DNS data is at various points of the DNS hierarchy. The study's main objective is to provide insights and guidance for future examinations of the DNS name collision data that will be used by ICANN for risk analysis and assessments of TLD string applications. This study, referred to as "A Perspective Study of DNS Queries for Non-Existent Top-Level Domains", focuses on two key measurements: (1) comparing traffic received at each root server identifier and (2) comparing traffic received at public recursive resolver(s) and the root server system. The former measurement provides insights into the ability of name collision DNS data to be collected and analyzed by using a single or subset of root servers, while the latter provides insights into the completeness of DNS measurements taken only at the root by examining DNS name collision traffic at the recursive layer of the DNS hierarchy. The findings from this study indicate that measurements taken from any single root server identifier are largely representative of what is observed at the whole of the root server system; however, there are notable differences in DNS traffic observed by recursive resolvers and at the root server system. These findings are significant in terms of how future guidance and advice should be applied to name collision risk assessments.

Terminology

- Root Server Identifier (RSI)¹ - is the DNS name associated with a root server operator that appears in the root zone and root hints file. For example, c.root-servers.net is the root server identifier associated with the root server managed by Cogent at the time this document was published.
- Day-In-The-Life (DITL)² - a large-scale data collection project undertaken every year since 2006. This data has historically been the primary measurement asset for name collision studies.
- Delegation³ - The introduction of a TLD into the Internet's authoritative database, known as the Root Zone.

Background

Preceding the round of new gTLDs in 2012, numerous studies were conducted by JAS Global Advisors, Interisle, ICANN, Verisign, and other researchers using various types of DNS data to measure and assess name collision risks⁴. The primary data used was root server DNS traffic data collected by DNS-OARC's DITL project. The DITL data provided the most complete view/collection of DNS traffic to the root servers despite being limited to a small number of days per year. The DITL data helped form the guidance issued by JAS Global Advisors to assess the risk of the applied-for TLDs based on query volume and other metrics observed at the root.

The next round of new gTLD applications will require name collision risk assessments by the applicants and ICANN. However, DITL and root data may not be adequate or even available to assure accurate and complete assessments due to anonymization efforts by root server operators and general changes within the DNS ecosystem that raise concerns about availability and accuracy. This study aims to understand the distribution of DNS name collision traffic throughout the DNS hierarchy and provide insights into where and how DNS data can be collected and assessed and data limitations within the context of name collisions.

¹ Root Server System Advisory Committee (RSSAC), "RSSAC 026 – RSSAC Lexicon" 14 March 2017, <https://www.icann.org/en/system/files/files/rssac-026-14mar17-en.pdf>.

² DNS Operations, Analysis, and Research Center (DNS-OARC), "Day In The Life of the Internet (DITL)," accessed 26 January 2022, <https://www.dns-oarc.net/>.

³ ICANN, New Generic Top-Level Domains, "Delegated Strings," accessed 26 January 2022, <https://newgtlds.icann.org/en/program-status/delegated-strings>.

⁴ ICANN, "Mitigating the Risk of DNS Namespace Collisions Final Report by JAS Global Advisors," ICANN Announcements, 30 November 2015, <https://www.icann.org/en/announcements/details/mitigating-the-risk-of-dns-namespace-collisions-final-report-by-jas-global-advisors-30-11-2015-en>.

Studies

“A Perspective Study of DNS Queries for Non-Existent Top-Level Domains” consists of two main studies: the comparison of traffic among RSIs and the comparison of name collision traffic observed by root and recursive resolvers. Together these two studies help provide insights into how risk assessments of name collisions should be evaluated based on the availability of DNS traffic data.

RSI Comparison: This study uses root server data collected by the 2020 DNS-OARC DITL to compare recursive resolver traffic received by each RSI. Using the source IP address and its number of queries issued, various measurements comparing the overlap and distribution of these sources to the various root server identifiers are calculated. Further analysis looking at A and J root server traffic data compares the top name collision strings based on two previously established critical diagnostic measurements, query volume and source diversity.

Public Recursive Resolver and Root Comparison: This study aims to examine the relatively opaque and widely inaccessible data for name collision analysis traffic to public open recursive resolvers. This study uses root server data at A, J, and L root servers and a single public open recursive resolver to compare the top name collision strings based on two critical diagnostic measurements established in the NCAP case study for .CORP, .HOME, and .MAIL: query volume and source IP address diversity.

Study 1: Root Server Identifier Comparison

Data

In order to compare RSIs, data was sourced from the DNS-OARC DITL 2020. At the time, the data for 2021 was not yet available. The 2020 DITL data was collected from May 5th to the 7th, 2020. The contributing root server identifiers were A, B, C, D, E, F, H, I, J, K, L, M. Note that B, E, and F data files are very “small” in terms of data stored in the 2020 DITL fileshare.

Processing DITL data can be cumbersome and computationally expensive (both in time and resources). Fortunately, this study was able to primarily rely on a derived aggregated data set previously generated by Casey Deccio, who was hired to serve as the NCAP technical investigator for name collision reports sent to ICANN. The data included the following fields (note: the aggregation ignored TCP queries):

- IP Address
- Number of queries
- Number of priming queries (i.e., NS . queries)
- Root letter

Notable Limitations of the Data

Two of the root server identifiers, L and I, anonymize the source IP address. Unfortunately, this limits the ability to use those RSI’s data. For example, the I-root data actually takes the source IP address and anonymizes all of them into the 10.0.0.0/8 IP address space. L-root anonymized the source IP address across the whole IPv4 range. IP anonymization does not work for most of our measurements, thus both I and L RSI’s were excluded from this study’s measurements. Furthermore, the size and completeness of B, E, and F RSI’s data was inadequate for this study’s required measurements and these RSIs were also excluded. These exclusions reduced the original twelve RSIs down to seven.

Measurements

The following twelve measurements⁵ were taken against the data:

1. Query volume per RSI
2. Unique source IP address at each RSI
3. Distribution of query volume per source IP to all of the root server system
4. Identifying top talkers⁶ that constitute a large percentage of overall traffic
5. Measuring overlap of top talkers at each RSI

⁵ Further investigation could explore the number of sites per RSI and other ratios related to queries and source IP addresses.

⁶ Top talkers are recursive resolvers that issue the largest amount of DNS queries to the RSS.

6. Comparing the set of IPs at each RSI to the other RSIs
7. How many RSIs must be analyzed to reach 100% of the top talkers
8. How many RSIs does a typical top talker IP query
9. Are there any geospatial outliers within the top talker set of IPs
10. How evenly do top talkers distribute the query volume over RSIs
11. Is there a geographical bias for various countries to favor a subset of RSIs
12. What variation exists in the Top-N non-existent TLDs per RSI

Query Volume per RSI

The first baseline comparison of RSI traffic is the number of queries each receives. As shown in Figure 1 below, the number of queries received at each RSI varies; accordingly, this measurement provided insights into data collection issues with B, E, and F and why they were ultimately excluded from further analysis. The total query count for these root server letters (B, E, and F) combined was less than 0.8% of the entire 2020 DITL.

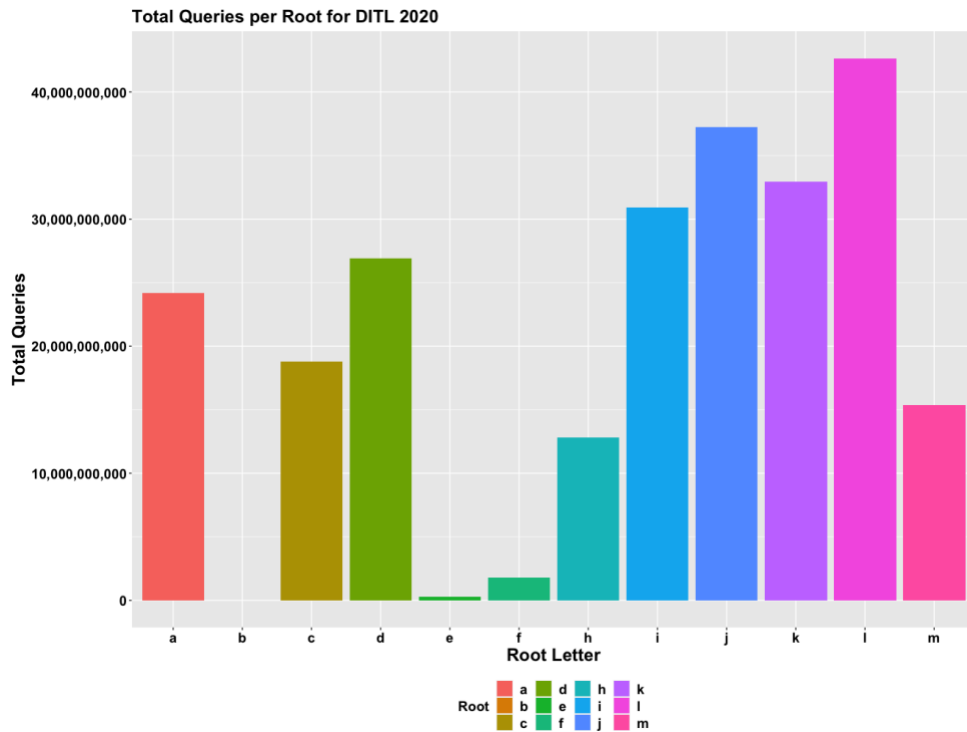


Figure 1 - Query Volume per RSI

Top Talkers

A second fundamental measurement was to understand the number of unique IP addresses seen at each RSI. This is useful to understand if we should expect IP affinities, which would have a direct impact on

any future name collision analysis that uses a subset of RSIs. Figure 2 below shows the number of unique IPv4 and IPv6 addresses seen at each RSI. That distribution, on the included RSIs, is relatively even (median = 8.52M, mean = 7.61M, standard deviation = 1.8M). In total, 15.51M unique IPv4 and 1.56M IPv6 addresses were observed.

The focus of this study is to understand how similar the sets of source IP addresses and queries for non-existent TLDs are across different RSIs. There are numerous similarity measurement approaches but a simple and often reliable measurement is the Jaccard index that is a statistic used in understanding the similarities between sample sets. Already from Figure 2 we can see that any type of set measure of unique IP sources will have significant variance. In order to support the focus of this study, a smaller number of IP addresses, that are representative of the entire RSS, is needed for the use of set similarity. To that end, understanding the profile of how many queries each IP sent is needed.

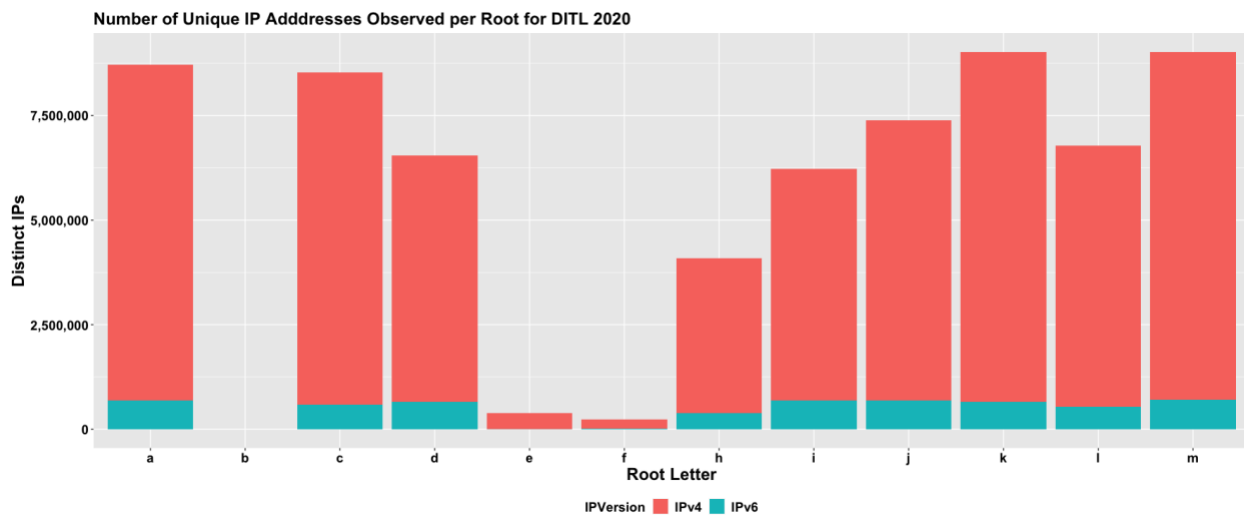


Figure 2 - Unique IP Addresses Observed per RSI

Query volume from each IP is typically not equally shared across the RSIs. To understand the query volume distribution over the set of IP addresses observed in the 2020 DITL collection, a cumulative distribution measurement was made by ranking IP addresses in ascending order by the number of total queries that IP sent to the RSS. Figure 3 below depicts this distribution measurement relative to the total percentage of IP addresses observed during the 2020 DITL. A typical Power Law Distribution⁷ was observed:

- 15% of IP addresses issued only 1 query.
- 27% of IP addresses issued 2 or fewer queries.
- 50% of IP addresses issued 10 or fewer queries.

⁷ Wikipedia contributors, "Power law," *Wikipedia, The Free Encyclopedia*, accessed January 26, 2022, https://en.wikipedia.org/w/index.php?title=Power_law&oldid=1065551786.

- 98% of IP addresses issued 10,000 or fewer queries.

It is unclear as to what those source IP addresses that issue so few queries actually are. A typical recursive resolver with even a minimal amount of users, clients or other stub systems behind it might be expected to generate more than 10 queries to the RSS over a period of 48 hours. It is possible these are odd pieces of software, spoofed source IPs, or a variety of other things.

From Figure 3 we observe that the vast majority of the traffic collected during 2020 DITL come from IP addresses that issue a significant amount of queries. This is further confirmed in Figure 4.

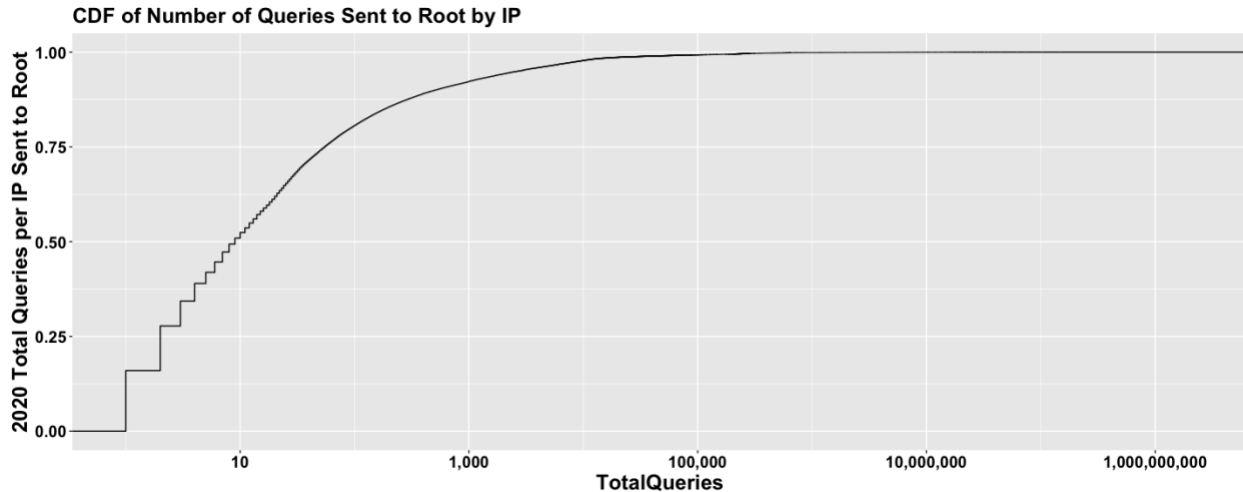


Figure 3 - Cumulative Distribution of the Number of Queries Sent to RSS by IP Address

This insight helps inform traffic comparisons across RSIs. Any measurement of similarity will likely be very skewed by the nature of having so many IP addresses that account for negligible amounts of RSS traffic. Therefore, additional measurements were made to determine top talkers, i.e., those IP addresses that constitute a large percentage of the traffic, and how they are distributed across RSIs (if they are distributed at all). This is important because it will provide us a more consistent and accurate measurement of how RSIs compare to each other based on IP addresses that constitute the majority of the query volume (and accordingly, name collision queries) on the RSS. This measurement does come at the cost of disregarding the longtail of low querying source IPs but will facilitate the intended measurement of this study. Additional analysis would need to be conducted to better quantify any impact thresholding imparts on the findings here within.

Figure 4 below shows a distribution of the number of the top querying IP addresses relative to the total percentage of 2020 DITL queries received. This measurement shows that 90% of the total 2020 DITL can be represented by only looking at the top 115K IPs. Likewise, 95% of the total 2020 DITL can be represented by the top 250K IPs. The remaining 5% of query volume observed in 2020 DITL is distributed in the long tail of millions of IPs. From this point on, our analysis will be performed on the top

115K (0.67%) of IP addresses, which we hereafter refer to as the “top-talkers”. Additional analysis of the long tail of IP addresses and extending the top-talker threshold is contained in Annex 2.

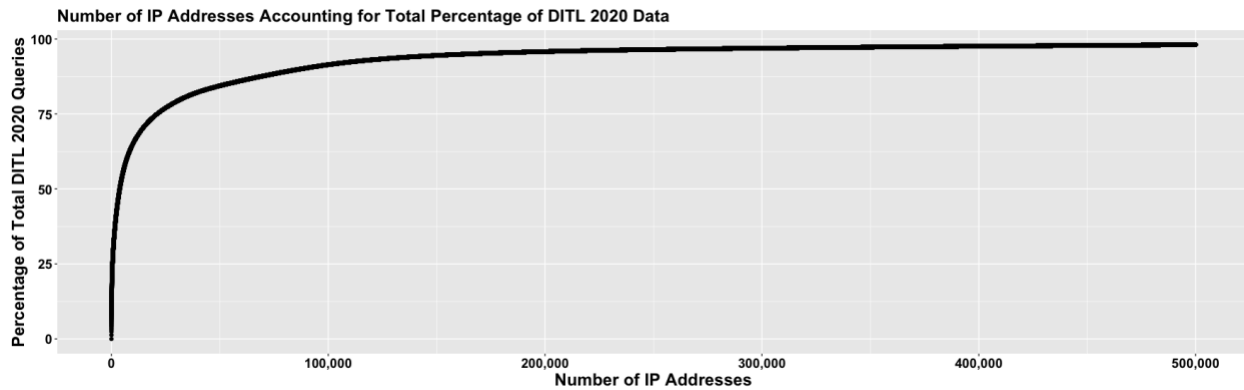


Figure 4 - Number of IP Addresses Accounting for Total Percentage of DITL 2020

The next measurement was tailored to better understand how these top talking IP addresses are distributed over the RSIs. Figure 5 below shows the percentage of the top talking IPs observed at a given RSI. On average, each RSI observed 96% of the top talkers that account for 90% of total traffic. That percentage drops to 94% when using the 95th percentile top talkers. Based on these findings, only the 90th percentile top talkers were used for the remaining measurements in this study. Similarly, Figure 6 shows a histogram of the number of RSIs queried by the 90th percentile top talkers. This distribution indicates the vast majority of these IP addresses (89%) are seen by all RSIs - a key indicator that any RSI may be representative of the general RSS.

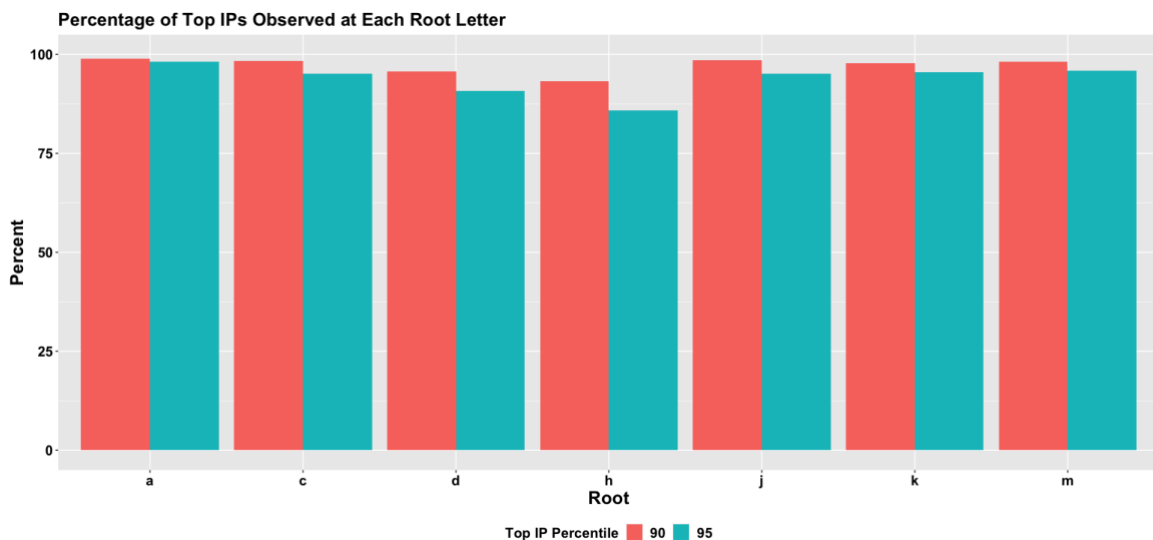


Figure 5 - Percentage of Top IPs Observed at each RSI

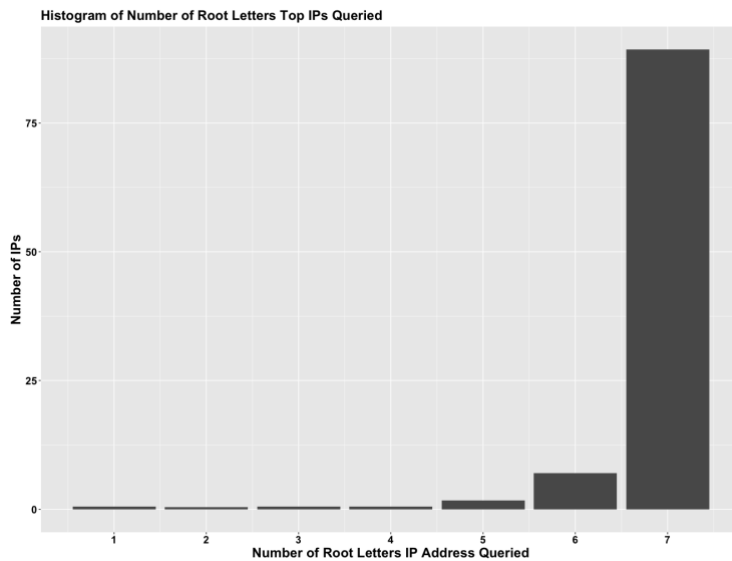


Figure 6 - Histogram of Number of Root Letters Top IPs Queried

A more detailed measurement of how top talking sources are observed at any two RSIs is depicted in Figure 7 below. The figure shows one-half of a similarity matrix that utilizes the Jaccard index, a similarity measurement that is further clarified in the Appendix, to measure the amount of overlap between two RSIs and the top talkers. From a source diversity perspective, any root letter, in general, sees a very high percentage of top talkers compared to any other root. On average 96% of top talkers are observed at any two roots.

Top talkers are widely seen at all root letters. Data from any combination of three RSIs will include 99.5% of top talkers, though all RSIs must be included to reach 100% of top talkers.

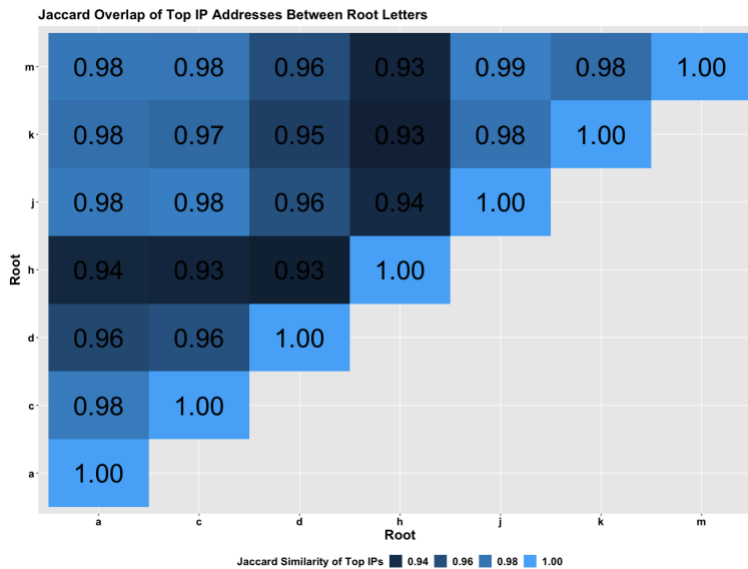


Figure 7 - Jaccard Overlap of Top IP Addresses Between RSIs

Geographic Relevance

The preceding measurements provided insights into the distribution of top talkers over the RSIs. The following measurements continue to compare the distribution of these top talkers from spatial and geographic means. Spatial representation of the IPv4 space is achieved via the use of a tool called IPv4 Heatmap.

IPv4 Heatmap is a program⁸ that generates a map of IPv4 address data using a space-filling Hilbert Curve. Each pixel in the image represents a single /24 network and is assigned one of 256 colors. Pixel colors range from blue (1 host) to red (256 hosts), while black represents no data (0 hosts). Figure 8 below is an example of how an IPv4 spatial distribution can be visualized.

⁸ Duane Wessels, ipv4-heatmap, last updated 1 March 2021, <https://github.com/measurement-factory/ipv4-heatmap>.

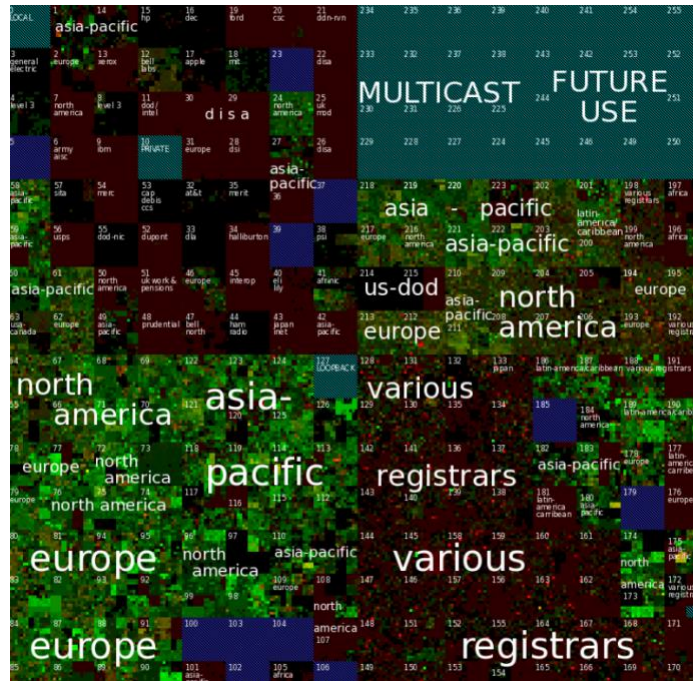


Figure 8 - Example of Hilbert Curve IPv4 Visualization

After exclusion of various RSIs due to IP anonymization or minimal data, seven RSIs were bucketed into the 256 color range by increments of $256/7$. Each of those colors was then assigned in increasing order to represent the number of RSIs an individual IP address queried during the 2020 DITL. Blue dots are top talker IP addresses that only queried 1 root while red dots represent top talker IP addresses that queried all 7 RSIs. As seen in Figure 9 below, there is a heterogeneous distribution across IPv4 address space. There are some notable exceptions in which several netblocks have a concentration. A few interesting groups of IP addresses, which queried only one or a few RSIs, appear in a small number of netblocks (e.g. 178.0.0.0, 172.0.0.0, etc.). It remains unclear as to what those resolvers are or their purpose without a more thorough analysis of their specific queries. Overall, these measurements indicate no large biases of source IP addresses showing specific RSI affinity.

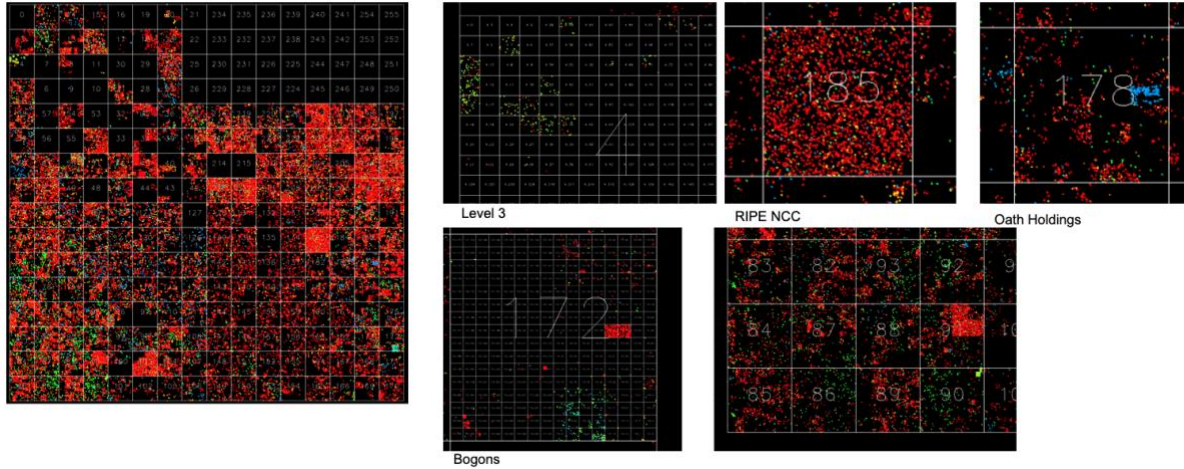


Figure 9 - IPv4 Spatial Distribution of Top Talker IP Addresses

Expanding into geospatial measurements, we next used the Gini coefficient to measure how much inequality a top talker IP has for the distribution of root letters. We also geolocated top talker IP addresses to determine country-root inequality. The Gini coefficient measures the inequality among values of a frequency distribution (for example, levels of income). Gini values are bound between 0 and 1, in which a value of 0 would indicate the values are evenly distributed and 1 would indicate complete inequality.

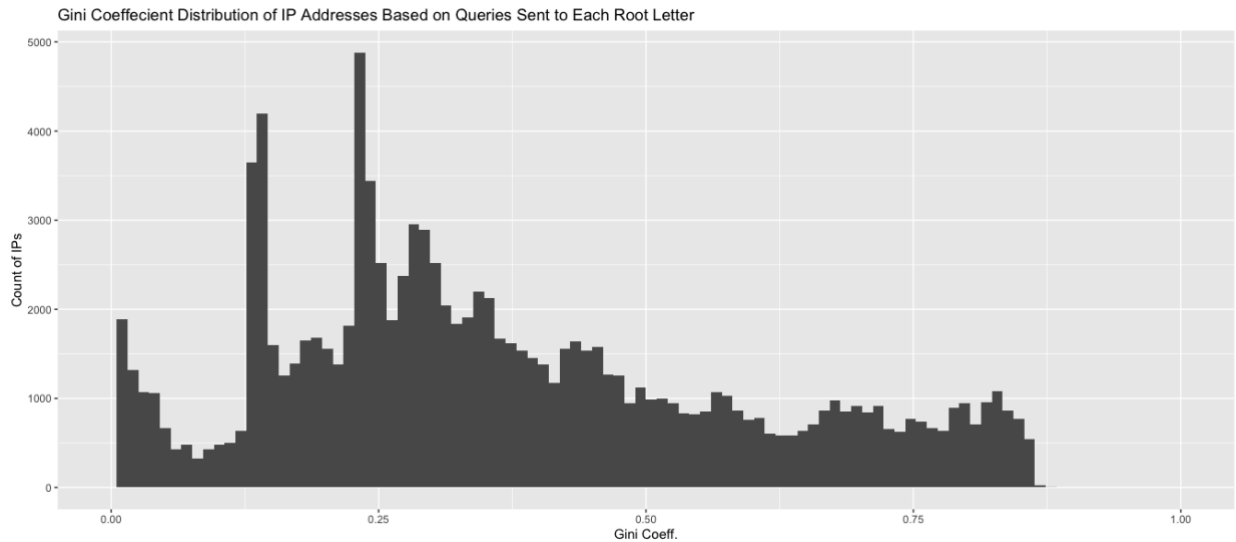


Figure 10 - Gini Coefficient Histogram of Top Talker IPs

Using the 90th percentile top talker IPs, each IP address Gini coefficient was calculated based on the number of queries the IP sent to each of the seven RSIs. Figure 10 above shows the distribution of those 115K Gini coefficients. While it appears to be multi-modal, the majority of the IP addresses resulted in values nearer to zero, indicating that these top talkers are distributed their query load over all of the

participating RSIs. Likewise, the top talker IPs were mapped to countries using the Maxmind GeoIP database and the country traffic for each RSI was calculated. Figure 11 shows a geographical plot coloring in which the shading of the country is based on its Gini coefficient.

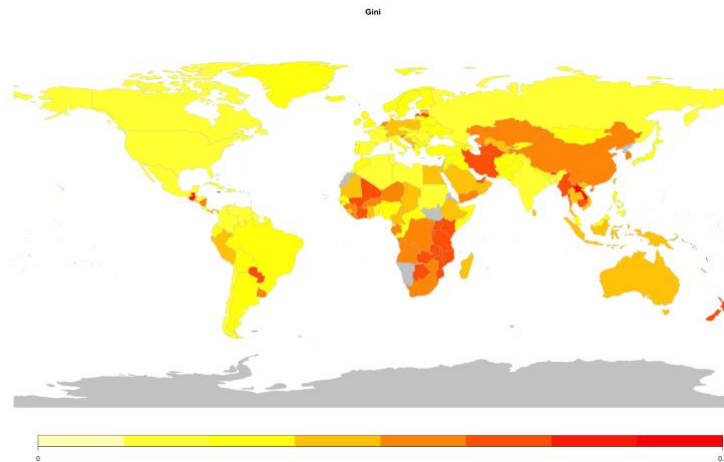


Figure 11 - Country to RSI Traffic Distribution

The overall per country Gini was an average of 0.32. Certain regions of Africa, Asia, and island countries have elevated Gini values and stronger affinities to certain root letters (it is expected this is likely due to placement/peering). An example of this bias/affinity can be seen in Figure 12 below, in which 2.9% of K-root traffic originated from Iran while other RSIs observed rates closer to 0.3%. Overall, this measurement helps confirm there is no large geographical bias of RSIs.

```
> country_root[CC == "IR",]
  CC Root CountryRootTotal RootTotal CountryRootPercent
1: IR  a      110692882 21548638126      0.51368853
2: IR  c       67909454 17360762933      0.39116630
3: IR  d      125520262 25331103790      0.49551833
4: IR  h        4251817 12051048679      0.03528172
5: IR  j        59228457 34957035326      0.16943215
6: IR  k       910077945 30783508659      2.95638147
7: IR  m        50587079 13627881554      0.37120281
> |
```

Figure 12 - Gini Coefficient for RSIs in Iran

Non-Existent TLDs with Highest Query Count

The previous analysis shows with a high level of confidence that traffic to any RSI is largely representative of what any other RSI may be observing at a particular moment in time. This is important because it provides some confidence that future name collision measurements could be taken by any RSI without requiring an RSS-wide collection. In addition to looking at how representative traffic is received from querying recursive resolvers to RSIs, the following measurements will look at the similarity of the

names. Specifically, the following figures and tables will examine what variation exists in the top N non-existent TLDs on a per RSI basis.

In order to understand how top non-existent TLDs compare at each RSI, the top 10,000 TLDs based on query volume were compared at A and J RSIs using the 2020 DITL data. If a TLD was observed at one RSI but not at the other RSI, a rank value of zero was associated with that TLD at the other TLD. Thus any TLD depicted in Figure 13 in which the dot is at $x=0$ or $y=0$ means that particular TLD was not seen in the top 10,000 by the other RSI. Figure 13 shows that TLDs under rank ~1,000 are often (812 of 1,000 TLDs with a correlation coefficient of 0.64) seen at the other root; however, as the rank increases (e.g., the total traffic volume decreases), the correlation of a TLD's rank at one RSI diminishes (TLDs above the 1,000 rank have a correlation coefficient of 0.004).

These initial findings suggest that the most queried for non-existent TLD strings will be seen at another RSI proportional in terms of query volume rank. However, as the amount of query volume decreases for a non-existent TLD, the likelihood in which that string is seen at another RSI's top-N list decreases. This helps support and inform the DNS community that the publication of top-N strings could be beneficial to future TLD applicants. It also further illustrates that not all name collision traffic is evenly distributed in which an accurate assessment of risk can be done via RSS data alone.

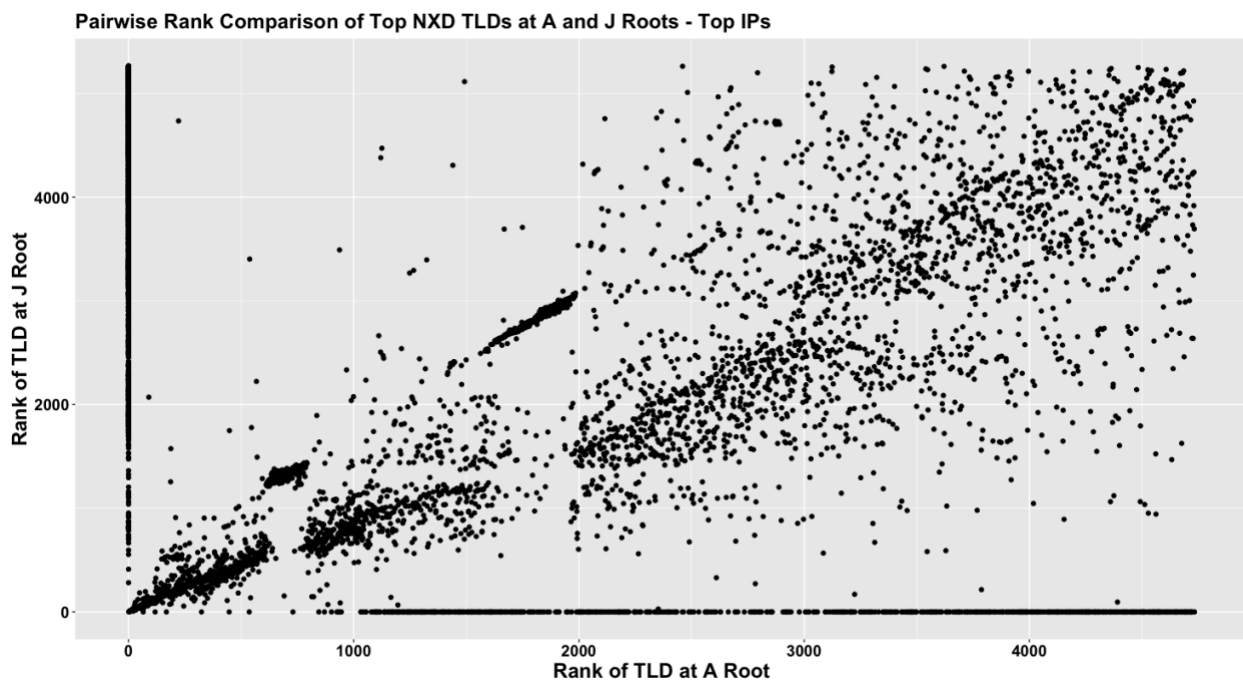


Figure 13 - Rank Comparison of A and J Top NXD TLDs

Figure 14 shows a more focused scatterplot depiction of the top 1,000 TLDs. This data was measured on November 15, 2021 at A and J RSIs. The TLD strings were also required to match the regular expression

[a-z0-9]{3,63}⁹. Again a stronger rank correlation is expressed at low ranks (correlation coefficient of 0.6 less than rank 200) and diminishes as the ranks approach 1,000 (correlation coefficient of 0.1 greater than rank 200)

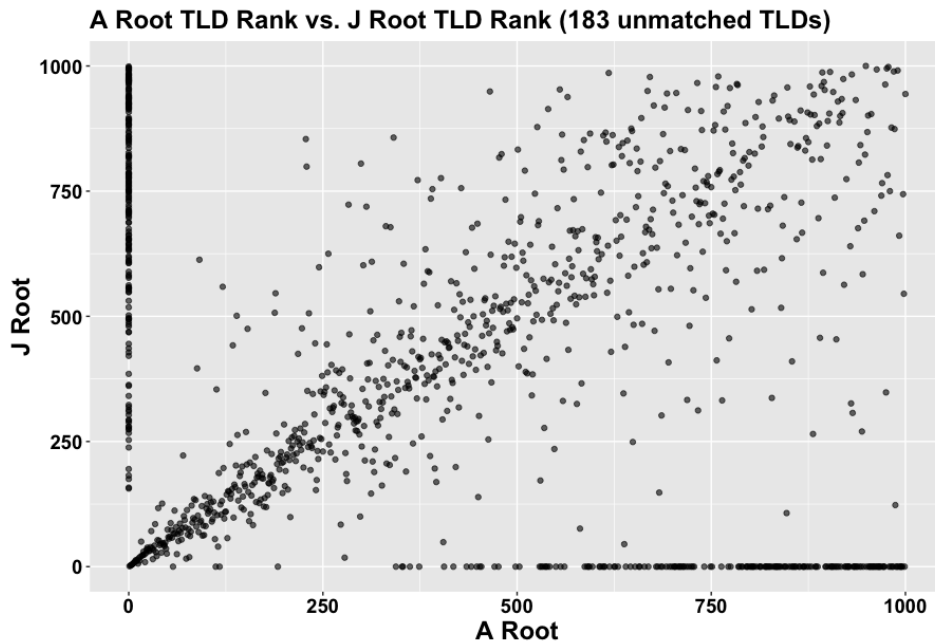


Figure 14 - Rank Comparison of A and J Top NXD TLDs

Figure 15 is a standard Venn diagram showing the overlap of the Top non-existent TLDs plotted in Figure 14. The overlap of A and J RSI Top non-existent TLDs was 817 strings, with 183 strings only being observed at one of the two RSIs. Again, this supports our earlier findings that any RSI will likely be representative of major name collision issues expressed in the whole of the RSS.

⁹ This regular expression was used because it will loosely match some of the ASCII label technical and policy string requirements set out in section 2.2.1.3.2 (“String Requirements”) of the evaluation procedures in the gTLD Applicant Guidebook (v. 2012-06-04) for DNS Stability. Many of the top non-existent TLD strings contain non-alphanumeric characters that do not match this requirement.

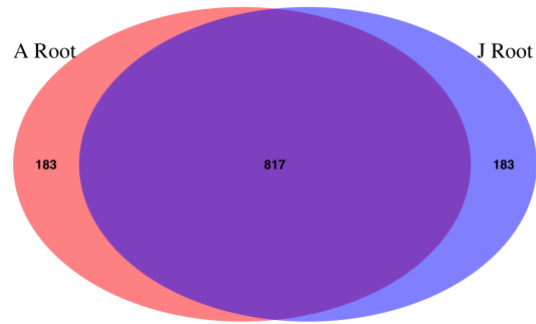


Figure 15 - Venn Diagram of TLD Overlap between A and J RSIs

Study 1 Key Observations:

In the analysis that we have conducted, we have observed the following:

- Ninety percent of RSS queries are sent from a relatively small set of IP addresses (115K).
- Top-talking IPs are broadly seen at all root letters. Queries from 89% of top-talking IP addresses are observed at all root servers.
- Some geographic affinity/preference to certain root letters does occur.
- Initial research using two RSIs show top non-existent TLD strings between letters appear to generally correlate for the top 1K strings and that lower volume non-existent strings appear to be more root letter dependent.

Study 2: Public Recursive Resolver and Root Comparison

Since the 2012 round of TLD delegations, several new technologies and recommended best practices within the DNS ecosystem now have a significant impact on the volume and fidelity of DNS queries observed at name servers in the DNS hierarchy. The emergence of popular open recursive resolvers has also dramatically shaped the DNS ecosystem since the new gTLD delegations. These recursive services may provide a richer and more complete understanding of name collisions if they can be utilized for analysis. Therefore Study 2 was designed to investigate the differences of name collision strings at the RSS level as well as the public recursive resolver level.

Data

In order to understand how DNS traffic compares at various layers of the DNS hierarchy, query data sent to several root server identifiers and one public recursive resolver were collected and measured in such a way that would facilitate the examination of top non-existent TLD observed in queries. The data was measured using two sorting functions that reflect the importance of our critical diagnostic measurements: (1) Query Volume and (2) IP Address diversity. Two lists of the top 1000 non-existent TLDs matching the regular expression `[a-z0-9]{3,63}` were generated based on the two sorting functions. The resulting aggregated data was used to measure how recursive and root server query volume compare by examining rank ordering as well as general TLD string overlap.

Notable Limitations of the Data

While concerted efforts were made to obtain recursive resolver data from numerous sources, only one recursive resolver operator provided the data. The limiting factor appears to be data privacy concerns. To that end, the recursive resolver that did provide the data will not be identified and herein simply referred to as the “public recursive resolver” (PRR). Without obtaining data from other public recursive resolvers, it is unclear how each recursive resolver compares to another. It is likely due to their underlying user-base, deployment size, and internal DNS protocol optimizations, that each recursive resolver represents a unique vantage point of the DNS; however, without additional data this will remain only a hypothesis. The measurements presented in this study, while only looking at one PRR, do provide a novel and previously unknown understanding of name collisions via passive DNS telemetry data used for quantifying and assessing name collision risks at multiple collection points within the DNS hierarchy.

Measurements

The following five measurements were conducted against the data:

1. Query volume distribution of RSIs and the PRR
2. Rank correlation between RSI and PRR based on query volume
3. String overlap between RSI and PRR based on query volume
4. Rank correlation between RSI and PRR based on source diversity
5. String overlap between RSI and PRR based on source diversity

Total Query Volume per TLD Distribution

A baseline measurement comparing query volume of the top 1,000 non-existent TLDs at two RSIs, A and J roots, and the PRR is depicted in Figure 16 below. The distributions appear similar in nature, forming a power-law distribution in which the top non-existent TLDs express query volumes that are several magnitudes higher than the other TLDs. All three distributions seem to “flatten out” into the long tail distribution after the top 50 TLDs.

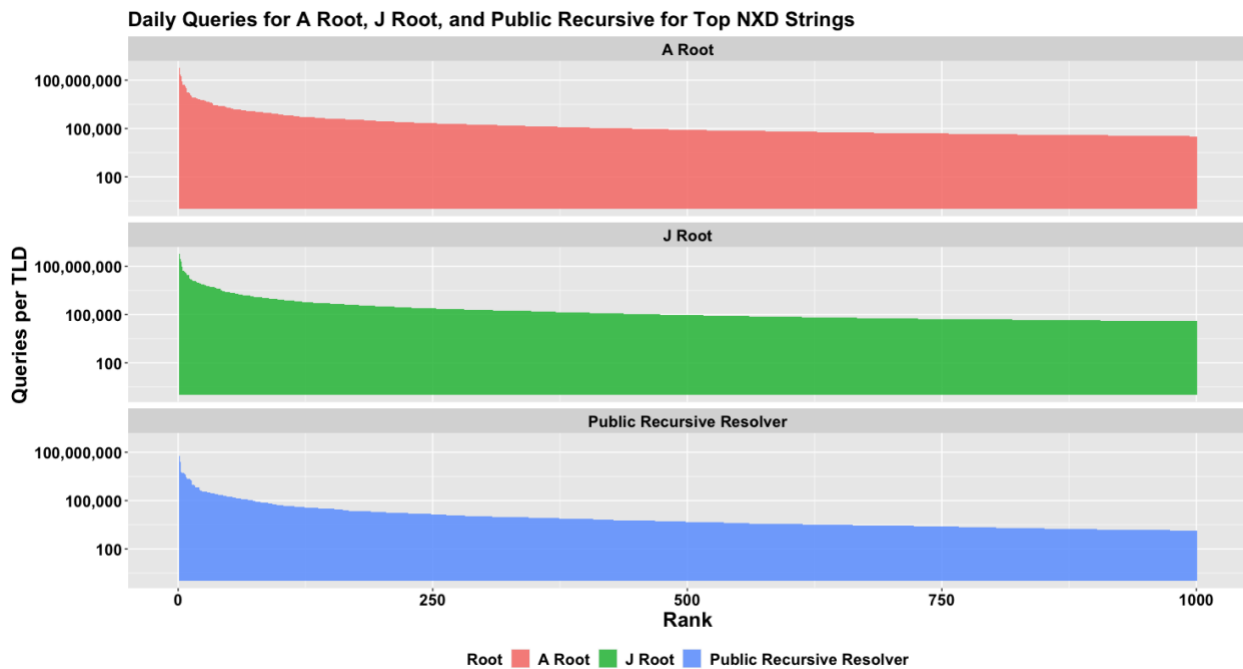


Figure 16 - Daily Queries for A and J RSIs and the PRR for Top NXD TLDs

A and J Root Servers Compared to a PRR Using Total Query Volume per TLD Ranking as a Function

While the initial query volume distribution shown in Figure 16 may have shown some similarities, no other strong similarities were found between the RSIs and the PRR data. Figure 17 below shows a simple scatter plot of the top RSI TLD rankings vs. those of the PRR. Unlike the rank scatter plots comparing top RSI TLD rankings relative to another RSI, the RSI to PRR plot shows no correlation between the two DNS data sets (e.g., there is no “diagonal” line with a slope of ~ 1).

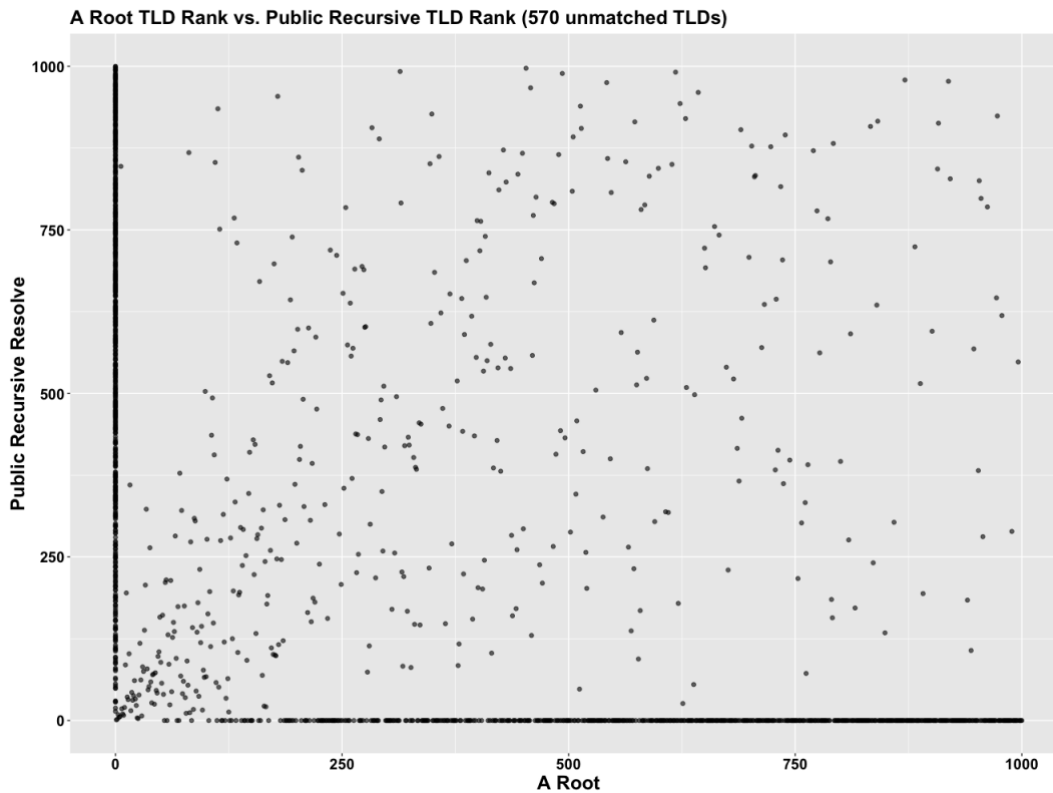


Figure 17 - Rank Correlation of Top TLDs at A Root and Public Resolver based on Query Vol.

This lack of correlation shown in Figure 17 is better explained by looking at the Venn diagram that examines the set overlap of the top 1,000 non-existent TLDs. Only 430 strings were both observed at the RSI and the PRR. This is significantly different from the overlap previously seen between RSIs in which ~ 800 of the strings overlap.



Figure 18 - Venn Diagram showing TLD overlap of A Root and PRR based on Query Vol.

Figure 19 below is another examination of a ranking scatter plot at a second RSI. Again no correlation is observed between the RSI and the PRR. This is again reconfirmed by the Venn diagram in Figure 20, in which only 417 of the top non-existent TLDs were observed by both the RSI and the PRR.

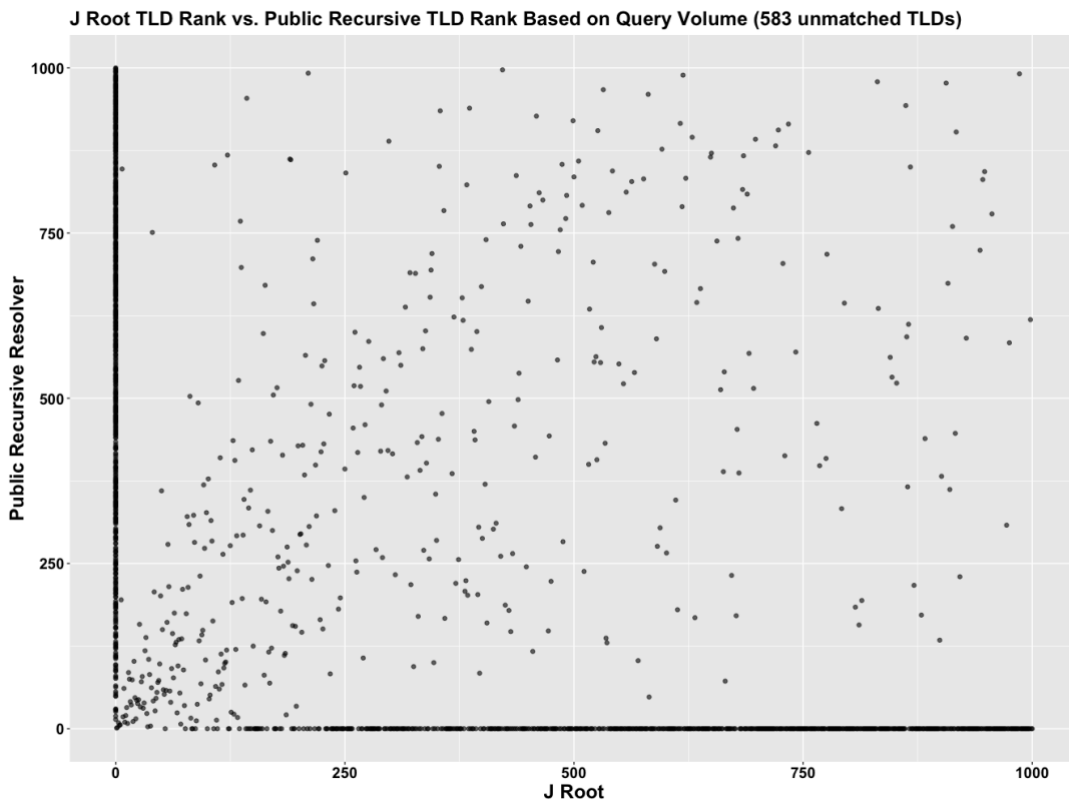


Figure 19 - Rank Correlation of Top TLDs at A Root and the PRR based on Query Volume

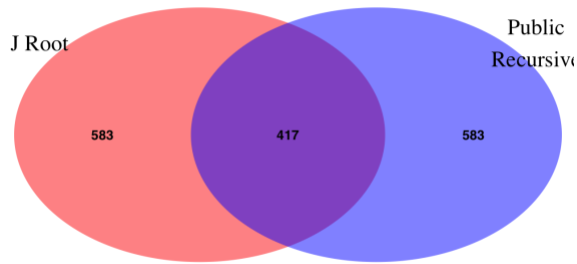


Figure 20 - Venn Diagram showing TLD overlap of the J RSI and PRR based on Query Vol.

These initial comparisons of top strings based on query volume observed at RSIs and the PRR reveal there is a significant difference in DNS queries. The small overlap of top strings between the two data sources further suggests that an accurate and complete picture and risk assessment of collision strings is not possible from RSS data alone. Figure 21 below show the top 50 non-existent TLD strings observed at A root and the PRR.

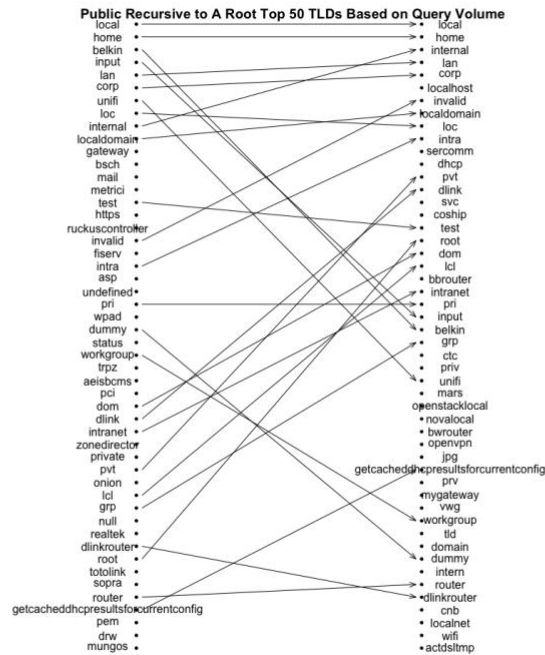


Figure 21 - PRR to A Root Top 50 TLDs Based on Query Volume

A, J, and L Root Servers Compared To Public Recursive Using Distinct Source IPs per TLD Ranking Function

The top 1,000 non-existent TLDs were identified at each of the RSIs and the PRR based on the number of unique IP addresses observed per TLD. An initial measurement looking at TLD string overlap via a Venn diagram is shown in Figure 22 below. The PRR still observed 311 strings which none of the RSIs observed in their top 1,000. This measurement shows greater overlap between RSIs and the PRR than top strings by query volume. However, the significant dissimilarity between the PRR TLDs with the greatest source IP diversity and those of the RSIs means that name collision strings cannot be measured or assessed properly based on only using data from the RSS.

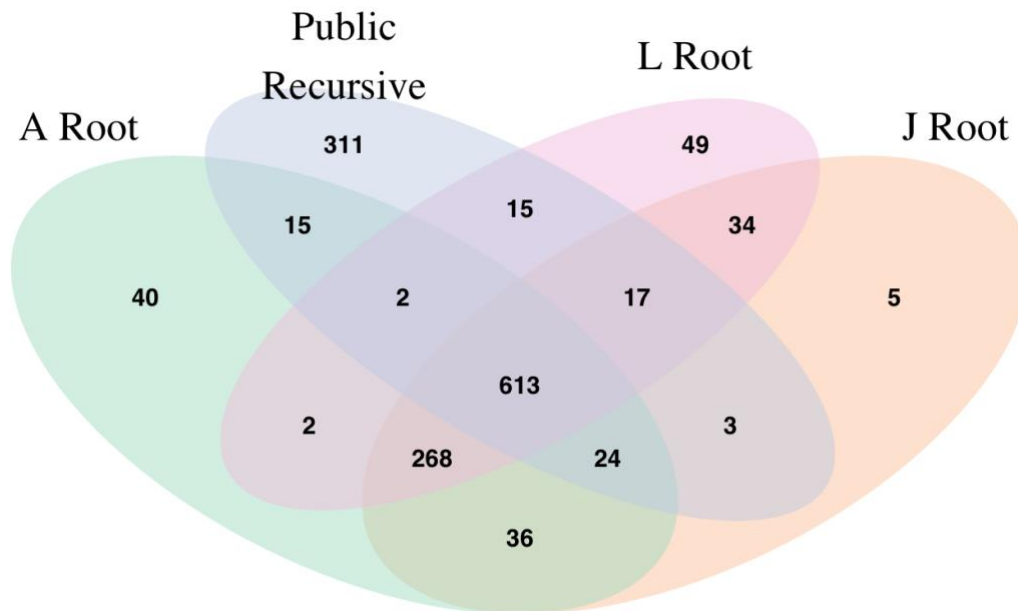


Figure 22 - Venn Diagram Comparing Overlap of Top TLDs at A, J, and L RSIs and PRR based on Source IP Address Diversity

Examining a rank scatter plot between an RSI and the PRR does indicate a slightly better correlation of TLD rankings; however, this correlation appears slightly weak (correlation coefficient of 0.36), at best,

and mainly for the top-ranking strings that had large source diversity measurements (i.e., TLD rankings under 100).

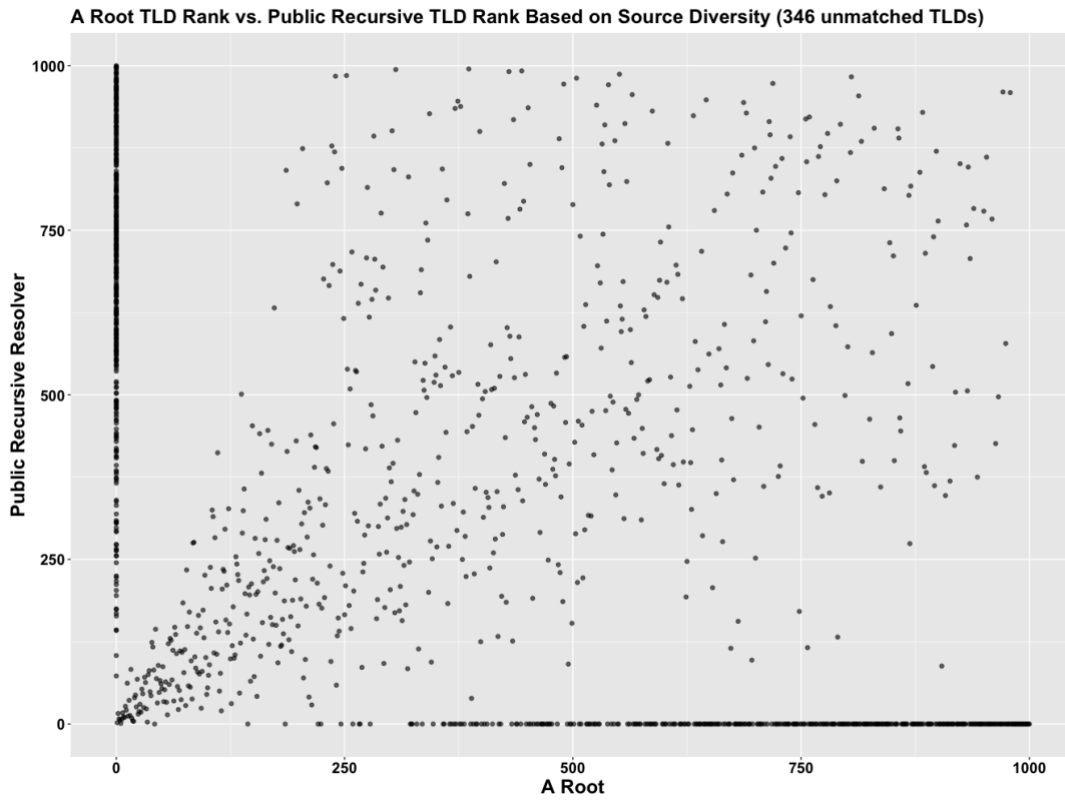


Figure 23 - Rank Correlation of Top TLDs at A Root and PRR based on IP Diversity

Study 2 Key Observations:

In the analysis that we have conducted, we have observed the following:

- Obtaining non-existent TLD PRR data was very difficult and that challenge will likely continue into the future.
- Data provided from the PRR was highly constrained due to data privacy requirements. These constraints only allowed for simple measurements based on non-existent TLD query volume and source diversity.
- Many unknowns about the PRR functionality and deployment makes correlating RSS to PRR measurements difficult. Regardless of these limitations, this is the first analysis of PRR to RSS data in the context of name collision assessments.
- Initial results from one PRR indicate there is a difference in top non-existent TLDs using either query volume or source diversity measurements.
- Many non-existent TLDs observed at the PRR are not in the top non-existent TLDs seen by RSIs based on query volume (~ 40%) and source diversity (~ 30%).
- Significant differences between the PRR non-existent TLDs and those at the RSIs suggest that name collision strings cannot be measured or assessed reliably using data only from the RSS.

Key Findings

The two studies in this analysis provide two key findings that will help the NCAP provide guidance and advice to ICANN as to how future risk assessments of name collision strings should be evaluated.

Finding 1: The IP addresses representing the most queries across the root server system (i.e., the "top talkers") are likely to query any given root server in the course of two days, and queries from those IP addresses for non-existent TLDs are likely to be found with similar prevalence on different RSIs.

Implications:

- Future efforts to analyze name collision risks need not pin themselves to require DITL-like collections. Non-existent DNS queries for top querying and top source diversity TLDs appear to be comparable and representative at any RSI. PRR data further indicates that there is a very different view of the top non-existent TLDs based both on query volume and source diversity.
- ICANN, as the operator for the L RSI, is well-positioned to instrument, collect, analyze, and disseminate name collision measurements to subsequent gTLD applicants both prior to submission and during the application review.

Finding 2: Name collision traffic observed at the root is not sufficiently representative of traffic received at recursive resolvers to guarantee a complete and or accurate representation of a string's potential name collision risks and impacts.

Implications:

- Name collision traffic observed via root server telemetry data should be considered the minimal recorded value.
- Given the current and likely future state of the DNS ecosystem, a complete and accurate risk assessment of a string's name collision potential cannot be determined prior to the string's temporary delegation.

Annex 1: Statistical Methods

Jaccard Index

The Jaccard index, also known as the Jaccard similarity coefficient, is a statistic used for gauging the similarity and diversity of sample sets.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}.$$

Jaccard index measurements are bound between 0 (identical sets) and 1 (completely distinct sets).

Note: This measurement is only accounting for set presence. It does not consider the magnitude/volume of queries sent - it is only if the IP appears in both sets.

Gini Coefficient

$$G = \frac{2 \sum_{i=1}^n i y_i}{n \sum_{i=1}^n y_i} - \frac{n+1}{n}.$$

The Gini coefficient measures the inequality among values of a frequency distribution (for example, levels of income). Gini coefficient measurements are bound between 0 (even distribution) and 1 (completely uneven, e.g., one member receives all traffic).

Annex 2: Non-Top-Talkers

As seen in Figure 3, nearly 92% of all source IP addresses seen during 2020 DITL received less than 1,000 queries over a span of 48 hours. Further investigation was conducted to better understand the long-tail of source IP addresses. Specifically, an investigation was conducted to better understand the domain names these low query volume source IP addresses requested.

Annex 2 Table 1 examines the top queried names at A and J root server during the 2020 DITL based on the set of IP addresses that received only one query over the seven roots analyzed in this study. Nearly 15% (2.7 million) IP addresses of the total 17 million received only one query during 2020 DITL. A and J roots observed 732,711 of those 2.7 million single querying IP addresses. The most frequently observed name was a root priming query at 160,740 unique source IPs accounting for 22% of total single querying IPs. Other notable query names include RFC 8145 trust anchor signals. Two other noteworthy aspects appear in these top single query names: 1.) Most of the names are legitimate and delegated in the global DNS and would result in a rcode of 0. Based on RSSAC002 data, it shows that most traffic at the RSS is for non-existent domains. 2.) Many of these top names appear to be multi-labeled domains under delegated suffixes such as qq.com and in-addr.arpa.

Rank	Qname	Number IPs	Rank	Qname	Number IPs	Rank	Qname	Number IPs
1	.	160740	21	local	3192	41	--	1481
2	sy.eu.angsrvr.com	20494	22	glborigintest.canarytest.net	2652	42	tv	1461
3	im.mielse.com	14582	23	moiaawsorigin.clo.footprintdns.com	2647	43	m.root-servers.net	1383
4	ad.afy11.net	13421	24	qq.com	2545	44	p.qpic.cn	1379
5	chat.grindr.com	11664	25	teredo.ipv6.microsoft.com	2513	45	livem.l.qq.com	1379
6	com	10813	26	support0.bigo.sg	2465	46	tv.l.qq.com	1377
7	gsm1.g4c5j.com	10747	27	home	2051	47	qq.m.cn.miaozhen.com	1364
8	www.google.com	9119	28	_xmpp-client._tcp.wallapop.com	1985	48	vgdt.gting.cn	1363
9	216.58.202.4.in-addr.arpa	9000	29	ns2.dnsimple.com	1933	49	mtcls.qq.com	1362
10	xmpp-prod.monksoftware.it	6845	30	ns4.dnsimple.com	1932	50	dns.weixin.qq.com	1353
11	net	5788	31	api.surfeasy.com	1929	51	appmedia.qq.com	1353
12	_ta-4f66	5341	32	ns1.dnsimple.com	1900	52	cm.l.qq.com	1352
13	_sip._udp.smart.0038.net	4652	33	ns3.dnsimple.com	1865	53	nadia.ns.cloudflare.com	1352
14	ns-635.awsdns-15.net	3997	34	dion.ns.cloudflare.com	1785	54	la.gting.com	1341
15	ns-1891.awsdns-44.co.uk	3904	35	a1-14.akam.net	1666	55	wb.qq.com	1341
16	ns-1192.awsdns-21.org	3884	36	hostname.bind	1598	56	matchweb.sports.qq.com	1336
17	ns-399.awsdns-49.com	3878	37	reachit.lenovo.com	1561	57	aq.qq.com	1335
18	providers.cloudsoftphone.com	3701	38	sso.cloudsoftphone.com	1517	58	vweixinthumb.tc.qq.com	1334
19	_ta-4a5c-4f66	3634	39	fire-base.com	1512	59	vpic.video.qq.com	1331
20	cx-messenger.imbeepro.es	3292	40	ns3.msft.net	1508	60	http.qq.com	1331

Annex 2 Table 1 - Top Names Queried from source IPs only with one request in 2020 DITL

Annex 2 Table Table 2 examines the ranking of the most popular names being queried by single query sources but aggregates the name to the second level domain. Again we see that the top querying names all would return an rcode of 0 and that they account for almost 70% of all single query sources at A and J roots.

Rank	Qname (SLD)	Number IPs
1	qq.com.	198905
2	.	160740
3	in-addr.arpa.	31893
4	gting.com.	25111
5	angsrvr.com.	20523
6	google.com.	16991
7	mielse.com.	14582
8	co.uk.	13496
9	afy11.net.	13422
10	grindr.com.	11665

Annex 2 Table 2 - Top second level names queried from source IPs only with one request

With such a large percentage of these single query source IPs issuing requests for resolvable names as opposed to the overall norm of non-existent names being sent to the RSS, additional analysis investigating the top NXDomain names was examined in Annex 2 Table 3. Here we observe a significant amount of those names are for RFC 8145 trust anchors. Many of the other names contain DNS service discovery labels (e.g. _dns-sd._udp) and some type of encoded\escaped non-existent label.

Some previous research¹⁰ looking at RFC 8145 queries to the root identified certain pieces of software that linked to DNS software libraries would cause inadvertent queries into the public DNS. Based on the names seen in Annex 2 Tables these names might suggest a similar root cause in which a piece of software (e.g. a QQ mobile app) inadvertently issues DNS queries.

¹⁰ [Roll, Roll, Roll your Root: A Comprehensive Analysis of the First Ever DNSSEC Root KSK Rollover](#)

Rank	Qname	Number IPs
1	150.109.167.160m\183\128\156\023\001\233\172\186\158\244vah\026\001	7368
2	_ta-4f66	5341
3	_ta-4a5c-4f66	3634
4	local	3192
5	49.51.82.122\023r\170j\002\180\128\156\023\001\233\172\186\158,vwukq\001	2256
6	home	2051
7	49.51.82.122\023r\170j\002\180\128\156\023\001\233\172\186\158\$[=rq\016	2002
8	lb_dns-sd_udp.\168\147z\001\248\025\022\001\192\168\001m	1767
9	dr_dns-sd_udp.0\162z\001\248\025\022\001	1713
10	r_dns-sd_udp.\168\147z\001\248\025\022\001\192\168\001m	1566
11	.-	1481
12	perforce	1295
13	db_dns-sd_udp.\168\147z\001\248\025\022\001\192\168\001m	1164
14	b_dns-sd_udp.\168\147z\001\248\025\022\001\192\168\001m	1102
15	150.109.181.132m\178\162\128\156\023\001\233\172\186\158\244vah	1089
16	b_dns-sd_udp.0\162z\001\248\025\022\001	1077
17	r_dns-sd_udp.\146z\001\248\025\022\001\192\168\001m	1031
18	dr_dns-sd_udp.\168\147z\001\248\025\022\001\192\168\001m	1000
19	r_dns-sd_udp.\184k^\001\248\025\022\001	965
20	b_dns-sd_udp.\146z\001\248\025\022\001\192\168\001m	939
21	db_dns-sd_udp.0\162z\001\248\025\022\001	924
22	ns.zyxel-usg	919
23	db_dns-sd_udp.@j^\001\248\025\022\001\192\168\001m	908

Annex 2 Table 3 - Top NXDomains queried from source IPs only with one request

Rank	Qname	Number IPs	Rank	Qname	Number IPs
1	.	27166942	21	m.root-servers.net	457922
2	net	4022360	22	ns7.cloudflare.com	456293
3	com	3672953	23	ns6.cloudflare.com	454773
4	gsm1.g4c5j.com	2678614	24	pdns196.ultradns.biz	448244
5	ns3.msft.net	1655606	25	wpad.zyxel-usg	447540
6	ns1.msft.net	1630519	26	biz	433063
7	local	1348070	27	ns2.afnic.net	424529
8	org	970844	28	www.google.com	414777
9	au	779287	29	ari.gamma.aridns.net.au	413913
10	www.microsoft.com	768918	30	dns10.ovh.net	407591
11	manus.authdns.ripe.net	765465	31	ns10.ovh.net	407491
12	uk	707717	32	ari.delta.aridns.net.au	405893
13	_ta-4f66	666081	33	u1.amazonaws.com	388165
14	ns-2027.awsdns-61.co.uk	590959	34	ari.alpha.aridns.net.au	385307
15	ns-1384.awsdns-45.org	587544	35	pixels.change.me	375795
16	\000	582250	36	u2.amazonaws.com	363989
17	ns-749.awsdns-29.net	568703	37	ari.beta.aridns.net.au	358587
18	localdomain	541666	38	https://app-measurement.com/sdk-exp	357873
19	info	516986	39	e.root-servers.net	343246
20	arpa	510894	40	de	341469

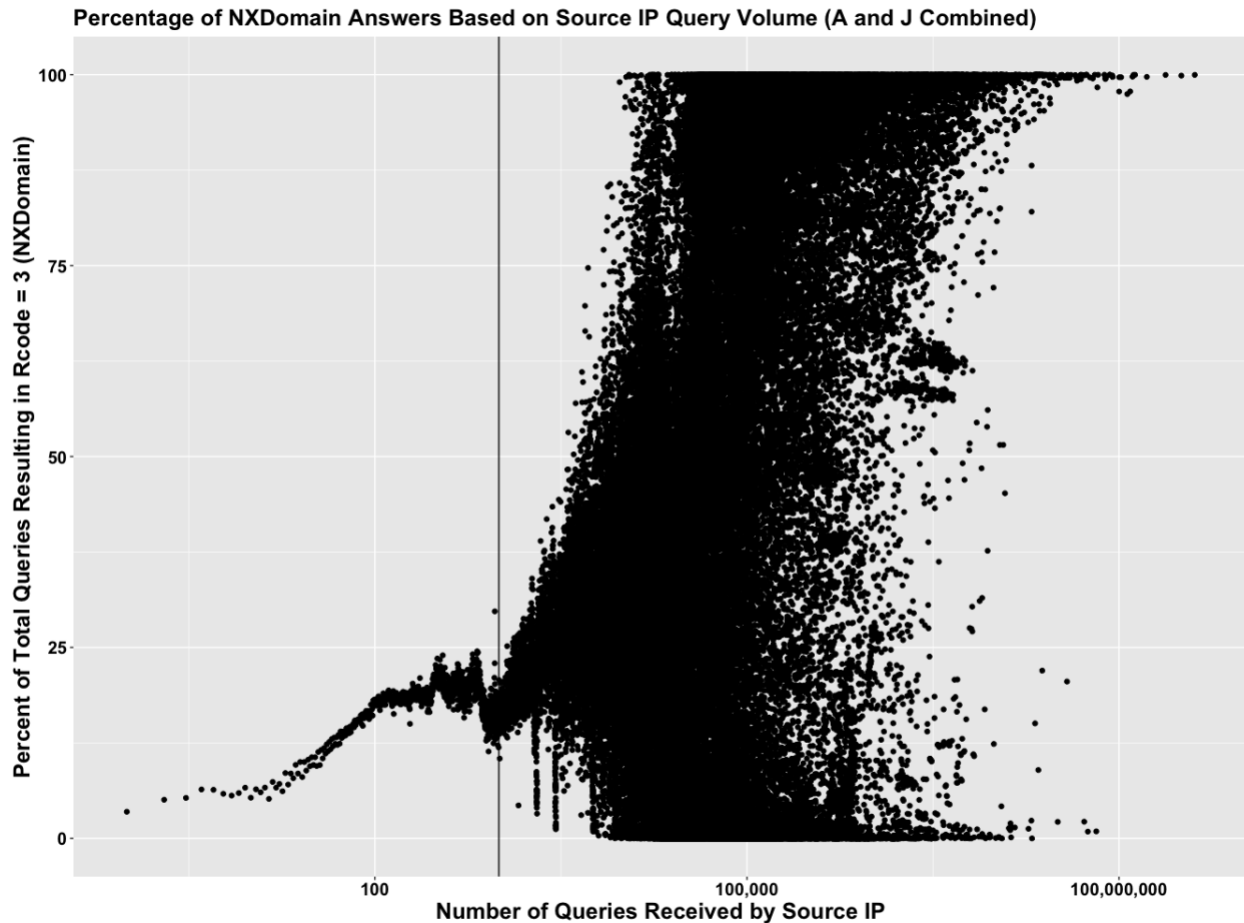
Annex 2 Table 4 - Top names queried from source IPs only with 2 to 1000 requests

A further examination of the long tail of names being queried by source IPs that sent between 2 and 1,000 queries during 2020 DITL is shown in Annex 2 Table 4. These names reaffirm the previously observed

behavior for the single query source IP address, that the majority of names being requested are not for non-existent names and the ones that are associated with RFC 8145 trust anchors.

The next examination of these low querying source IP addresses focuses on the abnormal proportion of NOERROR (rcode:0) to NXDOMAIN (rcode:3). RSSAC002 data shows that A and J NXDOMAIN rates to be around 55% of total queries. This was clearly not the case for these single query source IP addresses. In order to better understand the ratio of NOERROR to NXDOMAIN, a measurement was created to calculate the return code percentages based on the number of queries an IP sent - e.g., For all source IP addresses that sent one query what percentage were NXDOMAIN. For all source IP addresses that sent two queries, three queries, etc.

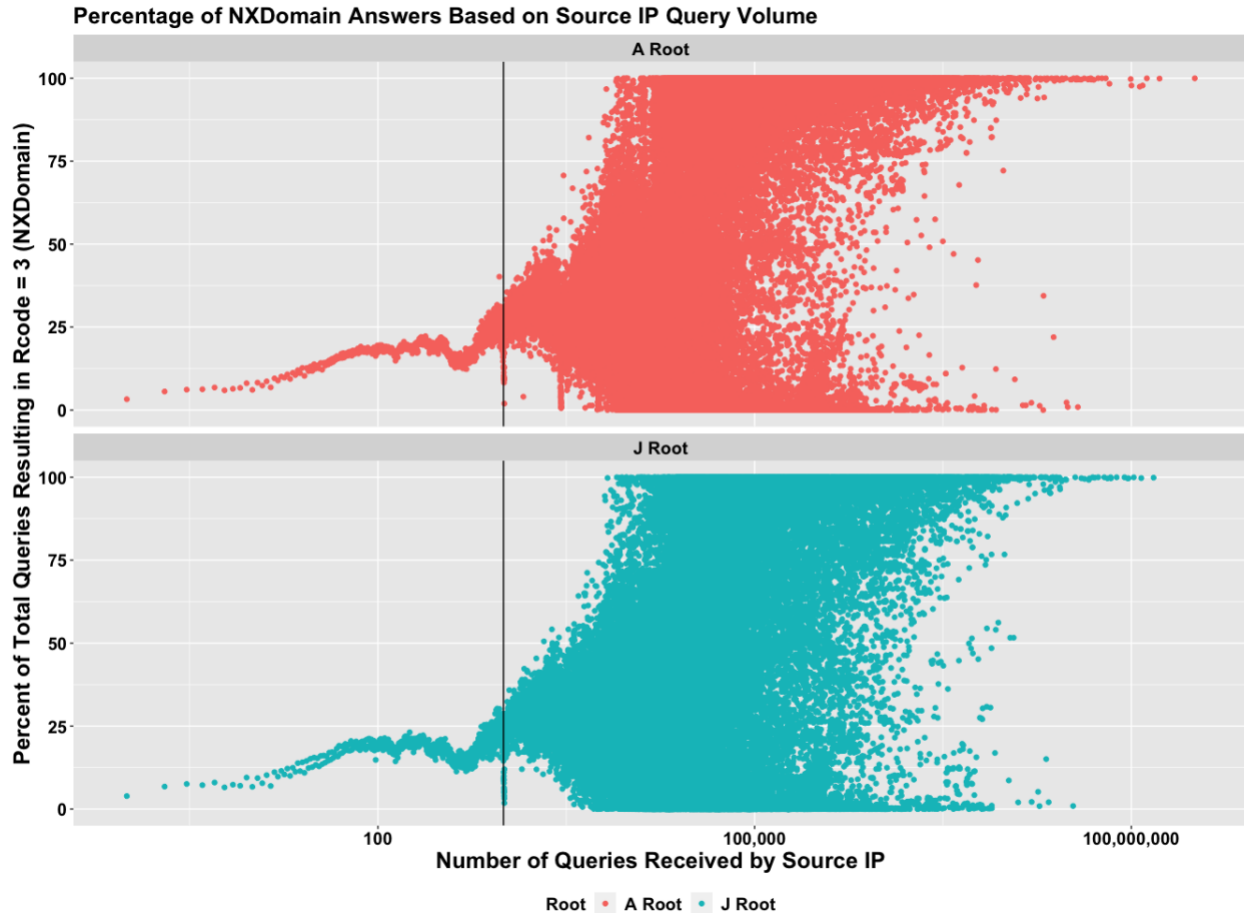
Annex 2 Figure 1 illustrates the NXDOMAIN percentage based on source IP total query volume during 2020 DITL. A very obvious dichotomy of NXDOMAIN percentage rates are observed in source IP addresses that issue smaller amounts of queries (i.e. from 1 to ~1,000 which is marked with a vertical line) than those source IP addresses issuing much larger DNS query volumes. Based on NXDOMAIN rates, these lower querying source IPs behave significantly differently than higher query volume sources. This further motivates and supports the exclusion of these source IP addresses for RSS similarity measurements.



Annex 2 Figure 1 - Percentage of NXDOMAIN answers per source IP total query volume.

In Annex 2 Figure 2, A and J roots are separated to ensure this trend is observed at both of the RSIs. Due to processing and time limitations, this measurement was not feasible for the remaining 2020 DITL roots.

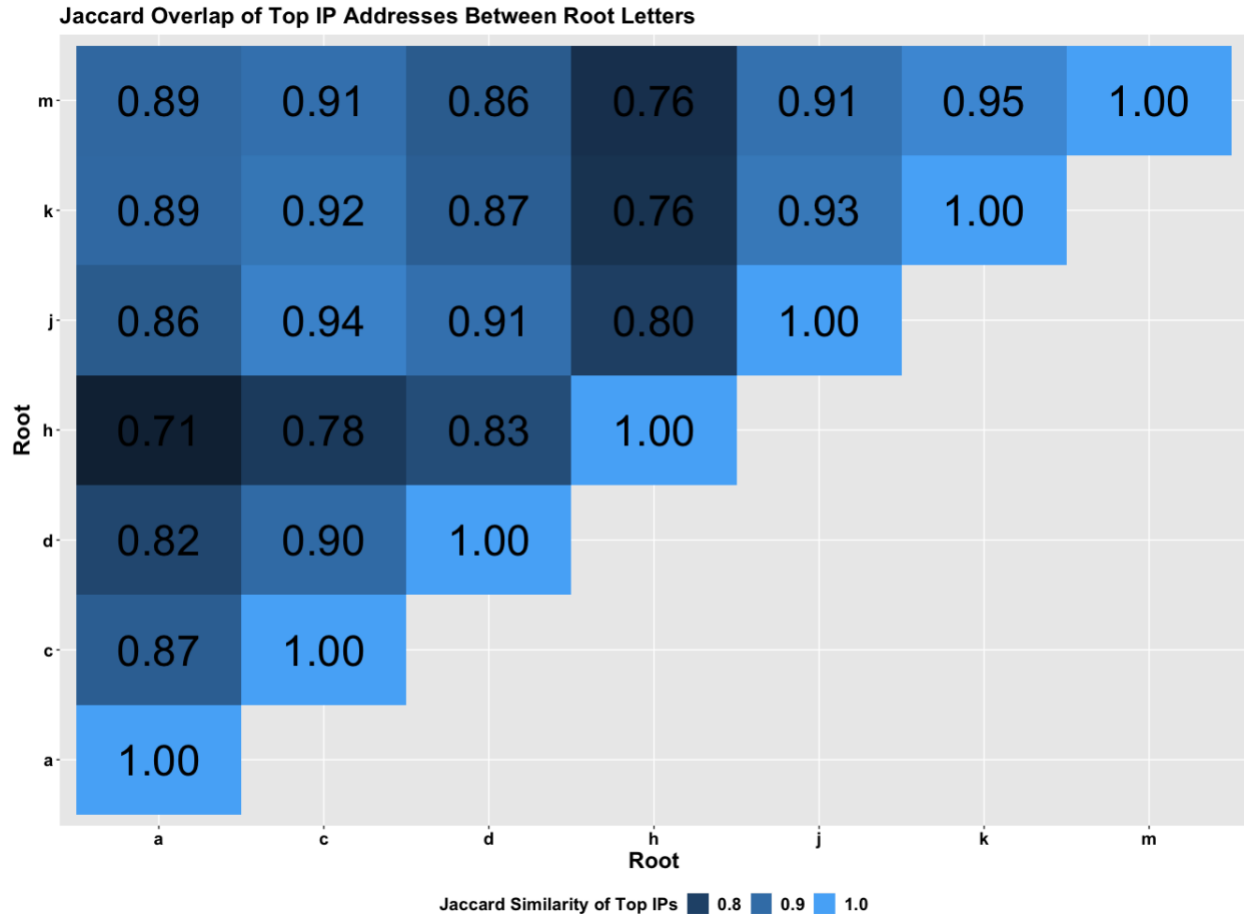
The vertical line in Annex 2 Figures 2 and 3 marks a threshold of source IP addresses that sent 1,000 queries during 2020 DITL. Within the context of name collision analysis, capturing queries that result in NXDOMAIN is crucial for risk assessment purposes. Using the 1,000 query threshold to separate that set of IPs, which clearly have a different behavior pattern, results in 98% of total NXDOMAIN responses to be captured.



Annex 2 Figure 2 - Pct. of NXDOMAIN answers per source IP total query volume by RSI.

The similarity measurements between RSIs in this study were done by using the top 115K (0.67%) of IP addresses. Based on the new data and insights presented in Annex 2, those measurements were re-evaluated using a threshold of 1,000 queries per source IP address. This new threshold changes the top-talkers to include 1.32 million IPs constituting 7.8% of total IPs. These new top talkers constitute nearly 98% of the total 2020 DITL queries.

85% of these 1.32 million IPs were on average seen at the seven RSIs. Annex 2 Figure 3 shows the Jaccard similarity matrix comparing each RSI with another using these new top-talkers. The overall pairwise similarity measured 0.86 (compared to the previous .96).



Annex 2 Figure 3 - Jaccard Overlap of Top IP Addresses Between RSIs