

Data Sensitivity Analysis

- Executive Summary** **2**
 - Terminology 3
 - Data Sensitivity in the Context of DNS Name Collisions 3
 - Studies 3

- Study 1: Root Server Identity Comparison** **5**
 - Data** **5**
 - Notable Limitations of the Data 5
 - Measurements 5
 - Query Volume per RSI 6
 - Top Talkers 7
 - Geographic Relevance 12
 - Top NXDomain TLDs per RSI for Top Talkers 16
 - Study 1 Key Observations:** **19**

- Study 2: Public Recursive Resolver and Root Comparison** **20**
 - Data 20
 - Notable Limitations of the Data 20
 - Measurements 22
 - Total Query Volume per TLD Distribution 22
 - A and J Root Servers Compared to a PRR Using Total Query Volume per TLD Ranking as a Function 23
 - A, J, and L Root Servers Compared To Public Recursive Using Distinct Source IPs per TLD Ranking Function 26
 - Study 2 Key Observations: 28

- Key Findings** **28**

- Annex 1: Statistical Methods** **30**
 - Jaccard Index 30
 - Gini Coefficient 30

Executive Summary

As part of the Name Collision Analysis Project Study Two goals and objectives, a study was commissioned by the NCAP Discussion Group to better understand how representative DNS data from various points of the DNS hierarchy is within the context of name collisions. 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, known as the Data Sensitivity Analysis, focuses on two key measurements: (1) comparing traffic received at each root server identity 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 identity 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 Identity (RSI) - thirteen identities, each of which is named with the letters 'a' to 'm', collectively administered by twelve root server operators. They are authoritative for the 'root-server.net' domain.
- Day-In-The-Life (DITL)¹ - a large-scale data collection project initially undertaken every year since 2006. This data has historically been the primary measurement asset for name collision studies.

Data Sensitivity in the Context of DNS Name Collisions

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 risks². 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 RSI's DNS traffic 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 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.

Studies

The Data Sensitivity Analysis project consists of two main studies: the comparison of traffic among RSIs and the comparison of name collision traffic between root and recursive

¹ <https://www.dns-oarc.net/oarc/data/ditl>

² <https://www.icann.org/en/announcements/details/mitigating-the-risk-of-dns-namespace-collisions-final-report-by-jas-global-advisors-30-11-2015-en>

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 identities 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 a relatively opaque and widely inaccessible data for name collision analysis - traffic to public open recursive resolvers. The top leaking query strings based on query volume and source diversity are relatively comparable to the top strings observed by root server identities.

Study 1: Root Server Identity 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 identities 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 RSIs, 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 measurements 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 talkers³ that constitute a large percentage of overall traffic
5. Measuring overlap of top talkers at each RSI
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 NXDomain 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.

³ Top talkers are recursive resolvers that issue the largest amount of DNS queries to the RSS.

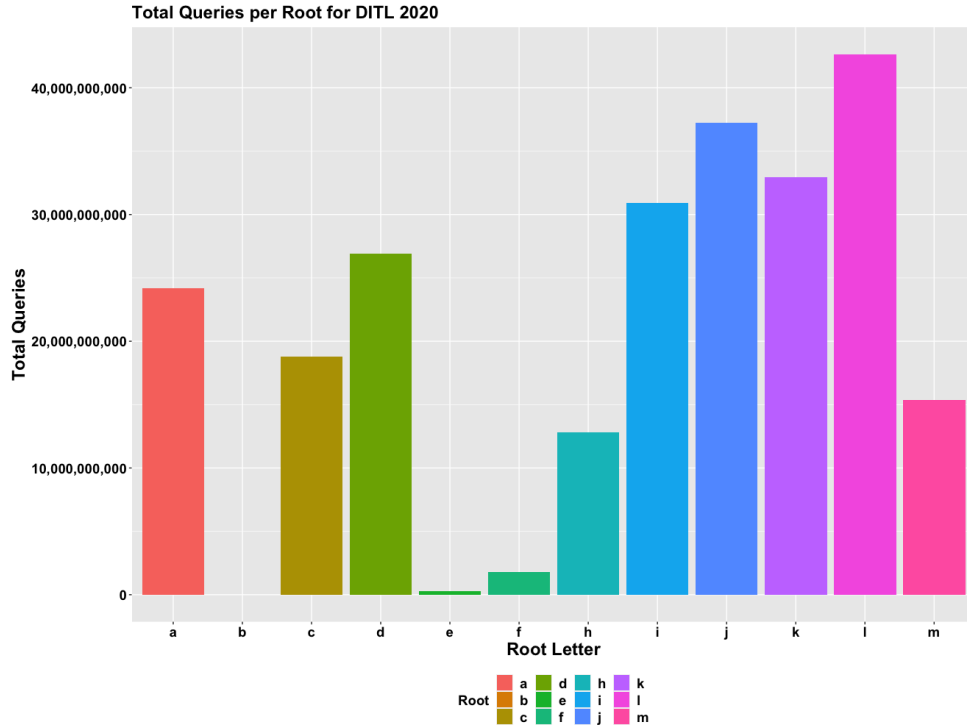


Figure 1 - Query Volume per RSI

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 non-excluded RSIs, is relatively even.

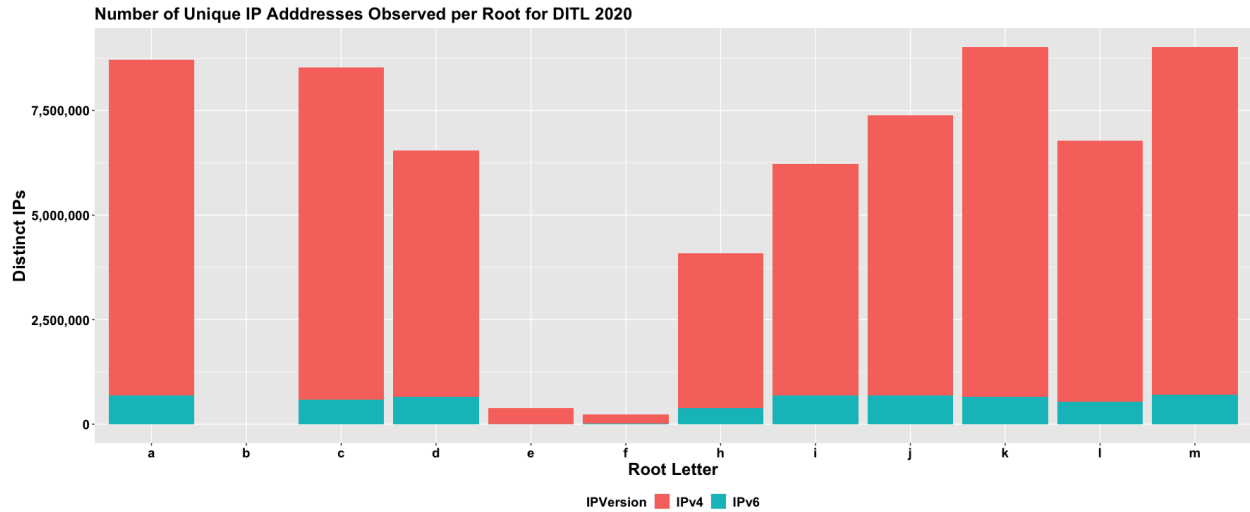


Figure 2 - Unique IP Addresses Observed per RSI

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 Distribution⁴ 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.
- 98% of IP addresses issued 10,000 or fewer queries.

⁴ https://en.wikipedia.org/wiki/Power_law

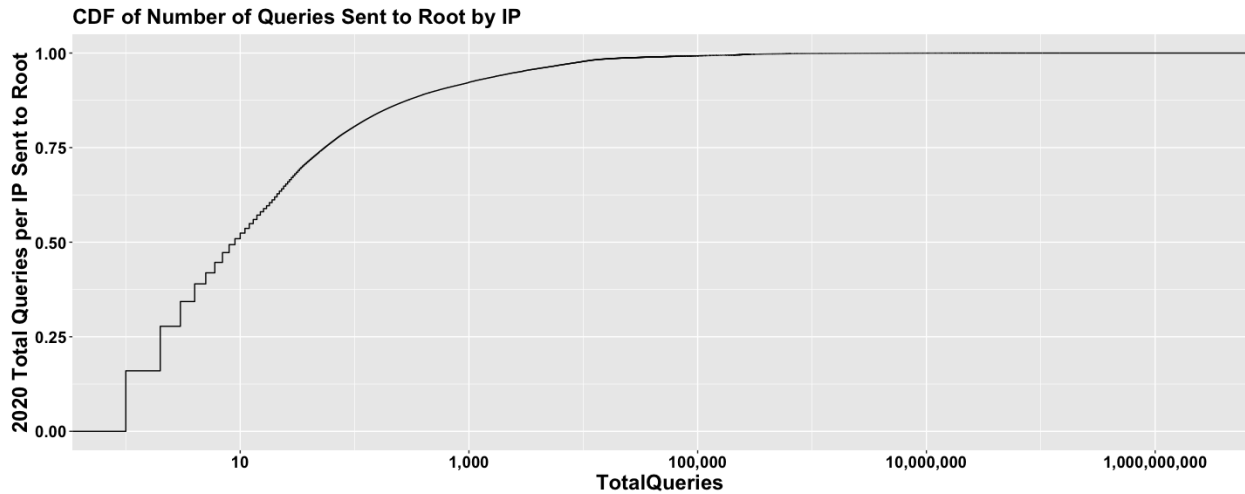


Figure 3 - Cumulative Distribution of the Number of Queries Sent to RSS by IP Address

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 systems 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 leakage) on the RSS.

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 115K IPs. Likewise, 95% of the total 2020 DITL can be represented by 250K IPs. The remaining 5% of query volume is distributed in the long tail of millions of IPs.

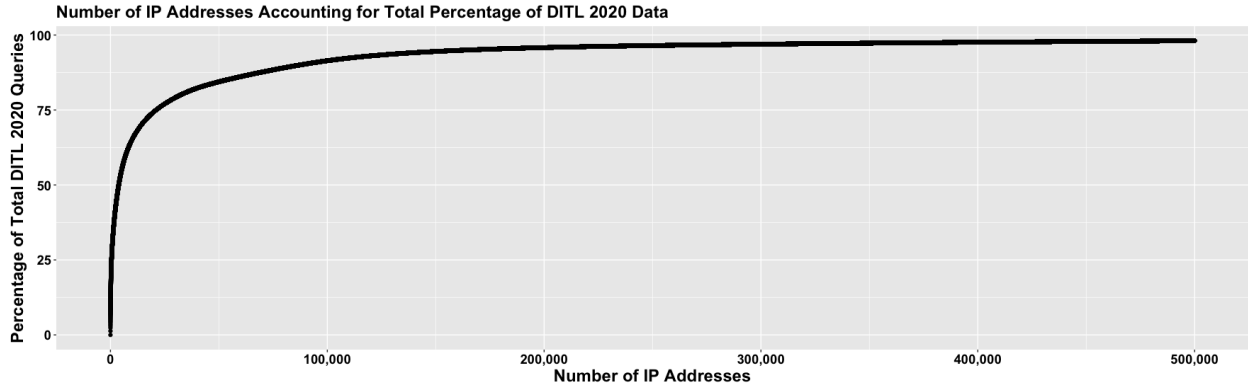


Figure 4 - Number of IP Addresses Accounting for Total Percentage of DITL 2020

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 are seen by all RSIs - a key indicator that any RSI may be representative of the general RSS.

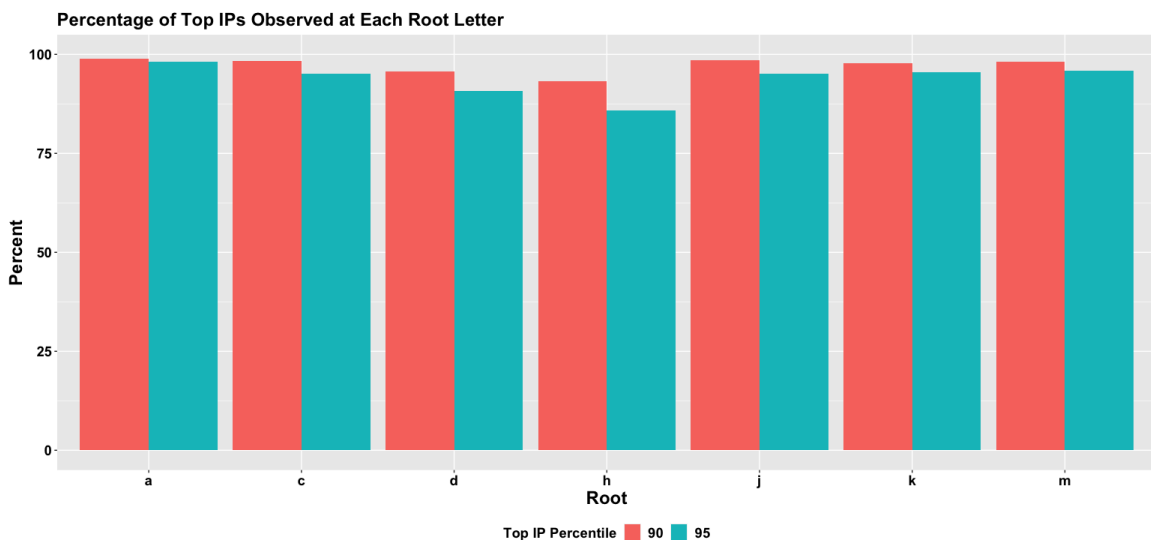


Figure 5 - Percentage of Top IPs Observed at each RSI

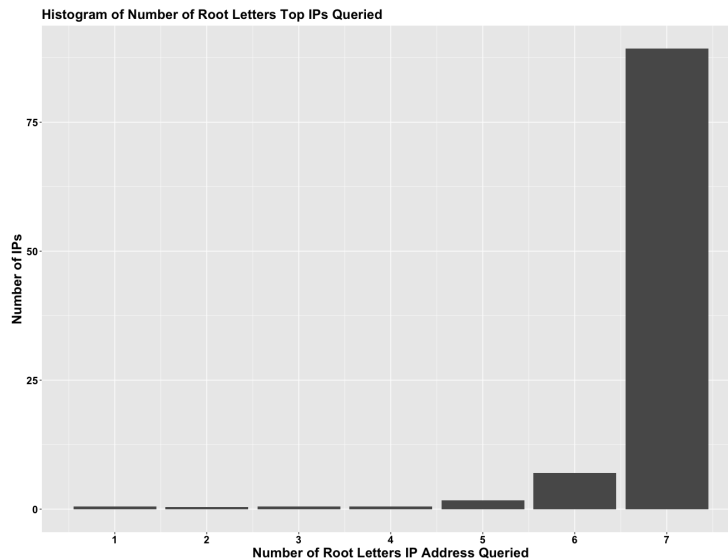


Figure 6 - Histogram of Number of Root Letters Top IPs Queried

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.

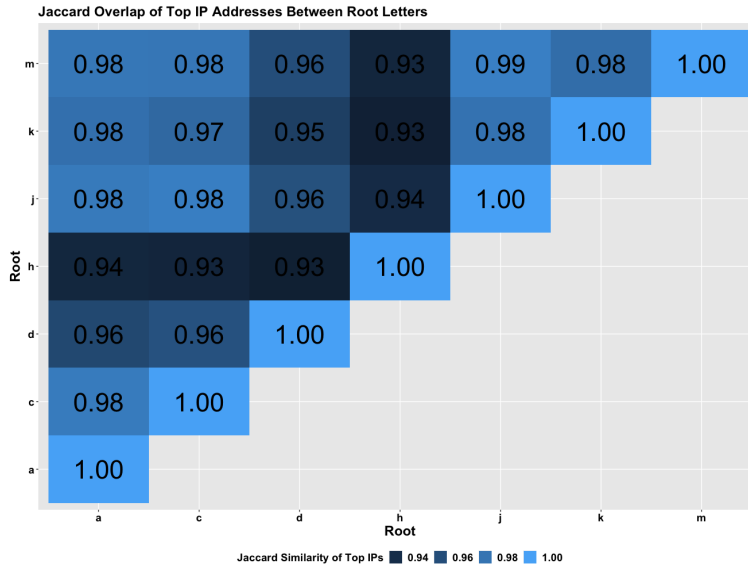


Figure 7 - Jaccard Overlap of Top IP Addresses Between RSIs

Geographic Relevance

The proceeding 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 program⁵ 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.

⁵ <https://github.com/measurement-factory/ipv4-heatmap>

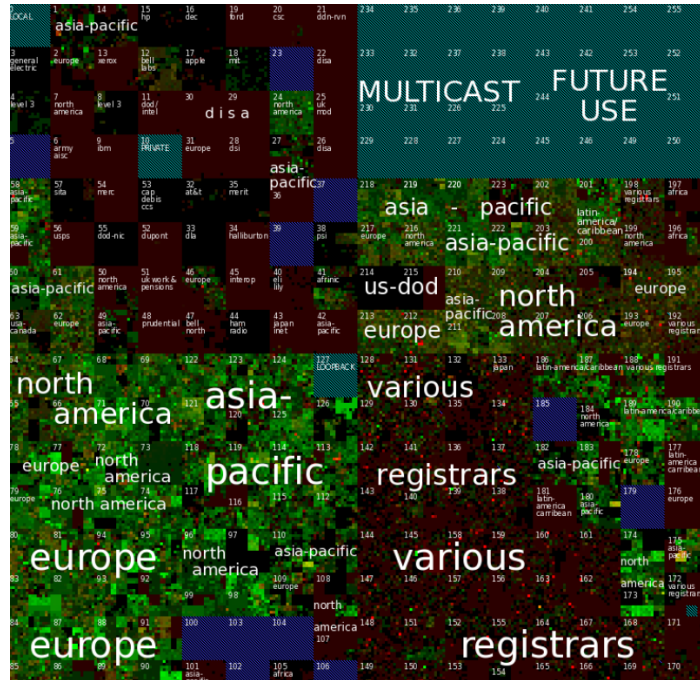


Figure 8 - Example of Hilbert Curve IPv4 Visualization

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.

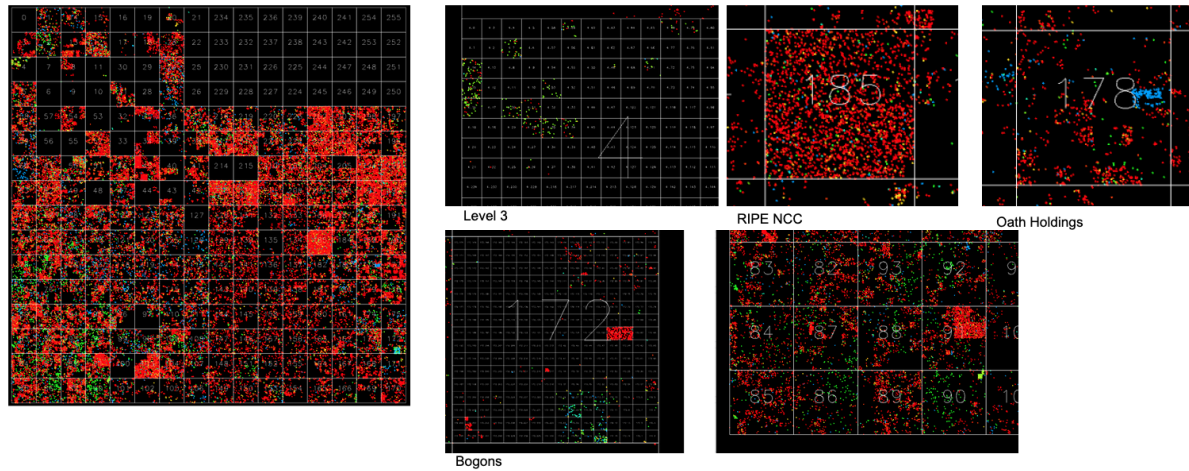


Figure 9 - IPv4 Spatial Distribution of Top Talker IP Addresses

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.

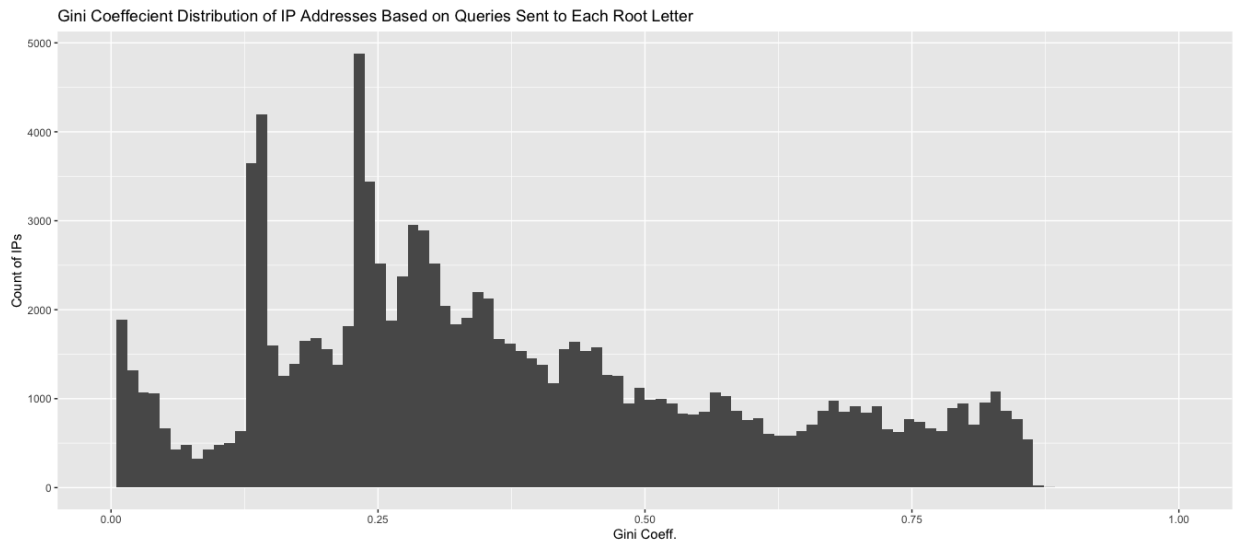


Figure 10 - Gini Coefficient Histogram of Top Talker IPs

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.

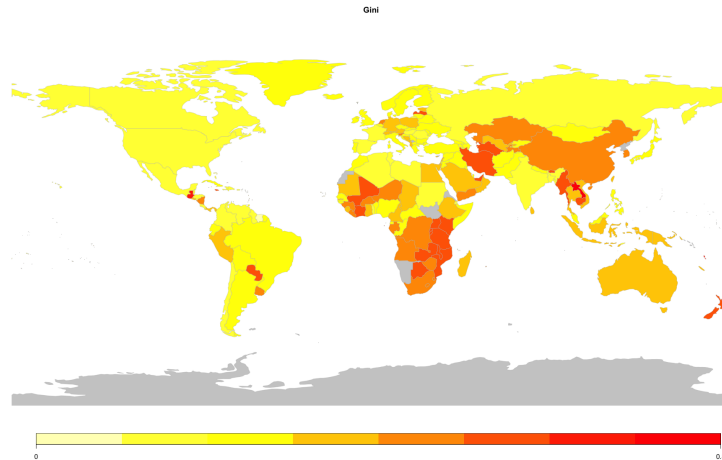


Figure 11 - Country to RSI Traffic Distribution

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

Top NXDomain TLDs per RSI for Top Talkers

The previous analysis shows with a high level of confidence that traffic to any RSI is 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 (NXDomain) TLDs on a per RSI basis.

In order to understand how top leaking NXDomain 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 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.

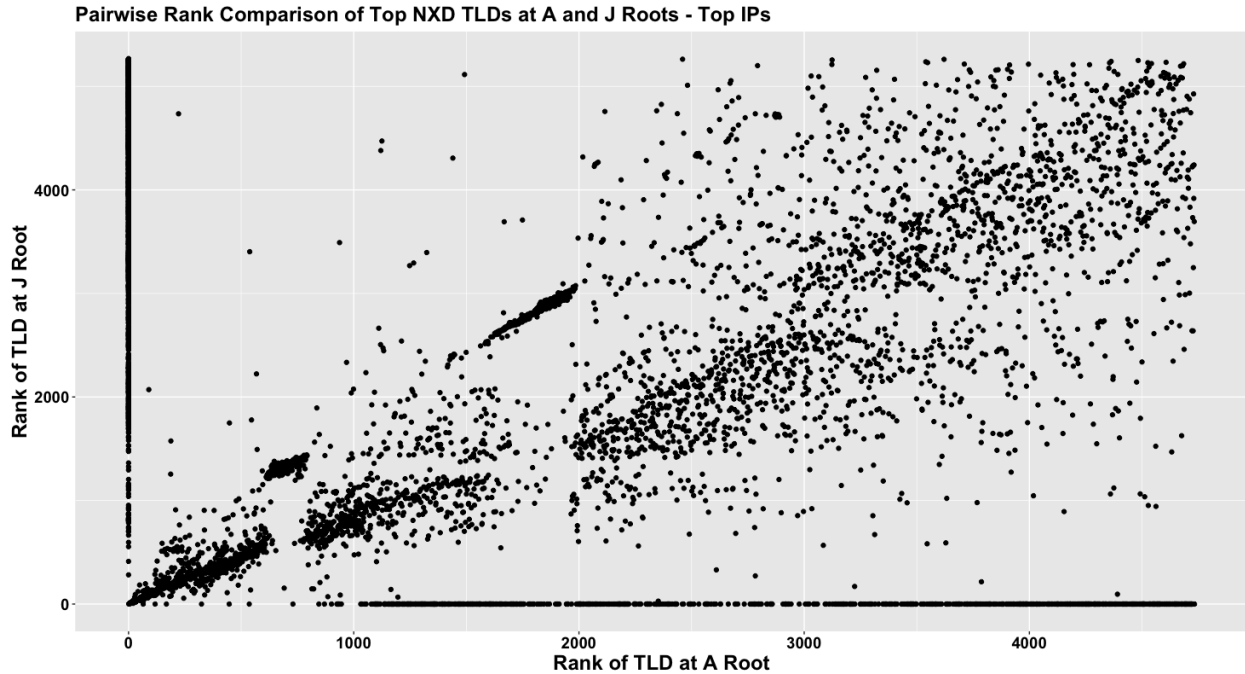


Figure 13 - Rank Comparison of A and J Top NXD TLDs

Figure 14 below 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 string were also required to match the regular expression $[a-z0-9]{3,63}$. Again a good rank correlation is expressed at low ranks and diminishes as the ranks approach 1,000.

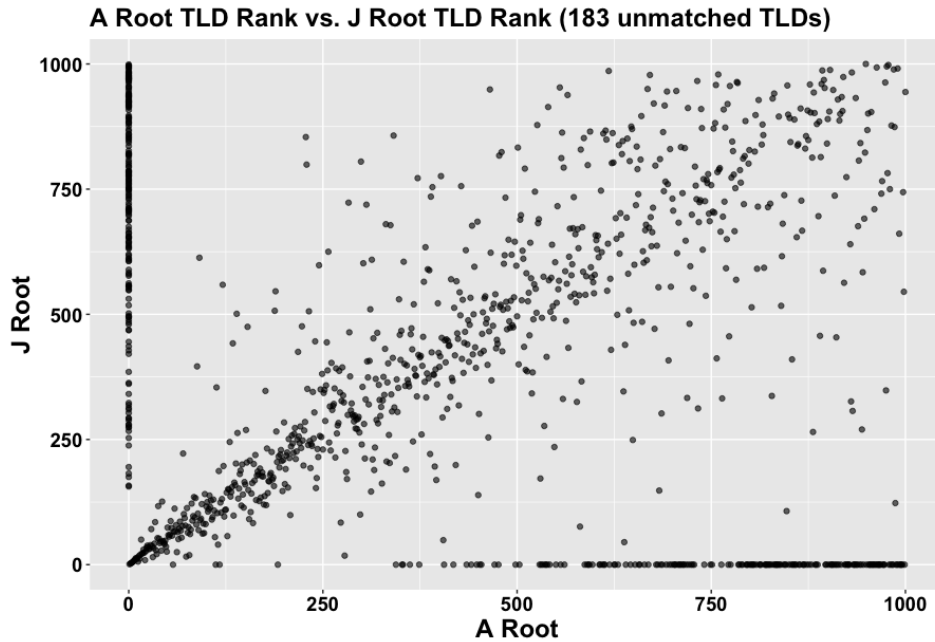


Figure 14 - Rank Comparison of A and J Top NXD TLDs

Figure 15 is a standard Venn diagram showing the overlap of the Top NXDomain TLDs plotted in Figure 14. The overlap of A and J RSI Top NXDomain 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.

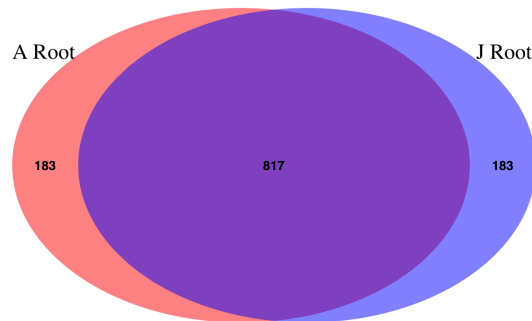


Figure 15 - Venn Diagram of TLD Overlap between A and J RSIs

Study 1 Key Observations:

The main result, as expected, is that all the roots see essentially the same set of queries.

- A large percentage of RSS queries are sent from a relatively small set of IP addresses (115K).
- Top-talking IPs are broadly seen at all root letters.
- Traffic is generally fairly evenly distributed across all root letters.
- Some geographic affinity/preference to certain root letters does occur.
- Top leaking strings between letters appear to generally correlate for the top 1K strings.
- Lower volume leaking 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 transpired and 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 from several root server identities and one public recursive resolver were collected and measured in such a way that would facilitate the examination of top NXDomain TLD strings. 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 strings 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 NXDomain 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 NXDomain 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.

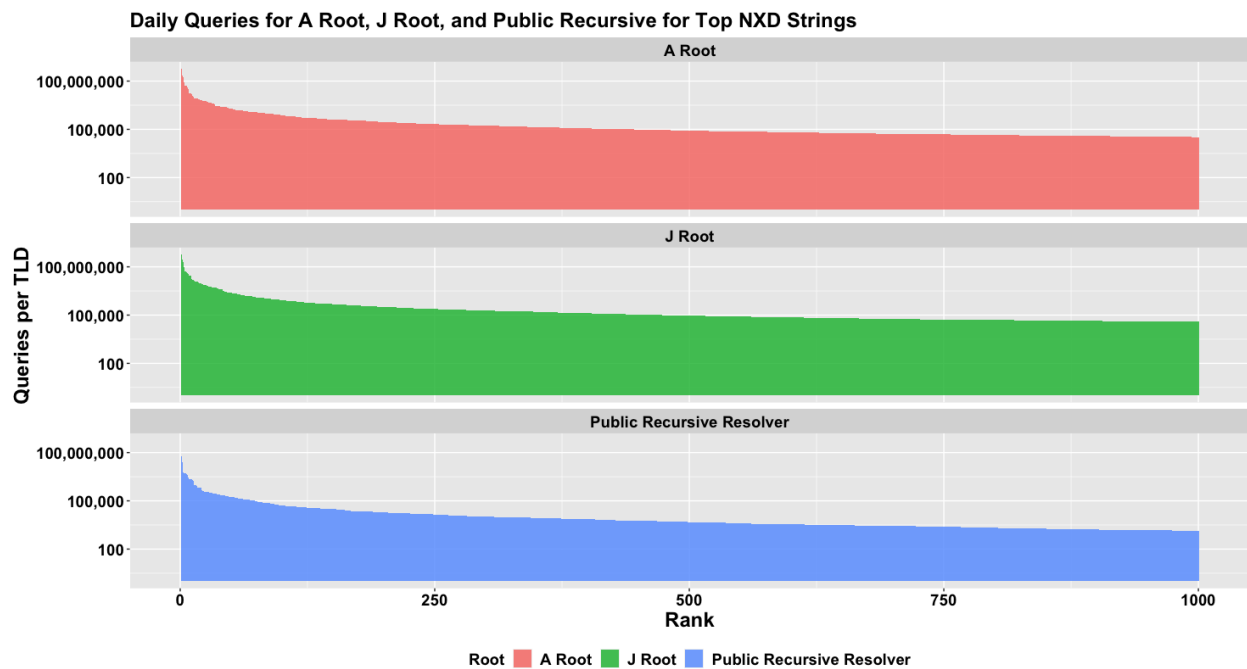


Figure 16 - Daily Queries for A and J RSIs and the PRR for Top NXD TLDs

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).

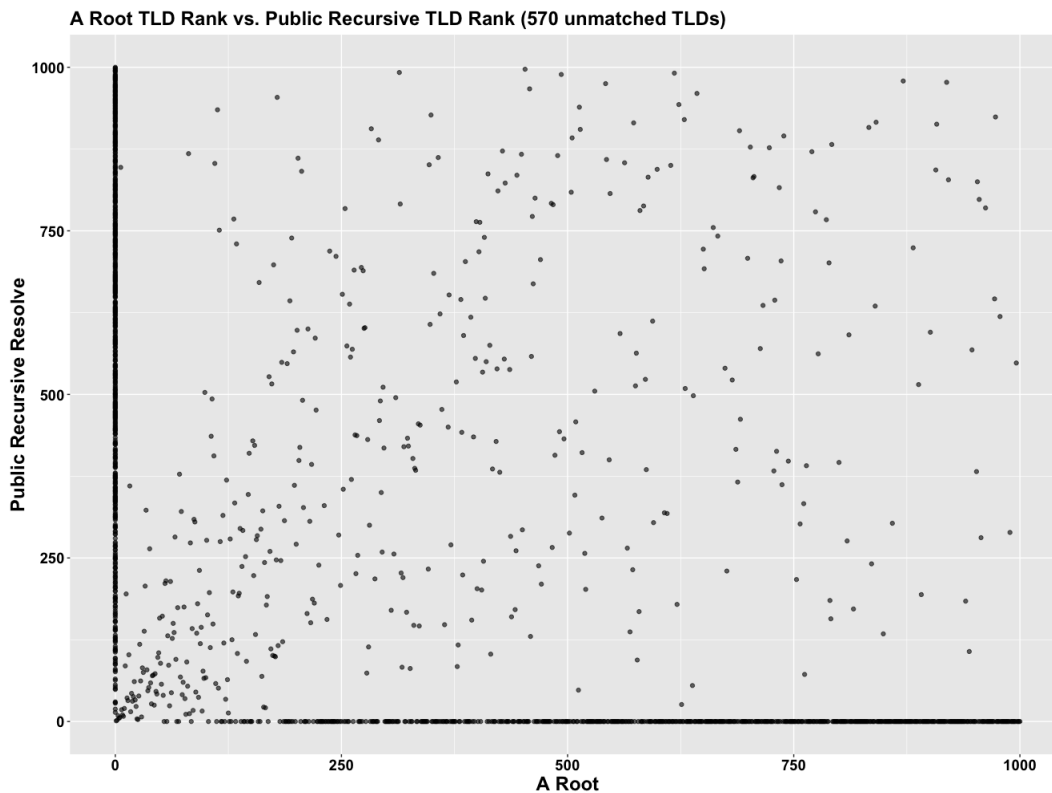


Figure 17 - Rank Correlation of Top TLDs at A Root and Public Resolver based on Query Vol.

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 NXDomain 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.

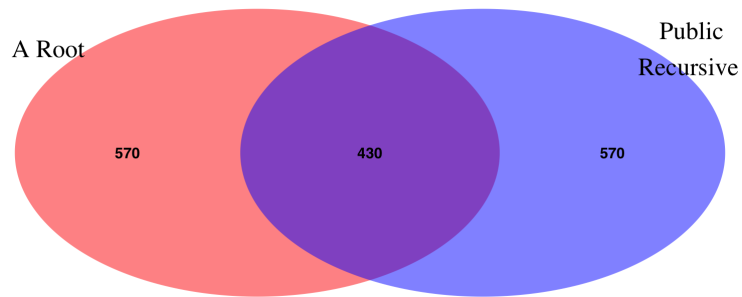


Figure 18 - Venn Diagram showing TLD overlap of A Root and PRR based on Query Vol.

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 NXDomain TLDs were observed by both the RSI and the PRR.

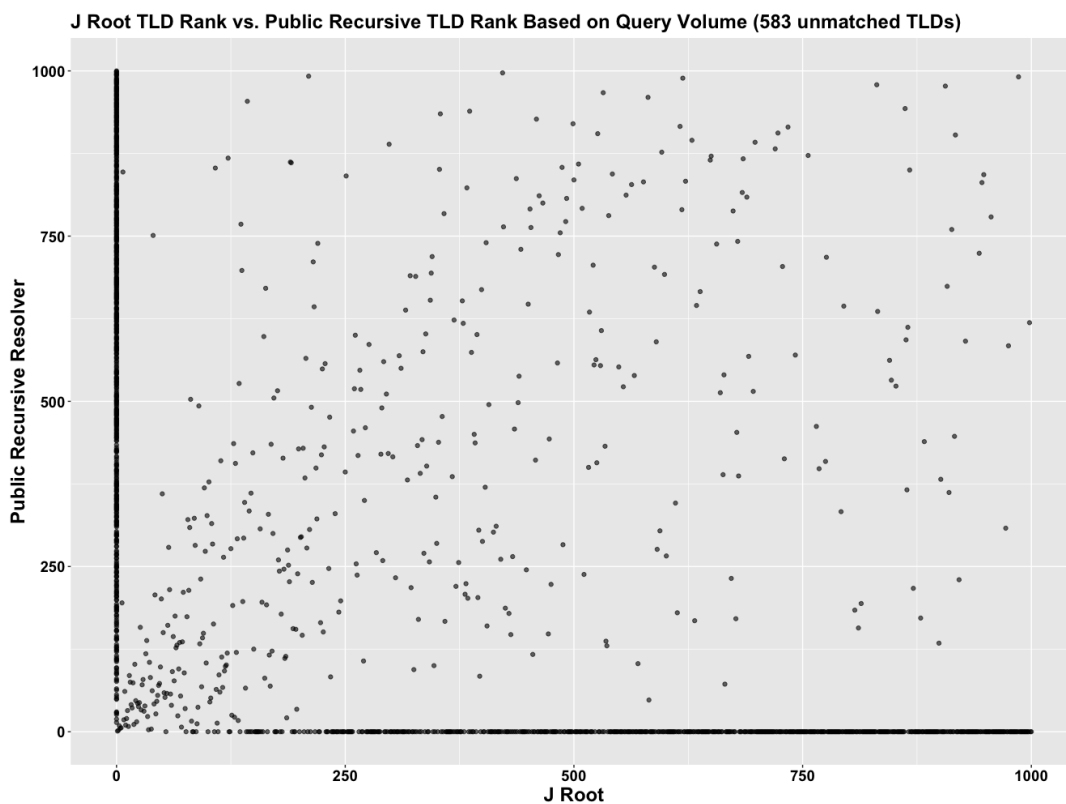


Figure 19 - Rank Correlation of Top TLDs at A Root and the PRR based on Query Volume

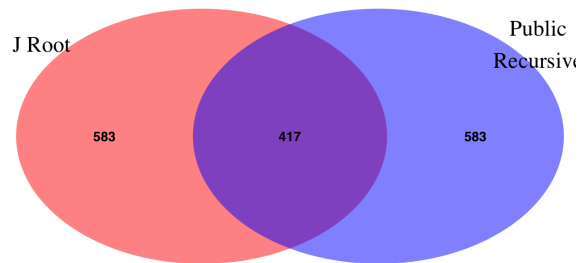


Figure 20 - Venn Diagram showing TLD overlap of the J RSI and PRR based on Query Vol.

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.

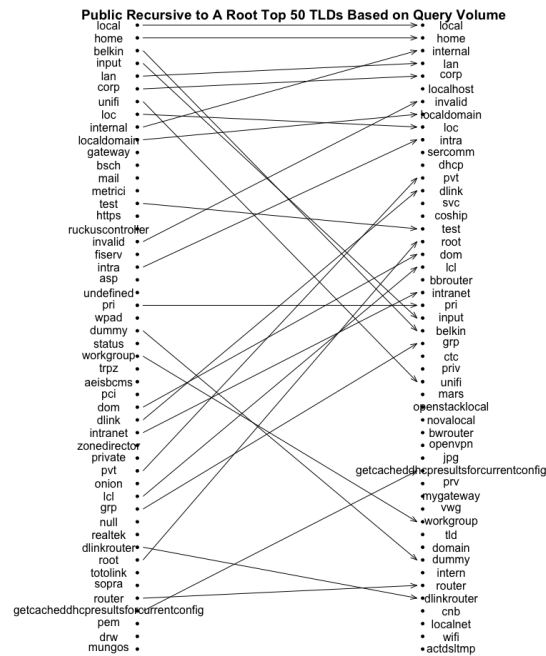


Figure 21 - PRR to A Root Top 50 TLDs Based on Query Volume

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

Using the secondary critical diagnostic measurement of IP source diversity, measurements were made between three RSIs and the PRR's top 1,000 NXDomain TLDs ranked by the number of unique IP addresses observed per TLD. An initial measurement looking at 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.

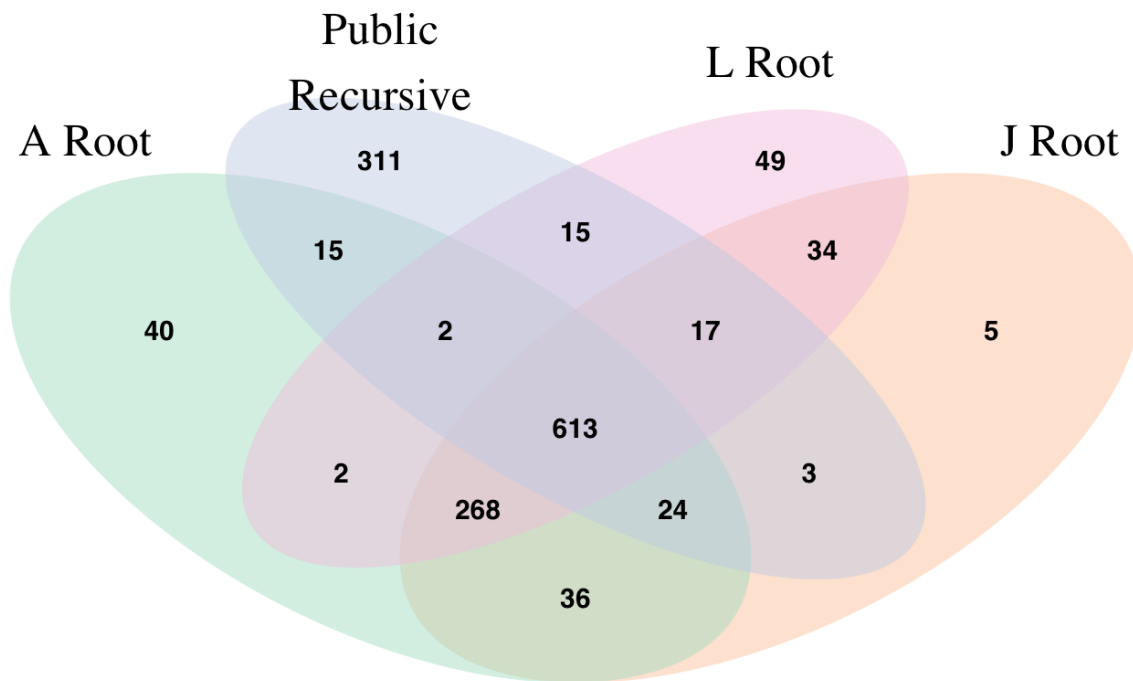


Figure 22 - Venn Diagram Comparing Overlap of Top TLDs at A, J, and L RSIs and PRR based on Source IP Address Diversity

Examining a rank scatter plot between an RSI and the PRR does indicate a slightly better correlation of TLD rankings; however, this correlation appears very weak, at best, and mainly for the top-ranking strings that had large source diversity measurements (i.e., TLD rankings under 100).

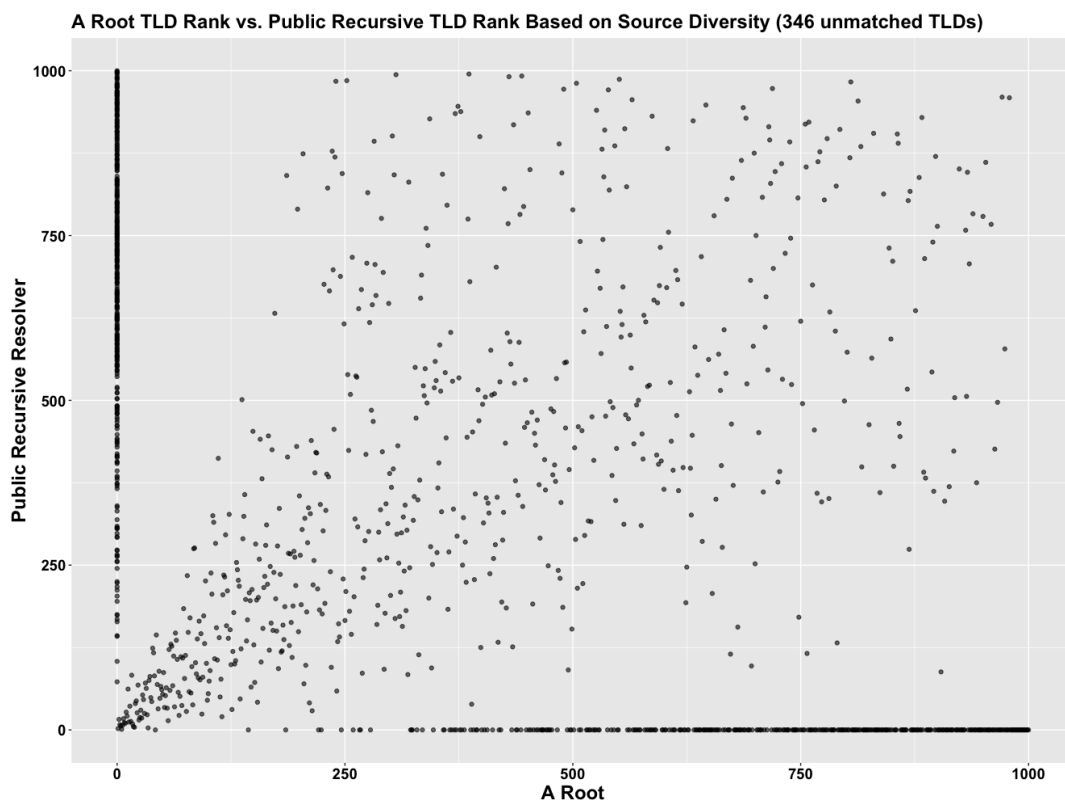


Figure 24 - Rank Correlation of Top TLDs at A Root and PRR based on IP Diversity

Study 2 Key Observations:

The main result is that the PRR observes name collision strings in a different manner than the root server system. This implies that name collision risk assessment and analysis based solely on root server traffic will be significantly incomplete and must incorporate a high degree of uncertainty. Until additional recursive resolver data can be obtained, it is unclear how recursive resolvers compare to each other and if sampling from any recursive resolver would be representative, analogous to the finding from Data Sensitivity Analysis Study 1 in which data from any RSI is likely to be representative of the whole root server system. However, given the unique aspects of public recursive resolvers implementing certain DNS protocol optimizations, their customer base, and other attributes, it is likely that each public recursive resolver has a special and unique vantage point of name collision DNS traffic. Furthermore, given the increasing difficulty of obtaining any public recursive resolver data for name collision analysis, in part due to legal and privacy concerns, it is likely that public resolver data will be even more obscure or scarce in the future.

Key Findings

The two studies in this Data Sensitivity 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: DNS traffic observed at any RSI is largely representative of traffic across the whole of the root server system at any given moment in time.

Implications:

- Future name collision risk assessments need not solely rely on yearly DITL data collection efforts.
- 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.
- A complete and accurate risk assessment of a string's name collision potential cannot be determined prior to the string's 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

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

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).