# Routing-Aware Partitioning of the Internet Address Space for Server Ranking in CDNs

(A more recent version of this manuscript is published in Elsevier Computer Communications, vol. 106, July 2017)

Gonca Gürsun

DEPARTMENT OF COMPUTER SCIENCE, OZYEGIN UNIVERSITY

## ABSTRACT

The goal of Content Delivery Networks (CDNs) is to serve content to end-users with high performance. In order to do that, a CDN measures the latency on the paths from its servers to users and then selects a best available server for each user. For large CDNs, monitoring paths from thousands of servers to millions of users is a challenging task due to its size. In this paper, we address this problem and propose a framework to scale the task of path monitoring. Simply stated, the goal of our framework is clustering IP addresses (clients) such that in each cluster the choice of best available server is same (or similar). Then, finding a best available server for one client in a given cluster will be sufficient to assign that server to the rest of the clients in the cluster.

To achieve this goal, first we introduce two distance metrics to compute how similar the server choices of any given two clients. Second, we use a clustering method that is based on interdomain routing information. We evaluate the goodness of our clusters by using the metrics we introduce. We show that there is a strong correlation between the similarity in how two destination clients are routed to in the Internet and the similarity in their server selections. Finally, we show how to choose representative clients from each cluster so that it is sufficient to learn the latencies from the CDN servers to the representative and find a best available server accordingly for the rest of the clients in the same cluster.

## 1. INTRODUCTION

A Content Delivery Network (CDN) is a collection of servers that deliver content to end-users on behalf of content owners. Today, significant amount of Internet traffic is served by CDNs. For example, one of the largest CDNs, Akamai, currently delivers 15-30% of all web traffic from a large distributed platform. This platform consists of over 160,000 servers in more than 100 countries and 1200 ISPs around the globe [32].

The main job of a CDN is replicating content across its geographically distributed server regions and redirect end-users to a *best available* server region at a given time. The goal is serving end-users with high performance [7], that is each user is redirected to a region to which it has low latency. One can expect that mapping an end-user to its geographically closest server region is sufficient. However, there are many cases where geographical closeness does not infer low latency [21]. Instead, the conditions in the network determine which server regions are best performing at a given time. The challenge is, since the network is dynamic, the conditions on its paths are to change frequently. Therefore, the CDN needs a monitoring solution that can keep pace with the variability of the Internet paths.

In large CDNs, such as Akamai, monitoring paths between all the server regions and end-users is challenging due to the scale of the task. At a given day Akamai sees more than 788 million unique IPv4 addresses [24]. It is not feasible to measure latencies from hundreds of thousands of servers to all these IPs. In the practice of traditional DNS-based mapping, end-users are represented by their local DNS resolvers and the path measurements are taken between the servers and the local DNS resolvers [23]. Although, this reduces the size of the task up to some extent, there is still need for clustering local DNS resolvers since there are millions of them in the Internet. Recently, with the increase in usage of public DNS resolvers [8, 9] and the adoption of EDNS [31], CDNs move towards end-user mapping [4]. In the case of end-user mapping, the users are not represented by their local DNS any more. That is, the need for partitioning the Internet address space to scale the path monitoring task is more crucial than ever. In this paper, we present a method that reduces the size of the path monitoring and server ranking tasks both in the case of DNS-based and end-user based mapping.

Our study has the following three stages.

**1. Clustering clients.** We seek to find a partitioning of the IP address space such that the clients in a given partition[1] orders the server regions from least latency to most in similar fashion. The reason why we are not only

---

[1] Throughout this paper, we use the terms *clustering* and *partitioning* interchangebly.

interested in finding the least latency region but also in ordering the regions is as follows.

In addition to high performance, *best available* server region is subject to some other constraints such as load balancing at the CDN, availability of the requested content at the server, allowance rules (enforced by ISPs and content publishers) on serving specific contents to users from specific regions etc. Therefore, for a given client, a best available server region is the one with the lowest latency that also satisfies the constraints. For that reason, the clients *rank* the server regions from least latency to most and then these rankings are used as input to the server mapping algorithm.

The clustering scheme we propose for our problem is called RS-CLUSTERING and it is introduced in [19]. RS-CLUSTERING is a method that groups BGP prefixes based on how similar the ASes in the Internet route to these prefixes. The key idea behind using this clustering scheme for our problem stems from the fact that routing is one main factor that impacts path latencies. Therefore, our hypothesis is that if traffic from the server regions to two client prefixes follow the same paths, then these clients rank the server regions similarly. We show that this intuition holds. Routing-aware clustering successfully partitions the address space and outperforms other clustering methods, such as the ones based on AS or geography.

**2. Evaluating the Goodness of Clusters.** Once the clusters are obtained, the next step is evaluating their goodness. In a good cluster, we expect that the server rankings of clients to be similar to each other. Such similarity can be defined in various ways. For instance, one can expect that in a good cluster, all clients have the exact same server region as their first (rank-1) choice. Or alternatively, one can expect that each server region is ranked in close positions by all clients in the cluster.

To capture these expectations we propose two metrics, called GEOMETRIC DISTANCE (*g-dist*) and PARTIAL SPEARMAN FOOTRULE DISTANCE (*ps-dist*) We use *g-dist* and *ps-dist* to measure the similarity between two server rankings. Using these metrics we show how to evaluate any given clustering scheme.

**3. Finding representative clients for each cluster.** Finally, we seek to find a client from each partition whose server ranking is a good representative of all the other clients in its partition. We scale the task of path monitoring by taking measurements only to the representative of each cluster. We first find a *consensus ranking* per cluster by aggregating the rankings of all clients in the cluster. Then, we show that assigning one client at random from the center of the cluster is almost as good as the consensus ranking.

**Roadmap.** The rest of the paper is organized as fol-



Figure 1: Simplified view of Akamai CDN architecture with four Akamai server regions, $\{T_1, T_2, T_3, T_4\}$, and six clients $\{X_1, X_2, X_3, X_4, X_5, X_6\}$. Dotted lines show the paths that are monitored for client $X_6$.

|       | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | $X_6$ |       | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | $X_6$ |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| $T_1$ |       | 10    | 42    | 20    | 5     | 34    | $T_1$ |       | 1     | 3     | 3     | 1     | 1     |
| $T_2$ | 90    | 15    | 36    |       | 13    | 40    | $T_2$ | 1     | 2     | 2     |       | 2     | 2     |
| $T_3$ | 240   | 140   |       | 3     | 29    | 80    | $T_3$ | 2     | 3     |       | 1     | 3     | 4     |
| $T_4$ | 260   |       | 15    | 10    |       | 65    | $T_4$ | 3     |       | 1     | 2     |       | 3     |
|       |       |       | (a)   |       |       |       |       |       |       | (b)   |       |       |       |

Figure 2: Latencies between the server regions and the clients (a) and the corresponding rank vectors (b) for the example in Figure 1.

lows. In Section 2 we describe the server mapping system of Akamai. In Section 3, we introduce the metrics we use to evaluate the goodness of clusters and follow by describing our dataset in Section 4. In Section 5, we set basis for the routing-aware clustering by investigating whether IP addresses can be pre-clustered to their BGP prefixes. In Section 6 we show how to group BGP prefixes further based on inter domain routing choices in the Internet. In Section 7 we show how to find representative nodes per cluster to scale the path monitoring task. We present related works in Section 8, discuss some issues related to our work in Section 9, and finally conclude in Section 10.

## 2. BACKGROUND

In this section, we first provide a high-level description of the server mapping system in Akamai's CDN. Next, we present the challenges in the system and the goals of our work.

### 2.1 Ranking Server Regions

The core component of Akamai's CDN is the mapping system. One main job of the mapping system is finding a *best available* server region for each client. In order to do that, first, the mapping system monitors the network conditions in the Internet to learn the latencies on the paths from the server regions to the clients. Next,

using these measurements, it generates a list per client by ranking the server regions from best performing to worst. Finally, these candidate lists are used as input to the server region assignment algorithm that matches one server region with each client.

Note that in practice server region assignment algorithm has various other inputs and constraints in addition to latency. For instance, delivery cost of traffic (that varies based on ISPs and the content), availability of the content in the servers, capacity of the servers, the allowance rules enforced by ISPs, the type of the content application are some of the constraints that effect the choice of a best available server for a given client. That is, a best available server for a client is not necessarily the least latency one. In fact, this is exactly the motivation behind ranking the servers instead of just finding the best performing server, i.e. the best performing one might not be the best available match in terms of the additional constraints. In practice, these constraints are applied after the ranking is performed. Their use and how the server matching is performed based on these constraints are not in the scope of this paper. Instead, we study the methods that lists the server rank lists based on path performance in a scalable fashion.

Figure 1 shows a simple example with six clients and four Akamai server regions. Both Akamai server regions and the clients are spread across the globe. For each client, the paths from server regions to the client are monitored as shown for client $X_6$ in the example. Let $\{34, 40, 80, 65\}$ be the latencies in milliseconds from the set of server regions $\{T_1, T_2, T_3, T_4\}$ to $X_6$, respectively. Then the ranking list of $X_6$ is $x_6 = \{1, 2, 4, 3\}$.

At the high-level, there are two main challenges in this context. First, due to the dynamic nature of the Internet, the mapping system works in real-time. Therefore rank lists need to be regenerated a few times in a minute. This requires constantly measuring path performance. Second, the scale of the problem is too large. There are millions of end-users and taking measurements from thousands of server regions to each of them is not feasible. Note that in the case of DNS-based mapping, end-users are represented by their DNS resolvers, i.e. measurements are taken between the server regions and the DNS resolvers to estimate the path performance to end-users behind each DNS resolver. Even in that case, the scale of the problem is too large since there are millions of DNS resolvers [2].

To overcome these challenges and make the monitoring scalable, the mapping system aims to cluster clients into groups and select one *representative* client per group such that measurements are taken for only the representative client. Then, all clients in the group are assigned to a server based on the ranking of the representative client.

## 2.2 Goals and Challenges

Our first goal in this work is finding a partitioning of the Internet address space so that, in each partition, the server region rankings of the clients are similar to each other. We expect a good partitioning to be stable, i.e. the clients do not migrate from one partition to the other one frequently over time. There are various ways of partitioning the address space, e.g. by mapping clients to their geography or autonomous systems. Although these mappings might work up to some extent, our intent is a partitioning that is driven by the network dynamics in order to capture the changes on the paths. For that reason, we propose using a clustering method that is based on the routing state in the Internet.

Our second goal is developing metrics to evaluate the goodness of a given partition. That is, given a pair of clients, our intent is quantifying their similarity in terms of ranking server regions. There are well-known metrics such as Kendall tau and Spearman footrule that computes the correlation between two rank vectors [6]. The assumption behind these correlation metrics is that the rank vectors are fully known, i.e. each region is assigned a rank position by each and every client. However, in practice, even for a single client, taking measurements from all the server regions to the client is not scalable. Therefore, the mapping system samples from the set of server regions for each client and take measurements only from this sampled subset. Such sampling is necessary in order not to overwhelm the clients with large number of requests. This results in learning the performance of the paths from only a subset of server regions per client. The subset vary for each client. For instance, in Figure 1, assume that measurements from only 3 server regions are available for some clients. Then, the known latencies and the corresponding rankings are as shown in Figure 2 (a) and (b), respectively. The empty entries in the tables are due to the unknown latencies. The missing values in the rank vectors make the well-known rank correlation metrics unsuitable. To overcome this challenge, in the next section, we propose metrics that measure similarity between server rankings in the case of unknown measurements.

Note that the actual server assignment algorithms is not in the scope of this paper. Instead, we study the server ranking problem which is an input to the server assignment algorithm. To that end, our goals are 1) finding a method that clusters IP addresses based on similarity in path performance they receive from server regions,

---

[2]Note that methods described in this paper are applicable to not only end-users but also DNS resolvers. Therefore, in the rest of the paper we will use the term *client* to refer to either *DNS resolver* (in the case of DNS-based mapping) or *end-user* (in the case of end-user mapping)

2) developing the metrics that evaluate the success of clustering, and 3) developing methods to select representative clients to scale the server ranking problem.

## 3. EVALUATION METRICS

Before we introduce our clustering method, we first propose two metrics to measure the similarity between two clients based on how they rank server regions. In the following sections, we use these metrics to evaluate how good (compact) clusters are.

A successful partitioning generates clusters where the server rankings of clients within a given cluster are similar to each other. We propose two similarity definitions: 1) we expect that in a good cluster, server regions are ranked in close positions. In other words, two rankings that are close to each other in $l_1$ norm should be grouped together. 2) top few rank positions are more important than the others since they have higher probability to be matched with the client. Therefore we expect that in a good cluster, all users have same server regions for their top few rank positions. That is, each rank position has a weight that is proportional to its order in the ranking. Note that this is a much more strict constraint than the previous definition and more sensitive to the missing measurements.

To capture each of these similarity definitions we propose two metrics, called PARTIAL SPEARMAN FOOTRULE DISTANCE (*ps-dist*) and GEOMETRIC DISTANCE (*g-dist*).

Let $x_1$ be a real-valued rank vector of length $m$ for client $X_1$, where $m$ is the total number of server regions. We define a ranking (or ordering) of its elements $\sigma_1$ such that $\sigma_1(i)$ is the rank position of element $i$ in the sorted $x_1$. That is $\sigma_1$ is a permutation of numbers from 1 to $m$. We say that in $x_1$, $i$ is preferred over $j$ if $\sigma_1(i) < \sigma_1(j)$.

Given two vectors $x_1$ and $x_2$ of same length, below we define two distance metrics between their ranking vectors $\sigma_1$ and $\sigma_2$, respectively.

### 3.1 Partial Spearman's Footrule Distance

One well-known distance metric for rankings is Spearman's footrule. It is the $l_1$ distance between two vectors $\sigma_1$ and $\sigma_2$ s.t. $\rho = \sum_{i=1}^{m} |\sigma_1(i) - \sigma_2(i)|$ [6]. By definition, Spearman's foot rule distance is maximum when the ordering in $\sigma_1$ is the reverse of the ordering in $\sigma_2$.

Note that Spearman's footrule requires the complete information on the rankings, i.e. for each entry $i$ of $x_1$ the ranking of $i$ must be known (likewise for $x_2$). In cases where $x_1$ and $x_2$ are only partially known, Spearman's foot rule can be modified as follows [12].

Let $\tau_1$ be the ranking of elements in $x_1$ such that $\tau_1(i)$ is the rank position of element $i$ only if $x_1(i)$ is known. Let $k$ be the number of known elements in $x_1$. Then for any unknown element $j$, one can set its rank to $l$, where $l > k$. Likewise $\tau_2$ is defined for $x_2$.

The intuition is that the servers whose latency are unknown are still considered (equally) but not preferred over the servers whose latency are known. This modified version of Spearman Footrule Distance have some nice properties (e.g. being a metric) as discussed in [12].

Then, the Partial Spearman's footrule distance between $x_1$ and $x_2$ is $PS-dist(x_1,x_2) = \sum_{i=1}^{m} |\tau_1(i) - \tau_2(i)|$.

In this work, we set all the unknown ranks to $l = k + 1$. Then, we normalize $PS-dist$ by $k \times (k+1)$ so that it is always between 0 and 1. We call the normalized distance as *ps-dist*.

For the example in Figure 2, $k = 3$. Therefore any unknown rank value is assigned 4. Then, $ps\text{-}dist(x_2,x_4) = \frac{8}{12}$ and $ps\text{-}dist(x_2,x_3) = \frac{6}{12}$.

Note that *ps-dist* computes the distance between two rankings without assigning weights to the rank positions. However, in some applications the higher rank positions matter more than the lower ones. Below we define another distance metric that assigns weights to the rank positions such that the distance between $x_1$ and $x_2$ is smaller when the elements in the higher rank positions are the same.

### 3.2 Geometric Distance

We define geometric distance (*g-dist*) between $x_1$ and $x_2$ as follows:

$$g\text{-}dist(x_1,x_2) = 1 - \sum_{i=1}^{m} I(i)\frac{1}{2^i} \qquad (1)$$

where $I$ is an indicator function s.t. $I(i) = 1$ if $x_1$ and $x_2$ both prefer the same element for the $i^{th}$ position, otherwise it is 0.

Note that $g\text{-}dist(x_1,x_2) \to 1$ as $m \to \infty$ and it drops proportionally to the importance of the rank position. For instance if the highest ranked element (i.e. rank-1) of $x_1$ and $x_2$ are the same then their distance is guaranteed to be less than or equal 0.5. Likewise, if both their rank-1 and rank-2 elements are the same then their distance is guaranteed to be less than or equal to 0.25. For the example in Figure 2, $g\text{-}dist(x_2,x_4) = 1$, and $g\text{-}dist(x_1,x_2) = \frac{1}{2}$.

One very important point to note is that we use neither *ps-dist* nor *g-dist* to cluster the clients. We only use them to evaluate the goodness of an already formed cluster. Although, our purpose is grouping clients whose server rankings are similar, the reason we do not cluster based on *ps-dist* and *g-dist* is as follows.

As we introduce in Section 2.1, the mapping system works at real time. That is, the paths need to be monitored a few times in a minute. However, the task of clustering itself should be run much less frequently. Therefore, the clusters should be stable over time. That is, the clients should not migrate from one cluster to the

other one frequently, at least until the next run of cluster generation. For that reason, the clusters should be formed based on a more stable metric than latency. We know that latency is prone to fluctuations due to many factors, such as queueing time, server response time etc. Moreover, in order to cluster based on *ps-dist* and *g-dist* we still need to know the latency between server regions and the clients which is the challenge that we tackle to solve in the first place. Therefore, we propose a clustering method that is based on inter domain routing which is more prone to frequent to fluctuations compared to latency on the paths.

Finally, in addition to *ps-dist* and *g-dist*, one can define a metric that considers the degree of the latencies instead of their orders. That is, one can categorize latency values as (really low, low, ok, high, really high) by defining a lower and upper boundary latency values for each category and then measure the similarity between categorical vectors. For brevity, in this paper, we only note that such similarity metric yields similar results to *ps-dist* and *g-dist*. We refer the reader to [18] for details.

## 4. DATASETS

In this study, we use traceroute measurements and BGP announcements that are collected in the Akamai's CDN. We collected each of them on two separate days, July 2, 2014 (Day-1) and January 24, 2016 (Day-2).

**1. Traceroute Measurements.** We collected traceroute measurements from Akamai server regions to local DNS clients. For Day-1, the measurements are taken from 2211 server regions to 20110 clients. For Day-2, the measurements are taken from 2073 server regions to 23004 clients. The Akamai server regions are spread across the globe. The DNS clients are located in six European countries (France, Germany, Spain, Italy, Switzerland, Belgium) and they belong to various ISPs.

In practice, there are limitations on the number of times a DNS client can be tracerouted at a given time period. Such limitations are set by the ISPs in order not to keep DNS clients busy. Therefore, each DNS client is tracerouted from a subset of the Akamai server regions. The number of server regions (known vector entries) per DNS client is 20 at minimum. Therefore we set $k = 20$ for *ps-dist*.

Each measurement from a server region to a DNS client consists of three consecutive ICMP packets and among these three, we use the one with the minimum latency. Using these latency values we generate a rank vector for each DNS client.

**2. BGP Announcements.** We use a collection of BGP tables collected from Akamai routers. For Day-1, the tables are collected from 233 peer routers and consist of over 1.7M BGP paths to over 37K prefixes located in the

six European countries (France, Germany, Spain, Italy, Switzerland, Belgium). Using this dataset, we map our DNS client IPs to their longest matching BGP prefixes. 20110 DNS servers in Day-1 map to 5491 unique prefixes.

For Day-2, the tables are collected from 297 peer routers and consist of over 790K BGP paths to over 48K prefixes located in our six European countries. Using this dataset, we map our DNS client IPs to their longest matching BGP prefixes. 23004 DNS clients map to 3272 unique prefixes.

## 5. BGP PREFIXES AS PRE-CLUSTERS

One can group clients by their longest matching network prefixes as advertised by the BGP system. Such grouping is atomic in terms of routing because BGP dictates that all clients matching a destination prefix are routed the same way at the inter-domain level.

Our intent is pre-grouping the clients by their prefixes [3] and using these prefixes as data points for clustering instead of using individual client IPs. Our aim is providing a faster and more efficient partitioning by reducing the number of data points to be clustered.

To that end, we ask whether pre-grouping clients by their prefixes yields good partitioning. That is, we test whether two clients from the same prefix are close to each other in terms of their server rankings. In order to do this, we first define *optimal partitioning* to set the benchmark for testing goodness of prefix clustering.

**Optimal Partitioning.** Let $X$ be the set of clients and $\mathscr{P}^r = \{\mathscr{P}_1^r, \mathscr{P}_2^r, \ldots, \mathscr{P}_n^r\}$ be a partitioning on $X$ such that every $X_i \in X$ belongs to one partition $\mathscr{P}_l^r$. Let $x_1$ and $x_2$ be the ranking vectors of $X_i$ and $X_j$, respectively. If two clients, $X_i$ and $X_j$, prefer the same servers in the same order as their top-$r$ choices then they belong to the same partition, i.e. $X_i \in \mathscr{P}_l^r$ and $X_j \in \mathscr{P}_l^r$.

For instance, in the example in Figure 2, for $\mathscr{P}^2$, $X_5$ and $X_6$ are in the same partition since they prefer $T_1$, $T_2$ in the same order as their top-2, whereas none of the other clients will be partitioned together as they all have different top-2 orderings.

The reason we call such partitioning *r-optimal* is that it does not allow two clients to be in the same partition if they don't agree on the same servers in the exact same order for their top-$r$ choices. Therefore it sets the tightest constraints and guarantees the most compact clusters for the particular choice of $r$.

Also notice that, by definition, optimal partitioning groups clients based on their *g-dist*. For instance, in a *1-optimal* partition, *g-dist* between any pair of client is not greater than 0.5. Similarly, in a *r-optimal* partition

---

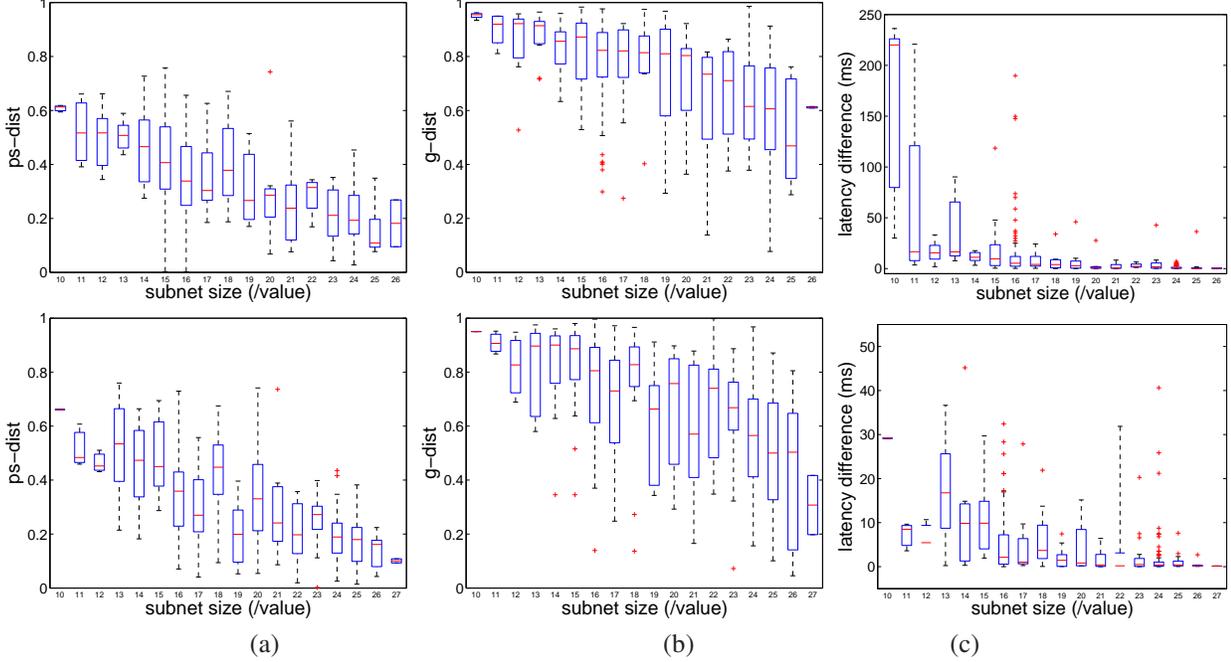[3]Throughout this paper, we will refer *longest matching BGP prefix* as *prefix*.

Figure 3: *ps-dist* (a), *g-dist* (b), latency difference (c) vs. BGP prefix length. In each box, the red line is the median, the upper and the lower end of the box are the first and third quantiles, respectively. The upper and lower ends of the whiskers are the maximum and minimum values in the data, respectively. The individual red points are outliers. The top row shows the results from `Day-1` and the bottom row shows the results from `Day-2`.
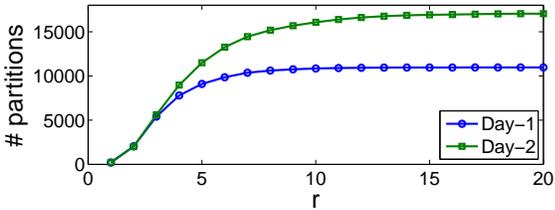


Figure 4: The number of optimal partitions for a range of *r* values.

for a slightly large value of *r*, *g-dist* between any pair of client is close to 0.

We apply optimal-partitioning on the set of 20110 (`Day-1`) and 23004 (`Day-2`) rank vectors that we described in Section 4. Figure 4 shows the number of partitions for *r* values from 1 to 20. The number of partitions increases as *r* increases. This is expected since for large values of *r*, there exists larger number of combinations for top-*r* rankings. For *r* values of 1,2,3, the number of partitions are around 200, 2000, 5400 respectively. For *r* values greater than 6, the number of partitions converges to a number slightly greater than 10000 (`Day-1`) and 14000 (`Day-2`). That is, on average, there exists only one or two clients under the same partition for *r* values greater than 6.

Next, for each partition $\mathscr{P}_i^r$, let $t_i^r$ be the number of clients in the partition and $\mu_l^r$ be the mean of all pairwise *ps-dist* values within the partition as written below:

$$\mu_l^r = \frac{\sum_{X_i,X_j:X_i,X_j \in \mathscr{P}_l^r} ps\text{-}dist(x_i,x_j)}{t_l^r * (t_l^r - 1)} \qquad (2)$$

Then $M^r = \{\mu_1^r, \mu_2^r, \ldots, \mu_n^r\}$ is the distribution of mean of pairwise distances for the partitioning $\mathscr{P}^r$.

**Comparing BGP prefix clusters with optimal partitioning.** One variable in BGP prefixes is the length, i.e. the number bits in the subnet mask. The longer a prefix, the less clients it has and the more likely that the clients are close to each other in the network.

In order to understand the role of the prefix length, first we cluster clients by their longest matching prefixes as found in the BGP dataset we described in Section 4. Note that each client is matched with only one prefix. Second, we compute the pairwise *g-dist* and *ps-dist* between each and every pair of clients that are from the same prefix. Then we compute the average of pairwise distances within each prefix. Third, we group these averages by the length of their prefixes. In our dataset, the length of BGP prefixes vary from /10 to /26 in `Day-1` and /10 to /27 in `Day-2`.

Figure 3 shows the statistics of each prefix length group as a separate box. The top row shows the results from `Day-1` and the bottom row shows the results from Day-

2. In Figure 3 (a) we see that there are client pairs from small prefixes (e.g. /19- /26) that are close to *1-optimal*. However, in large prefixes (e.g. /10, /11) the clients are far away from each other. In Figure 3 (b), we see the same trend, i.e. as the lengths of the prefixes increase, the distances within the prefix group decrease. This suggests that smaller subnets generate more compact clusters.

Next we investigate how large the latency can get for a client due to clustering. We compute the worst possible latency difference for each client in a given prefix as follows. For each client $X_i$ in a given prefix, we find the largest latency server with respect to $X_i$, say $T_j$ s.t. $T_j$ is the top-1 for some other client in the same prefix. Then, we compute how much the latency to $X_i$ increases if $T_j$ is assigned to $X_i$ instead of its top-1 server. Figure 3 (c) shows the average of such latency differences per prefix grouped by prefix length. The figure shows that the latency difference drops significantly as the prefix length grows and the difference is around 0 for small prefixes.

Next, we test grouping the clients by prefix against the optimal partitioning with a range of $r$ values. We compute the average pairwise *ps-dist* within each prefix (as described above). Then we divide the set of these average values into four groups by the length of their prefixes, /10-/15, /15-/18, /18-/24, /24-/26. We compare the distribution of values in each group with $M^r$ for $r = 1 \dots 5$ in Figure 5 and Figure 6 for Day-1 and Day-2, respectively. We see that the distributions of the prefixes from the /15-/18, /18-/24, and /24-/26 groups are very close the *2-optimal*, *3-optimal*, and *4-optimal*, respectively. To quantify the results, we run Kolmogorov-Smirnov tests to check the similarity of distributions from prefixes and their corresponding optimal partitions. The /18-/24 and /24-/26 groups passed the test at the 1% significance level. In addition, we note that 20110 clients in Day-1 map to 5491 prefixes and 23004 clients in Day-2 map to 3272 prefixes. That is, the mapping to BGP prefixes reduce the set of clients almost as much as *3-optimal* partitioning and better.

Finally, we compare the pairwise *ps-dist* and *g-dist* between clients that are from the same prefixes with the ones from different prefixes. We randomly sample 1000 client pairs that are from the same prefix and 1000 client pairs that are from different prefixes. Figure 7 shows the distribution of their pairwise distances. We see that for both metrics, the distance values between clients from the same prefixes are much lower compared to the the distance values between clients from different prefixes.

We conclude that grouping clients by their longest matching BGP prefixes generates successful clusters. Moreover, the similarity of server ranking within prefixes (especially the small ones) are close to the optimal. For
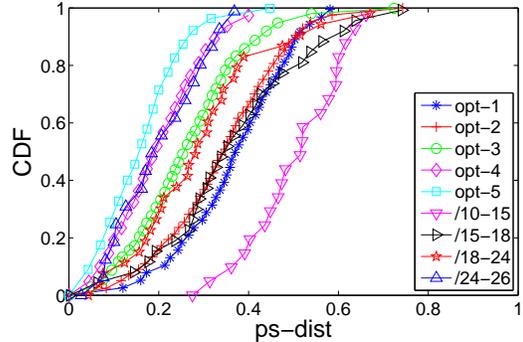


Figure 5: Day-1: Distributions of optimal partitioning for *r* values from 1 to 5 and avg. pairwise *ps-dist* within each prefix-length group.
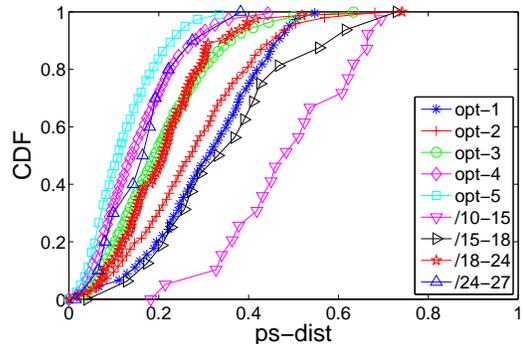


Figure 6: Day-2: Distributions of optimal partitioning for *r* values from 1 to 5 and avg. pairwise *ps-dist* within each prefix-length group.

that reason, one can use prefixes as data units for further clustering. We suggest dividing prefixes of small subnet lengths (e.g. /10- /12) into smaller subnets (e.g. /24s) in practice.

## 6. ROUTING-AWARE CLUSTERING

One factor that has large impact on the latency between two ends is the routing path between them. In this section, we apply a routing-aware clustering on the set of prefixes and study these clusters for the server ranking problem.

### 6.1 Correlation Between Routing State and Server Ranking

Our hypothesis is that if the routing paths to two destination prefixes, $p_1$ and $p_2$ are similar from a set of server regions, then $p_1$ and $p_2$ experience similar latencies from these server regions and therefore rank them similarly. To test this hypothesis, first we revisit the notion of *routing similarity* as introduced in [19].

Let $A$ be the set of all ASes in the Internet s.t. $A = \{a_1, a_2, \dots, a_t\}$ and $p_1$ and $p_2$ are announced by $a_1$ and $a_2$, respectively. We assume that for any $a_i$ there is a
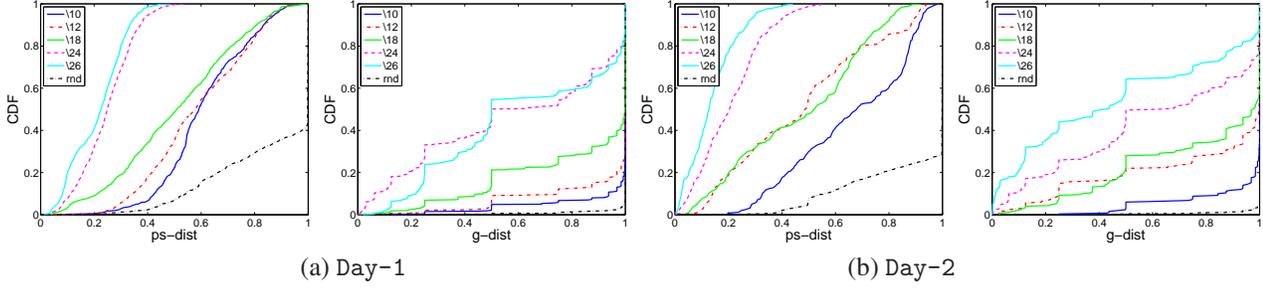
Figure 7: Distribution of pairwise distance values between 1000 pairs of clients that are from the same BGP prefix and 1000 pairs of clients that are from randomly selected, different BGP prefixes.

unique[4] $a_j$ which is the next hop AS on the path to $p_1$. That is, $nexthop(a_i,p_1) = a_j$, and $nexthop(a_1,p_1) = a_1$. Then, the *routing state* of $p_1$ is the collection of next hop choices to $p_1$ as described below:

$routestate(p_1) =$

$\quad \langle nexthop(a_1,p_1), nexthop(a_2,p_1),...,nexthop(a_t,p_1) \rangle$

Given $routestate(p_1)$ and $routestate(p_2)$, we measure the routing similarity between $p_1$ and $p_2$ by the number of ASes that prefer the same next hops to $p_1$ and $p_2$ as defined below:

$$rsim(p_1,p_2) = \#\{a_i \,|\, nexthop(a_i,p_1) = nexthop(a_i,p_2)\}$$

Similarly, the routing dissimilarity between $p_1$ and $p_2$ is called *Routing State Distance* (*rsd*) and it is defined below:

$$rsd(p_1,p_2) = \#\{a_i \,|\, nexthop(a_i,p_1) \neq nexthop(a_i,p_2)\}$$

**Applying *rsd* on the BGP dataset.** Having defined routing similarity, we seek to compute *rsd* for prefix pairs in our BGP dataset. However, there are two main issues. First, for a given prefix, we can not observe a *nexthop* from every AS in the Internet. To address this case when $nexthop(a_i,p_1)$ is not available, $rsd(p_1,p_2)$ is approximated by the fraction of *known* next hops in which $routestate(p_1)$ and $routestate(p_2)$ differ, times the total number of ASes in set $A$. This normalizes *rsd* so that it always ranges between zero and the total number of ASes in set $A$, i.e. $|A|$. We called the normalized version *RSD*.

The second issue is that for some AS-prefix pairs *nexthop* function is not uniquely defined, that is traffic destined for the same prefix may take different next hops e.g. when an AS uses hot-potato routing. We address this problem as in [17], i.e. by dividing each AS in set $A$ into a minimal set of sub-ASes such that for each (sub-AS, prefix) pair there is a unique next hop AS

---

[4]This assumption is relaxed later.

for the prefix. We call this extended version of $A$, $A'$. For the dataset of Day-1, $|A'| = 1460$, and for Day-2, $|A'| = 1404$. There are around 50 next hops in average per prefix in both days. Note that both implementation considerations are discussed in [17] in great detail.

Having addressed these two issues, we compute, *rsim* and *RSD* for each pair of prefixes in our BGP dataset. Notice that due to the normalization, unknown next hop, and multiple next hop issues described above, $RSD(p_1,p_2)$ is not simply $|A| - rsim(p_1,p_2)$.

Next, we seek to understand the correlation between routing similarity and server ranking of two prefixes. For each prefix pair we compute the minimum, average, and maximum *ps-dist* and *g-dist* between the clients of the prefixes as defined below. Note that the *g-dist* counterparts are defined likewise.

$$ps\text{-}dist_{min}(p_1,p_2) = \min_{\substack{x_1 \in p_1, \\ x_2 \in p_2}} ps\text{-}dist(x_1,x_2) \quad (3)$$

$$ps\text{-}dist_{max}(p_1,p_2) = \max_{\substack{x_1 \in p_1, \\ x_2 \in p_2}} ps\text{-}dist(x_1,x_2) \quad (4)$$

$$ps\text{-}dist_{avg}(p_1,p_2) = \frac{1}{|p_1||p_2|} \sum_{\substack{x_1 \in p_1, \\ x_2 \in p_2}} ps\text{-}dist(x_1,x_2) \quad (5)$$

Note that $ps\text{-}dist_{max}(p_1,p_1)$, $ps\text{-}dist_{min}(p_1,p_1)$ are the maximum and minimum *ps-dist* between two clients within $p_1$, and they can be written as $ps\text{-}dist_{max}(p_1)$ and $ps\text{-}dist_{min}(p_1)$, respectively. Similarly, $ps\text{-}dist_{avg}(p_1)$ is the average pairwise distances of clients in $p_1$. Note that *g-dist* counterparts are defined likewise.

Given a pair of prefixes, Figure 8 and Figure 9 show the relationship between the number of ASes that prefer the same next hops for these prefixes and their server ranking similarity. Prefix pairs are placed in buckets of size 10 according to their *rsim* values. Then, for each bucket, mean of the $ps\text{-}dist_{min}$, $ps\text{-}dist_{max}$, $ps\text{-}dist_{avg}$, $g\text{-}dist_{min}$, $g\text{-}dist_{max}$, $g\text{-}dist_{avg}$ values are plotted with 95% confidence interval.
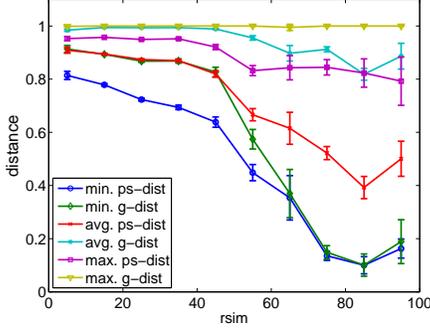
Figure 8: Day-1: *ps-dist* and *g-dist* vs. *rsim* between pairs of prefixes grouped by their *rsim* values.
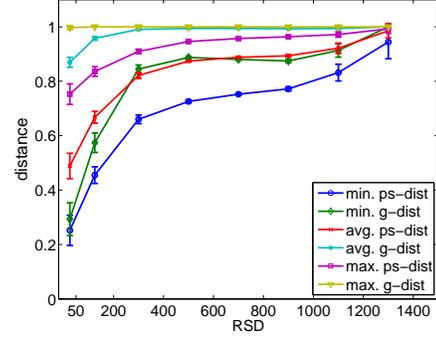


Figure 10: Day-1: *ps-dist* and *g-dist* vs. *RSD* between pairs of prefixes grouped by their *RSD* values.
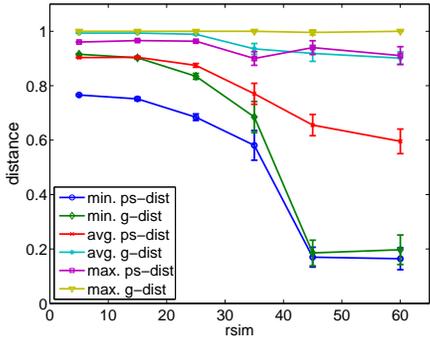


Figure 9: Day-2: *ps-dist* and *g-dist* vs. *rsim* between pairs of prefixes grouped by their *rsim* values.
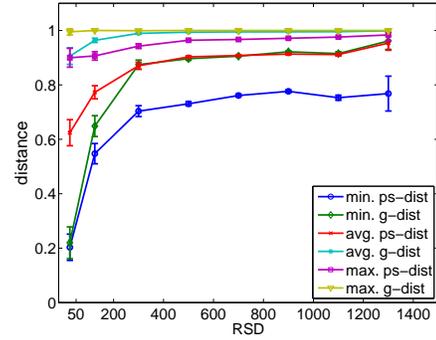


Figure 11: Day-2: *ps-dist* and *g-dist* vs. *RSD* between pairs of prefixes grouped by their *RSD* values.

Both Figure 8 and Figure 9 show that there is a strong correlation between the next hop choices for two prefixes and their server ranking. As the number of ASes that choose the same next hops increases, the server ranking distance between two prefixes decreases. The amount of decrease is larger for the *ps-dist_min* and *g-dist_min* metrics, i.e. there is at least a pair of clients, one from each prefix, that are very close to each other if the similarity between prefixes is high.

Similarly, in Figure 10 and Figure 11 prefix pairs are placed in buckets of 0-50, 50-200, 200-400, ... 1200-1400 according to their *RSD* values. Then, for each bucket, mean of the *ps-dist_min*, *ps-dist_max*, *ps-dist_avg*, *g-dist_min*, *g-dist_max*, *g-dist_avg* values are plotted with 95% confidence interval. Both figures show that as the *RSD* between two prefixes decreases, their max, min, and average *ps-dist* and *g-dist* decrease. In other words, two prefixes that are close to each other in the *RSD* space are close to each other in *ps-dist* and *g-dist* space too.

## 6.2  Clustering by Routing Similarity

Having shown the correlation between routing and server ranking similarity, next we seek clustering prefixes in *RSD* space. Intuitively, we are looking for a par-

titioning that minimizes the *RSD* between two prefixes that are from the same cluster. One extended version of this intuition is formalized as RS-CLUSTERING problem in [19] and solved by Algorithm 1.

---

**Algorithm 1** The Pivot algorithm .

---

A set of prefixes $P = \{p_1, \ldots, p_n\}$ and a threshold $\tau \in [0, t']$.
A partition $\mathscr{P}$ of the prefixes
1: pick a random prefix $p \in P$
2: create a cluster $C_p = \{p' \mid RSD(p, p') \leq \tau\}$
3: $P = P \setminus C_p$
4: Pivot$(P, \tau)$

---

The inputs of the algorithm are the set of prefixes, their pairwise *RSD* values, and a threshold parameter $\tau \in [0, t']$, where $t' = |A'|$ is the maximum possible value of *RSD*. The algorithm works as follows: Starting from a random prefix $p$, it finds all prefixes that are within the distance $\tau$ from $p$. All these prefixes are assigned in the same cluster, $C_p$ – centered at prefix $p$. We call $p$ the *pivot* of cluster $C_p$. Then the prefixes that are assigned to $C_p$ are removed from the set of prefixes $P$ and the
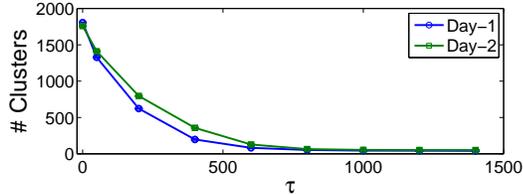
9

Figure 12: Number of clusters vs. $\tau$

`Pivot` algorithm is reapplied to the remaining subset of prefixes that have not been assigned to any cluster.

**Scaling with** RS-CLUSTERING. Notice that `Pivot` algorithm is a nonparametric algorithm where the number of clusters are not pre-set. In fact, the number of clusters is an output which is controlled by the choice of $\tau$. Setting $\tau$ to a larger value is likely to decrease the number of clusters but increase in-cluster *RSD*. That property of the `Pivot` algorithm provides flexibility in adjusting the scale of the problem in practice. Figure 12 shows the number of clusters we get by applying the algorithm on the 5491 prefixes (`Day-1`) and 3272 prefixes (`Day-2`) in our datasets. For each choice of $\tau$, we run the algorithm 10 times and show the average number of clusters for these runs with 95% confidence interval. The figure suggests that up to 600, as $\tau$ increases, the number of clusters decreases. Figure 10 and Figure 12, together, show the trade-off between the number of clusters and the closeness between prefixes within a given cluster. For the rest of the evaluation in this paper we set $\tau$ to 200. This results in around 600 clusters, i.e. reduces the scale of the problem by 90% for `Day-1` and by 82% for `Day-2`. The impact of $\tau$ and other properties of `Pivot` are discussed in [19, 15] in further detail.

Next, we analyze the goodness of the clusters. In order to do that, we first introduce the following notations. For a client $X_1 \in p_i$, if $p_i \in C_l$, then we write $X_1 \in C_l$. In addition, we use the definitions $ps\text{-}dist_{max}$, $ps\text{-}dist_{min}$, $ps\text{-}dist_{avg}$, and their $g\text{-}dist$ counterparts that are introduced in Section 6.1.

**Evaluating** RS-CLUSTERING . We evaluate the goodness of a cluster with two metrics: (1) *Growth of the diameter:* Diameter of a cluster is the maximum distance between any two clients of the cluster. We measure the growth of the diameter as the ratio of the maximum server ranking distance between any clients in the RS-cluster over the maximum diameter of the prefixes that are in that RS-cluster. Diameter growth of a cluster $C_l$, $dg(C_l)$ is defined below for *ps-dist*. Note that the *g-dist* counterpart is defined likewise.

$$dg_{ps\text{-}dist}(C_l) = \frac{\max\limits_{X_1,X_2 \in C_l} ps\text{-}dist(x_1,x_2)}{\max\limits_{p_i \in C_l} ps\text{-}dist_{max}(p_i)} \quad (6)$$

(2) *Growth of the average pairwise distances:* It is the ratio between average pairwise distances within a cluster and the maximum average pairwise distances of its member prefixes. Let $s_l$ be the number of clients in $C_l$. Average pairwise growth of a cluster $C_l$, $ag(C_l)$, is defined below for *ps-dist*. Note that the *g-dist* counterpart is defined likewise.
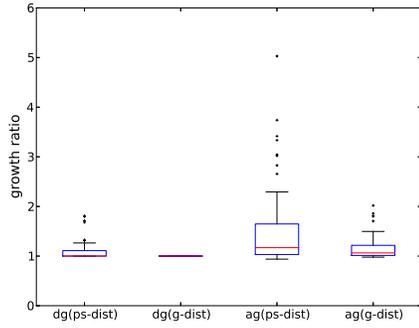
$$ag_{ps\text{-}dist}(C_l) = \frac{\frac{1}{s_l*(s_l-1)} \sum_{X_1,X_2 \in C_l} ps\text{-}dist(x_1,x_2)}{\max\limits_{p_i \in C_l} ps\text{-}dist_{avg}(p_i)} \quad (7)$$

By definition, when *dg* and *ag* values for a cluster $C_l$ are near 1, it means the diameter and average pairwise distances within $C_l$ did not grow further than the diameter and the average pairwise distances of its least compact prefix, respectively. Figure 13 shows *dg* and *ag* statistics for each cluster. The figure shows that the *dg* and *ag* values for both *ps-dist* and *g-dist* are around 1, that is the clusters did not grow out of their prefixes significantly. That means RS-CLUSTERING generates compact clusters with respect to its data units.
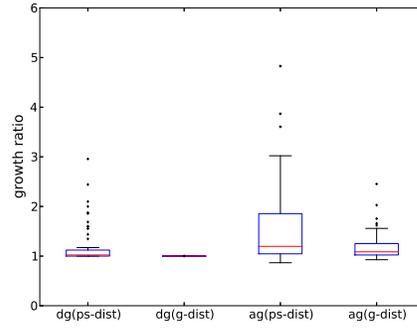
Next, in order to test the goodness of clusters further, we compare the average pairwise *ps-dist* and *g-dist* across prefixes that belong to the same cluster with the ones that belong to different clusters. We call the former in-cluster and the latter out-cluster distances. Figure 15 shows that there is a great difference in the distribution of distances between prefixes that are from the same cluster compared to the ones that are from different clusters.

In order to investigate the cluster statistics even further, we look at top 10 clusters in more detail. For each prefix $p_1 \in C_l$, we compute $ps\text{-}dist_{avg}(p_1,p_2)$ for all $p_2 \in C_l$. Also, we compute $ps\text{-}dist_{avg}(p_1,p_3)$ for all $p_3 \notin C_l$. We call the the former in-cluster, and the latter out-cluster. The hypothesis is that the average distance of a prefix should be much closer to another prefix that is in the same cluster compared to the ones that are outside the cluster. We compute in-cluster and out-cluster averages for all prefixes in the largest 10 clusters from `Day-1` and plot the statistics in Figure 14. The figure shows that across all clusters, in-cluster distances are lower compared to the out-cluster distances.

In addition Table 1 shows some statistics for these 10 clusters including their sizes (the number of prefixes per cluster), average *RSD* values between prefixes, the number of unique ASes that the prefixes in the clusters belong to, and the geo locations of the prefixes in each cluster. One thing to note is that for 7 out of these 10 clusters, the prefixes are from one single country. Also note that, the variety of ASes within a cluster ranges from 1 to 23. In fact, looking closer into the set of all clusters, we find that 80% of the clusters are com-

10

(a) Day-1                                     (b) Day-2

Figure 13: In each box, the red line is the median distance, the upper and the lower end of the box represent the first and third quantile, respectively. The upper and lower ends of the whiskers represent the maximum and minimum values in the data, respectively. The individual red points are outliers.
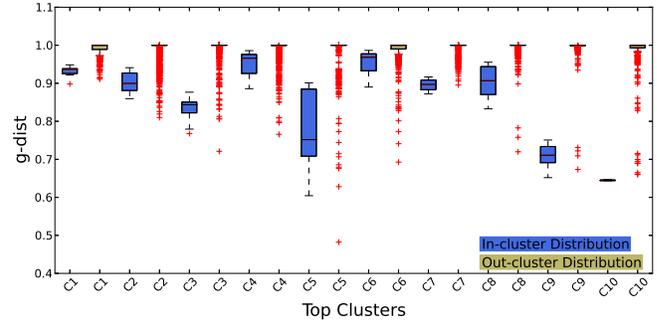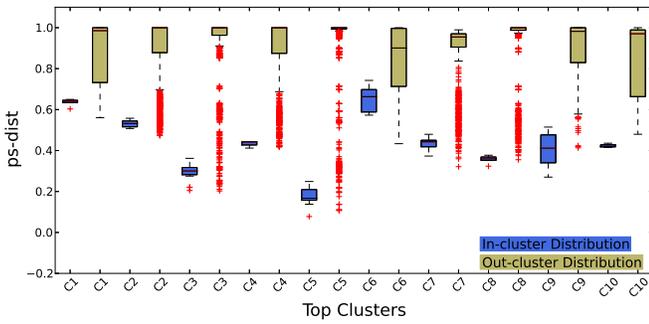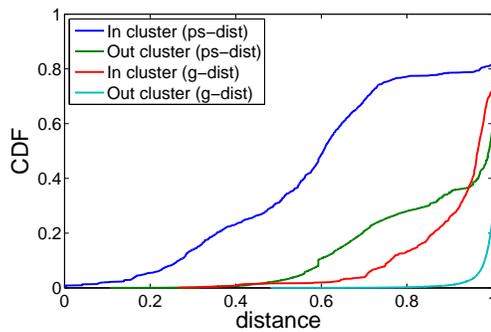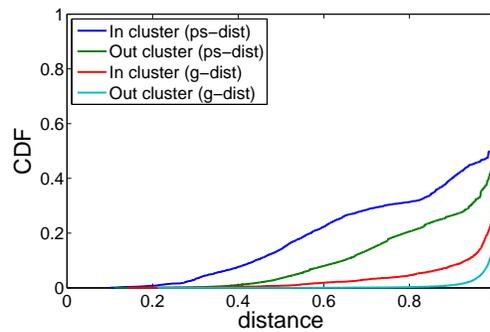


Figure 14: (a) *ps-dist* (b) *g-dist*. In each box, the red line is the median distance, the upper and the lower end of the box represent the first and third quantile, respectively. The upper and lower ends of the whiskers represent the maximum and minimum values in the data, respectively. The individual red points are outliers.



(a) Day-1                                     (b) Day-2

Figure 15: in-cluster vs. out-cluster distances

11

posed of prefixes from the same country and 38% of them are composed of prefixes from the same AS. To that end we ask the following question next : how does RS-CLUSTERING compares with clustering prefixes by their ASes or geographic locations ?

## 6.3 RS-CLUSTERING **vs. clustering by AS and Geography**

First we compare RS-CLUSTERING with clustering by country. We map each prefix to its country[5] . Then, for each prefix, we compute $ps\text{-}dist_{avg}$ and $g\text{-}dist_{avg}$ with every other prefix in its RS-cluster and country cluster separately. Then, we take the mean of these averages and plot them in Figure 16 (a-b) for Day-1 and in Figure 17 (a-b) for Day-2. Each point in the figure represents a prefix. We plot $x = y$ line for comparison, i.e. if the prefix is above the line it indicates that the prefix is closer to the other prefixes in its RS-cluster than the ones in its country cluster. Figure 16 (a-b) and Figure 17 (a-b) show that clustering by routing similarity is not same as clustering by geo location. They also show that clustering by routing similarity results in better clusters as the majority of the prefixes are above the $x = y$ line both in $ps\text{-}dist$ and $g\text{-}dist$.

Second we compare RS-CLUSTERING with clustering by AS. We map each prefix to its AS. The set of 5491 prefixes from Day-1 map to 1397 unique ASes. The set of 3272 prefixes from Day-1 map to 1206 unique ASes. For each prefix, we compute $ps\text{-}dist_{avg}$ and $g\text{-}dist_{avg}$ with every other prefix in its RS-cluster and AS-cluster separately. Then, we take the mean of these averages and plot them in Figure 16 (c) and Figure 17 (c). Each point in the figures represents a prefix. We plot $x = y$ line for comparison, i.e. if the prefix is above the line that indicates the prefix is closer to the other prefixes in its RS-cluster than the ones in its AS-cluster. In Figure 16 (c) and Figure 17 (c), we see that almost all prefixes are above the $x = y$ line for both $ps\text{-}dist$ and $g\text{-}dist$. In fact, we see a group of prefixes that are on the y-axis. For any one of these prefixes, the $ps\text{-}dist$ (or $g\text{-}dist$) between itself and the others in its RS-cluster is 0, whereas the $ps\text{-}dist$ (or $g\text{-}dist$) between itself and the others in its AS cluster is greater than 0.4.

In summary, Figure 16 and Figure 17 show that RS-CLUSTERING is different than clustering by geo and AS. In fact, RS-CLUSTERING outperforms clustering clients by their countries or ASes. In addition, by RS-CLUSTERING we can control the number of clusters to be generated, whereas clustering by AS and geo fix the number of clusters by definition. In that sense, RS-CLUSTERING is a flexible and accurate way of clustering.

---

[5]We use Akamai's EdgeScape tool for IP to geo mapping [11]

## 7. FINDING REPRESENTATIVE CLIENTS

Having partitioned the address space, we seek to find a *representative* client from each partition such that instead of measuring the paths to all clients in the partition, we only measure the ones that are to the representative. Then we rank the servers according to that representative client.

One can apply various methods to select a representative. For instance, one method is randomly choosing a client from the each cluster. The assumption is that if the cluster is compact enough, any client will be a good representative. Second method is choosing a client from the *pivot* prefix of each cluster as described in Section 6. Pivot prefix is the prefix that is the center of its cluster, i.e. it is guaranteed that all other prefixes in the cluster are at most $\tau$ away in *RSD* space. Given the high correlation between *RSD* and *ps-dist* (*g-dist*) we expect that clients from the pivot prefix are good representatives.

In order to investigate the effectiveness of these methods, first we find a single ranking for each cluster that best describes the rankings of all clients in the cluster. In other words, we find a *consensus* for each cluster and evaluate the goodness of the representative client by testing it against the *consensus*.

The problem of aggregating a set of rank vectors and finding consensus is known as *rank aggregation* problem [1, 10]. The problem is studied extensively and proposed solutions are subject to the definitions of what properties the consensus should have. For our application, we employ two of the proposed solutions, the Borda count [3] and the Plurality method [10].

The Borda count is a score-based method. Each candidate's (server's) score is the sum of the rank values assigned by every voter (client). Once all votes are counted, the candidates are reordered from the lowest to the highest rank and the lowest rank candidate is the winner. The nice property of Borda count is that a Spearman Footrule optimal solution can be computed in polynomial time [10] and therefore it sets a fair consensus ranking for our *ps-dist* metric.

The Plurality method is another voting scheme where the candidates are simply ordered by the number of rank values they are assigned to. For instance, a candidate with the most rank 1 assignment is ranked 1, similarly, a candidate with the most rank 2 assignment is ranked 2 and so on. By definition, this method sets a fair consensus for our *g-dist* metric.

To test how good the representatives are, we first compute the consensus vector (by both the Borda count and the Plural Method separately). Then, we compute the *ps-dist* (*g-dist*) between the consensus and all clients in its cluster, and compute their average. Next, we select a

Table 1: Statistics for 10 large clusters of `Day-1`.

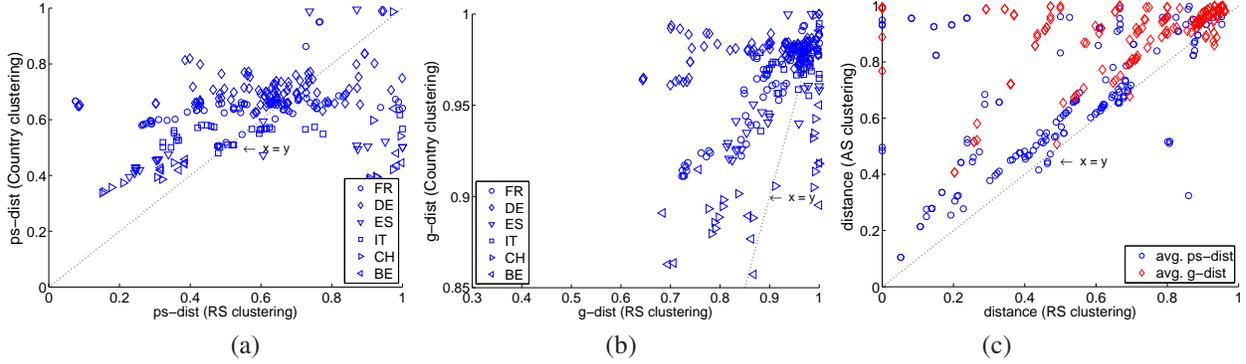| | $C_1$ | $C_2$ | $C_3$ | $C_4$ | $C_5$ | $C_6$ | $C_7$ | $C_8$ | $C_9$ | $C_{10}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| Size of cluster ($C$) | 102 | 102 | 102 | 85 | 64 | 52 | 48 | 35 | 26 | 18 |
| Avg. *RSD* ($C$) | 59.73 | 398.36 | 16.99 | 167.31 | 108.27 | 83.88 | 93.51 | 42.65 | 31.42 | 88.39 |
| Countries | FR | CH | DE | DE | FR,ES,IT,BE | FR | DE | DE | DE,IT | FR,DE |
| num. unique ASes | 23 | 6 | 10 | 16 | 17 | 1 | 7 | 4 | 5 | 2 |



(a)    (b)    (c)

Figure 16: `Day-1`: Comparing RS-CLUSTERING with clustering by country for (a) *ps-dist* (b) *g-dist*. (c) Comparing RS-CLUSTERING with clustering by AS for *ps-dist* and *g-dist*.
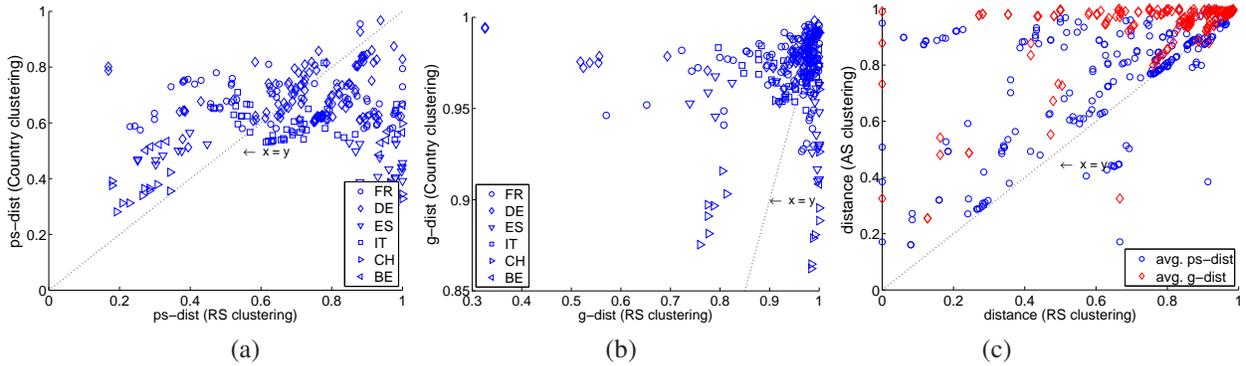


(a)    (b)    (c)

Figure 17: `Day-2`: Comparing RS-CLUSTERING with clustering by country for (a) *ps-dist* (b) *g-dist*. (c) Comparing RS-CLUSTERING with clustering by AS for *ps-dist* and *g-dist*.



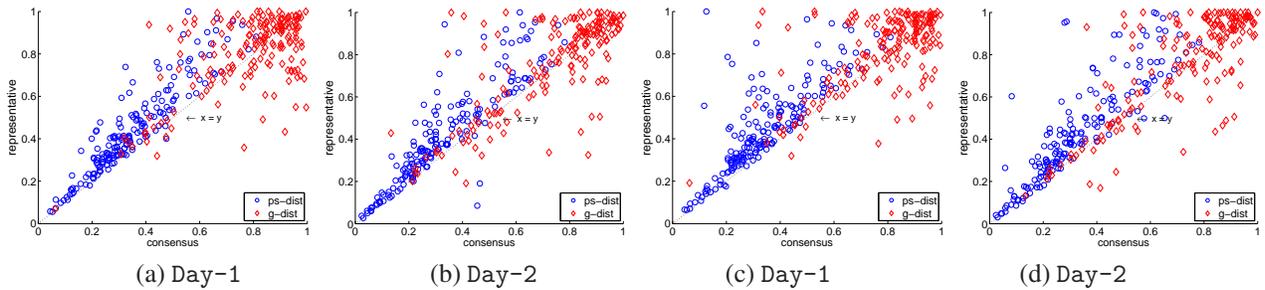(a) `Day-1`    (b) `Day-2`    (c) `Day-1`    (d) `Day-2`

Figure 18: The comparison of the average distances between (Borda count) consensus vector and all clients in the cluster with the average distances between representative client and all clients in the same cluster. The representative client is selected at random (a-b) from all clients in the cluster, (c-d) from the pivot prefix.

representative for the cluster and we compute the *ps-dist* (*g-dist*) between the representative and all clients in its cluster, and compute their average.

Figure 18 compares the consensus chosen by the Borda count with the representatives chosen (a-b) at random from the set of all clients in the cluster, and (c-d) at random from the clients that belong to the pivot (center) prefix of each cluster. Each point in the figure represents one cluster and the x and y axis are the average distances as described above. We also plot $x = y$ line for comparison, i.e. if the cluster is close to the line then it means that the representative client performs as good as the consensus ranking. The figure shows that most of the clients are very close to the line. The representative client performs slightly better for *ps-dist*. This is expected by the definition of the Borda count and its property of having an spearman foot rule optimal solution.

Figure 19 is generated exactly the same as described above but the consensus is chosen by the Plural method. The figure shows that again, most of the clients are very close to the $x = y$ line. The representative clients perform better for *g-dist* compared to the Borda count case. This is intuitive by the definition of the Plural method which simply orders the servers by the number of rank values they are assigned to. In other words, it tends to agree with the exact rankings of the majority for each position.

Next we investigate how large the latency can get for a client due to clustering. We compute the latency difference for each client in a given cluster as follows. For a given cluster, let $T_j$ be the top-1 server with respect to the cluster's representative client. Then, for any other client in the cluster, say $X_i$, we compute how much the latency to $X_i$ increases if $T_j$ is assigned to $X_i$ instead of its top-1 server. Figure 20 shows the distribution of such latency for all clients in the clusters. The figure shows that selecting representatives both random from the clients in the pivot prefix or any client in the cluster results in small latency difference for the rest of the clients. In fact, selecting representative at random from all clients slightly outperforms selecting from the pivot for Day−1. This is expected since the pivot prefix is not necessarily the largest one. In conclusion, we show that choosing a client at random from a cluster represents the other clients in the same cluster successfully and reduce the scale of the ranking task.

## 8. RELATED WORK

Grouping IP addresses has been of interest to various studies. [20] and [2] propose that IPs that are numerically close to each other to be grouped together. [20] clusters web client IPs by mapping them into their longest matching BGP prefixes so that the IPs within a

given cluster is under the same administrative domain. [2] assumes that IPs that are numerically close to each other show similar features (e.g. latency). Based on this assumption this work partitions a subnet as a supervised learning task so that in each subpartition the feature at interest is minimized. This work does not consider BGP routing information which already provides a natural grouping. [4] proposes grouping end-user IPs by their /24 prefixes in the case of end-user mapping. Unlike our work, in [20, 2, 4], the level of aggregation is limited to the prefix level or below and the state of the inter domain routing in the Internet is not considered.

Clustering IP addresses by their geographic locations is another way to partition the address space. However, [16] and [28] show that identifying the geolocation of an IP block accurately is hard. They claim that most geolocation tools are only reliable at the country level. Moreover, even in the case where geolocation based grouping is possible, it does not generate the desired clustering results we discuss in this paper, i.e. geographic closeness does neither infer topological closeness [14] nor low latency between servers and clients [21].

[19] proposes the *RSD* metric and the Pivot clustering to group prefixes according to the BGP path information. This work focuses more on the basic properties of the metric and shows how to uncover the factors that drive ASes to choose their interodomain next hops. However, unlike our work, this work does not study the goodness of clusters in terms of the path performances experienced by the clients in each cluster or how similarly the clients in a given cluster orders some set of servers.

In addition, the correlation between the routing dynamics and end-to-end path performance is widely studied. The findings in these studies suggest that we use a routing-aware clustering. [30] shows the routing and latency relation by comparing the default routing paths with alternate paths. [22] studies that the stability of paths between ISPs effect path performance. [33] shows that packet loss is significantly increased by the routing changes. [25] and [13] show that routing instabilities (e.g. routing loops and failures) can disrupt end to end connectivity.

Incorporating routing information into server redirection problem in CDNs is studied by [26, 27]. These works suggest that collaboration between CDNs and ISPs will be beneficial for mapping end-users to higher performing CDN servers. To that end, similar to our work, they also utilise what can be inferred from routing choices. However, notice that, the problem of server redirection studied in [26, 27] is different than the server ranking problem that we study in this paper. Our goal is *not* finding a best server mapping for end-users. Instead, our goal is grouping end-users such that within each group
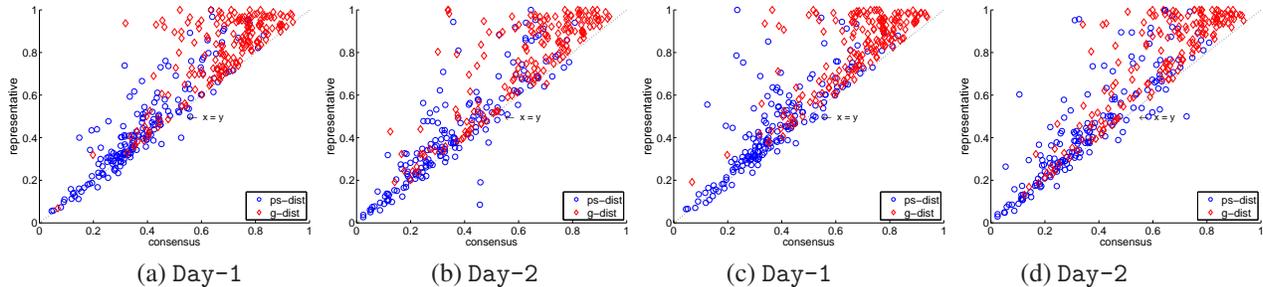
(a) Day-1      (b) Day-2      (c) Day-1      (d) Day-2

Figure 19: The comparison of the average distances between (Plural method) consensus vector and all clients in the cluster with the average distances between representative client and all clients in the same cluster. The representative client is selected at random (a-b) from all clients in the cluster, (c-d) from the pivot prefix.
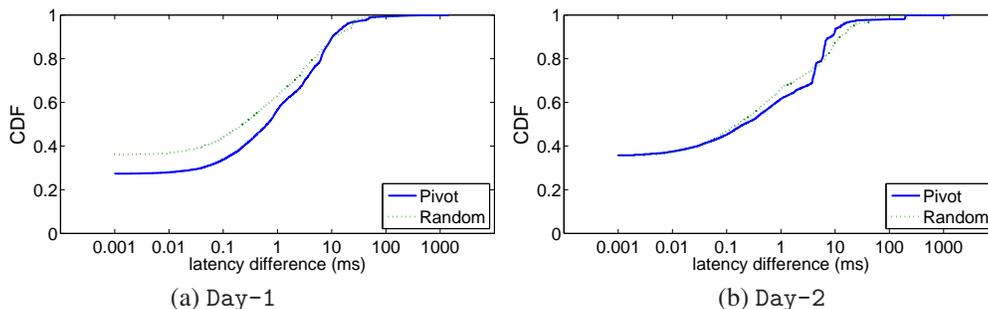


(a) Day-1           (b) Day-2

Figure 20: The distribution of latency differences due to assigning top-1 server of the representative

the order of best to worst performing servers are almost the same. In that respect, our work is complementary to [26, 27].

## 9. DISCUSSION

The RS-CLUSTERING has properties which makes it desirable for our clustering problem.

First, it is a lightweight algorithm. Notice that both computing pairwise *RSD* values and running Pivot algorithm is in $O(n^2)$ where $n$ is the number of prefixes. It is possible to reduce the run time of both jobs by pre-grouping the prefixes by some coarse grain geographic regions (e.g. their continents or even countries as we discussed in Section 6.2). For instance, it is very likely that for two clients from two different continents, neither their server choices nor the *RSD* of their prefixes will be close. Therefore, one can compute *RSD* and run the Pivot separately for each coarse geographic area.

Second, RS-CLUSTERING is a network-aware method that can capture the dynamicity of the Internet, yet it is relatively stable. [29] shows that small fraction of prefixes are responsible for most route changes and these are the ones that receive comparatively little traffic, whereas BGP is stable for popular destination prefixes. Moreover, [5] studies the temporal aspects of *RSD* and show that on any given day, approximately 1% of the next-hop decisions made in the Internet change. The

change goes up to 10% in a month. This shows that the clusters generated by *RSD* are valid for at least a day (possibly longer).

Third, RS-CLUSTERING provides flexibility for adjusting the number of clusters as we discuss in Section 6.2. The choice of $\tau$ in Pivot effects the number of the clusters generated. In this paper, we set $\tau$ empirically.

Fourth, the only input that RS-CLUSTERING requires is the set of BGP paths. Notice that RS-CLUSTERING does neither rely on the knowledge of the underlying topology nor the latency on the paths.

Fifth, in our experiments we show that RS-CLUSTERING reduces the problem by 90% for a region in Europe. We note that for other regions in the world, the percentage of reduction may vary based on how diverse the Internet paths in those regions. Remember that by definition of *RSD* the more ASes make similar next hop choices to a set of destination prefixes, the smaller their *RSD* values. The smaller the *RSD* values are, the more prefixes can be grouped together and the more the scale of the problem is reduced. Within a given region, one reason that ASes make similar next hop choices is because the next hop options are limited in the first place. In other words, in regions where path diversity is less, we expect to represent prefixes with only a few clusters and therefore reduce the problem significantly. In our study, we pick a region in Europe where the path diversity is relatively

15

rich compared to the other parts of the world. Therefore, we believe that the gain from RS-CLUSTERING will be more than 90% in other regions.

One direction for further analysis is the change in *ps-dist* and *g-dist* within RS clusters over time. We believe that such study can have various applications. For instance, one can identify the problematic clients/prefixes by observing the unexpected changes in the server rankings within an RS cluster.

## 10. CONCLUSION

In this paper, we introduce a framework to partition the Internet address space. Our goal is to scale the number of paths to be monitored between a CDN's servers and its clients. That is we aim to find a partitioning of clients such that in each partition the latencies from servers to the clients in the partition are ordered similarly. To achieve this goal, we first introduce two metrics (*ps-dist* and *g-dist*) that measures the similarity between two rank vectors even in the case of vectors are known only partially. Second, we show that for any given two clients, as the number of ASes in the Internet that prefer the same next hops to route to them increase, their server ranking similarity increase. Having shown the effect of inter domain routing on the server preferences of clients, we employ routing-aware clustering algorithm. We evaluate the goodness of the clusters by using our metrics *ps-dist* and *g-dist* and show that we obtain compact clusters. Finally, we show that one can successfully scale the task of server ranking by measuring a client at random from each cluster.

## 11. ACKNOWLEDGEMENTS

## 12. REFERENCES

[1] Nir Ailon, Moses Charikar, and Alantha Newman. Aggregating inconsistent information: Ranking and clustering. *J. ACM*, 55(5):23:1–23:27, November 2008.

[2] Robert Beverly and Karen Sollins. An internet protocol address clustering algorithm. In *Proceedings of the Third Conference on Tackling Computer Systems Problems with Machine Learning Techniques*, SysML'08, pages 5–5, Berkeley, CA, USA, 2008. USENIX Association.

[3] Jean C. de Borda. Mèmoire sur les èlections au scrutin. 1781.

[4] Fangfei Chen, Ramesh K. Sitaraman, and Marcelo Torres. End-user mapping: Next generation request routing for content delivery. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*, SIGCOMM '15, pages 167–181, New York, NY, USA, 2015. ACM.

[5] Giovanni Comarela, Gonca Gürsun, and Mark Crovella. Studying interdomain routing over long timescales. In *Proceedings of the Internet Measurement Conference (IMC)*, Barcelona, Spain, October 2013.

[6] P. Diaconis. *Group representations in probability and statistics.*, volume 11 of *IMS Lecture Notes-Monograph Series*. Institute of Mathematical Statistics, 1988.

[7] John Dilley, Bruce Maggs, Jay Parikh, Harald Prokop, Ramesh Sitaraman, and Bill Weihl. Globally distributed content delivery. *IEEE Internet Computing*, 6(5):50–58, September 2002.

[8] Google Public DNS. https://developers.google.com/speed/public-dns/.

[9] Open DNS. https://www.opendns.com.

[10] Cynthia Dwork, Ravi Kumar, Moni Naor, and D. Sivakumar. Rank aggregation methods for the web. In *Proceedings of the 10th International Conference on World Wide Web*, WWW '01, pages 613–622, New York, NY, USA, 2001. ACM.

[11] Akamai EdgeScape. http://www.akamai.com/dl/brochures/edgescape.pdf.

[12] Ronald Fagin, Ravi Kumar, and D. Sivakumar. Comparing top k lists. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '03, pages 28–36, Philadelphia, PA, USA, 2003. Society for Industrial and Applied Mathematics.

[13] Nick Feamster, David G Andersen, Hari Balakrishnan, and M Frans Kaashoek. Measuring the effects of internet path faults on reactive routing. In *ACM SIGMETRICS Performance Evaluation Review*, volume 31, pages 126–137. ACM, 2003.

[14] Michael Freedman, Mythili Vutukuru, Nick Feamster, and Hari Balakrishnan. Geographic Locality of IP Prefixes. In *Internet Measurement Conference (IMC) 2005*, Berkeley, CA, October 2005.

[15] Aristides Gionis, Heikki Mannila, and Panayiotis Tsaparas. Clustering aggregation. *Transactions on Knowledge Discovery from Data*, 1(1), 2007.

[16] Bamba Gueye, Steve Uhlig, and Serge Fdida. Investigating the imprecision of ip block-based geolocation. In *Proceedings of the 8th International Conference on Passive and Active Network Measurement*, PAM'07, pages 237–240, Berlin, Heidelberg, 2007. Springer-Verlag.

[17] Gonca Gürsun. Inferring hidden features in the internet. In *PhD Thesis, Boston University*, 2013.

[18] Gonca Gürsun. Routing-aware partitioning of the internet address space for server ranking in cdns. Technical report, Ozyegin University, 2015.

[19] Gonca Gürsun, Natali Ruchansky, Evimaria Terzi, and Mark Crovella. Routing state distance: A path-based metric for network analysis. In *Proceedings of the 2012 ACM Conference on Internet Measurement Conference*, IMC '12, pages 239–252, New York, NY, USA, 2012. ACM.

[20] Balachander Krishnamurthy and Jia Wang. On network-aware clustering of web clients. In *Proceedings of the Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication*, SIGCOMM '00, pages 97–110, New York, NY, USA, 2000. ACM.

[21] Rupa Krishnan, Harsha V. Madhyastha, Sridhar Srinivasan, Sushant Jain, Arvind Krishnamurthy, Thomas Anderson, and Jie Gao. Moving beyond end-to-end path information to optimize cdn performance. In *Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement Conference*, IMC '09, pages 190–201, New York, NY, USA, 2009. ACM.

[22] Craig Labovitz, Abha Ahuja, Abhijit Bose, and Farnam Jahanian. Delayed internet routing convergence. In *Proceedings of the Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication*, SIGCOMM '00, pages 175–187, New York, NY, USA, 2000. ACM.

[23] Erik Nygren, Ramesh K. Sitaraman, and Jennifer Sun. The akamai network: A platform for high-performance internet applications. *SIGOPS Oper. Syst. Rev.*, 44(3):2–19, August 2010.

[24] The Akamai State of the Internet Report. http://www.akamai.com/stateoftheinternet.

[25] Vern Paxson. End-to-end routing behavior in the internet. In *Conference Proceedings on Applications, Technologies,*

*Architectures, and Protocols for Computer Communications*, SIGCOMM '96, pages 25–38, New York, NY, USA, 1996. ACM.

[26] Ingmar Poese, Benjamin Frank, Bernhard Ager, Georgios Smaragdakis, Steve Uhlig, and Anja Feldmann. Improving Content Delivery with PaDIS. *IEEE Internet Computing*, 16(3):46–52, May-June 2012.

[27] Ingmar Poese, Benjamin Frank, Georgios Smaragdakis, Steve Uhlig, Anja Feldmann, and Bruce Maggs. Enabling content-aware traffic engineering. *SIGCOMM Comput. Commun. Rev.*, 42(5):21–28, September 2012.

[28] Ingmar Poese, Steve Uhlig, Mohamed Ali Kaafar, Benoit Donnet, and Bamba Gueye. Ip geolocation databases: Unreliable? *SIGCOMM Comput. Commun. Rev.*, 41(2):53–56, April 2011.

[29] Jennifer Rexford, Jia Wang, Zhen Xiao, and Yin Zhang. Bgp routing stability of popular destinations. In *Proceedings of the 2Nd ACM SIGCOMM Workshop on Internet Measurment*, IMW '02, pages 197–202, New York, NY, USA, 2002. ACM.

[30] Stefan Savage, Andy Collins, Eric Hoffman, John Snell, and Thomas Anderson. The end-to-end effects of internet path selection. In *Proceedings of the Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication*, SIGCOMM '99, pages 289–299, New York, NY, USA, 1999. ACM.

[31] Florian Streibelt, Jan Böttger, Nikolaos Chatzis, Georgios Smaragdakis, and Anja Feldmann. Exploring edns-client-subnet adopters in your free time. In *Proceedings of the 2013 Conference on Internet Measurement Conference*, IMC '13, pages 305–312, New York, NY, USA, 2013. ACM.

[32] Visualizing the Internet. http://www.akamai.com/html/technology/ visualizing_akamai.html.

[33] Feng Wang, Zhuoqing Morley Mao, Jia Wang, Lixin Gao, and Randy Bush. A measurement study on the impact of routing events on end-to-end internet path performance. *SIGCOMM Comput. Commun. Rev.*, 36(4):375–386, August 2006.