

# No-Jump-into-Latency in China's Internet! Toward Last-Mile Hop Count Based IP Geo-localization

Chong Xiang  
Shanghai Jiao Tong University  
China

Xinyu Wang  
Shanghai Jiao Tong University  
China

Qingrong Chen  
University of Illinois  
Urbana-Champaign  
USA

Minhui Xue  
The University of Adelaide  
Australia

Zhaoyu Gao  
University of Massachusetts Amherst  
USA

Haojin Zhu\*  
Shanghai Jiao Tong University  
China

Cailian Chen  
Shanghai Jiao Tong University  
China

Qihua Fan  
RTBAsia  
China

## ABSTRACT

Last-mile geo-localization plays an essential role in many location-based services, such as fraud detection and targeted advertising. In this study, we point out that round trip time (RTT) latency shows an extremely weak correlation with physical distance estimation in China's Internet, since a path between a vantage point and a destination can often be circuitous and inflated by queuing and processing delays. To sidestep the latency measurement, we perform a three-tier hop count based IP geo-localization mapping for China's Internet, on the assumption that each provincial router only serves a limited area. The mapping approach begins at the first tier using a single vantage point to fetch large-scale traceroute paths from the server to landmarks and target IPs. At the second tier, we try to find the last common routers along the traceroute paths of targets and landmarks and aggregate their hop count distances. At the third tier, we estimate the physical distances from hop count distances and provincial router radii, and geo-localize the targets to the nearest landmarks. Through large-scale experiments, we show that our approach is both cost-efficient and reliable, and can achieve last-ten-kilometer IP geo-localization for approximately 65% of the total 48874 pingable target IP addresses with a single ping server, and our hop count based approach completely outperforms the RTT based method.

## CCS CONCEPTS

• **Networks** → **Network measurement**; *Network performance analysis*.

\*Corresponding author. Email: zhuhaojin@gmail.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*IWQoS '19, June 24–25, 2019, Phoenix, AZ, USA*

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6778-3/19/06...\$15.00

<https://doi.org/10.1145/3326285.3329077>

## KEYWORDS

IP geo-localization, measurement, China's Internet

### ACM Reference Format:

Chong Xiang, Xinyu Wang, Qingrong Chen, Minhui Xue, Zhaoyu Gao, Haojin Zhu, Cailian Chen, and Qihua Fan. 2019. No-Jump-into-Latency in China's Internet! Toward Last-Mile Hop Count Based IP Geo-localization. In *IEEE/ACM International Symposium on Quality of Service (IWQoS '19)*, June 24–25, 2019, Phoenix, AZ, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3326285.3329077>

## 1 INTRODUCTION

Determining the location of an Internet host is essential for the QoS of many location-based applications, such as local weather forecasting [9], targeted advertising [16, 24, 26], spam filtering, location verification [17, 25], and digital rights management [16]. While coarse-grained IP geo-localization, such as country, state/province, or city-levels, may be sufficient to some of these applications, such as digital rights management, the precise fine-grained IP geo-localization is intensely desired by many other applications, such as online fraud detection and targeted advertising.

Existing round trip time (RTT) based IP geo-localization methods are usually expensive due to the requirement of a large number of active vantage points (ping servers) and ping requests to map RTTs to physical distances. To pinpoint a target IP, RTTs need to be actively measured by dozens of times to obtain the minimal time inflation. The performance is at the cost of ping capacity, which burdens small entities with tight budgets and limited measurement resources. In this paper, we try to ask *if it is possible for an entity with limited ping capacity (half a dozen PlanetLab nodes, RIPE Atlas probes, etc.) to perform IP geo-localization with high accuracy (Q1)*.

Furthermore, accurate IP geo-localization in China's Internet is very challenging due to the extremely weak correlation between RTT latency and physical distance. China's Internet is highly centralized with a handful of top-level ISPs, such as China Telecom, China Unicom, and China Mobile.<sup>1</sup> These ISPs are further hierarchically organized into national backbone networks and regional/provincial networks. The inter-connectivity among ISPs is quite weak due to the ISP barrier, which is unique and different from

<sup>1</sup>We hereafter refer to them as Telecom, Unicom, and Mobile, respectively.

North America and Western Europe [22, 27]. A path between a landmark and a destination can often be circuitous and inflated by queuing and processing delays, and such delays substantially overestimate the actual physical distance and increase geo-localization errors [13, 16, 22, 26]. We, hereby, try to ask further *if it is possible to boost IP geo-localization accuracy in China's Internet (Q2)*.

To answer above two questions, we use a new measure — network hop count — rather than RTT latency, for cost-efficient and reliable physical distance estimation in China's Internet. We first perform a large-scale measurement study on the physical distance's correlation with RTT and hop count, so as to justify the hop count based IP geo-localization approach. As we observe traceroute and hop count are more stable compared with RTT, a single traceroute request from one server would be sufficient for accurate physical distance estimation.

In doing so, we design a three-tier approach which begins at the first tier using a vantage point (*i.e.*, ping server) to obtain large-scale traceroute paths from the server to landmarks and target IPs, forming a network graph. At the second tier, we find the last common router along the traceroute paths of targets and landmarks and identify the shortest paths between nodes (landmarks and target IPs). We then assign an estimated reasonably-bounded coverage radius to each provincial router. Our insight is, within a province or city, the general Ethernet is limited to several kilometers over optical fiber, allowing Ethernet switches in different buildings to be connected [18]. At the third tier, we calculate the estimated physical distances between the targets and landmarks from shortest traceroute paths and the bounded provincial router radii, and geo-localize the targets to the positions of the nearest landmarks.

To evaluate our methodology, we apply our three-tier IP geo-localization approach to 10 provinces of China. For approximately 65% of the total 48874 pingable target IP addresses, we are able to geo-localize them within the error of ten kilometers, thus largely enabling a last-mile IP geo-localization mapping! Finally, the approach developed in this paper completely outperforms state of the art (*i.e.*, the street-level IP geo-localization method [24]) under the same experiment setting.

The contributions of this paper are as follows:

- We use a dataset of the total 244344 pingable IPs (including both landmarks and target IPs) across three major ISPs in China, *i.e.*, Telecom, Unicom, and Mobile. The ground-truth IPs associated with physical geo-locations reside across 10 provinces in China, with 50 m accuracy. We will also make all of the ground-truth data publicly available upon acceptance. At baseline, a large dataset of highly-accurate IPs associated with physical geo-locations is a significant resource and benchmark for future research.
- Through measurements, we show a better correlation between hop count and physical distance in China's Internet and further propose a cost-efficient and reliable hop count based IP geo-localization approach.
- We evaluate our three-tier IP geo-localization approach extensively across 10 provinces in China and show that we can achieve the last-ten-kilometer accuracy for around 65% of the total 48874 pingable target IP addresses. Our approach

outperforms the street-level IP geo-localization method [24] under the same experiment setting.

To the best of our knowledge, this is the first paper to perform the last-mile level IP geo-localization in China's Internet. The geo-localization approach developed in this paper can, at the very least, serve as a starting point for further study of this important area.

## 2 RELATED WORK

### 2.1 RTT-based IP Geo-localization

Various latency-measurement approaches have been proposed [11–13, 16, 19, 21, 22, 24–26]. These approaches primarily use the speed of light to convert RTT measurements to physical distance between a given landmark and a target end host. Representative approaches include GeoPing [22], GeoLim [13], TBG [16], Octant [26], and Spotter [19], each of which adds particular heuristics to adapt to the local Internet measurements where an end host potentially resides. Wang *et al.* [24] proposed to increase landmark density to the point where “street-level” geo-localization is feasible, by enlisting small businesses' web servers as additional landmarks, on the assumption that each server is physically located at the street address of the business. These RTT-based IP geo-localization approaches usually require a large number of ping servers and ping requests. Our hop count based approach, however, only needs one-time traceroute requests from a handful of probing nodes to achieve considerably high geo-localization accuracy.

### 2.2 Hostname Parsing Approaches

Location information are oftentimes encoded by the hostnames associated with router IPs. Using this observation, prior works [10, 15] attempted to map router hostnames to their locations through fixed mappings. Such research generally provides accurate IP geo-localization at the regional level, but can be inaccurate when there is ambiguity on the location-encoded hostname prefix. Furthermore, this technique fails unless location-encoded hostname information is available for the IP address.

### 2.3 IP Geo-localization in China's Internet

There have been three studies focused on China's Internet. Li *et al.* [20] pioneered many latency-measurement based geo-localization techniques and showed that the latency-measurement correlation is weak in China's Internet. Guo *et al.* [14] mined and extracted location information from websites to geolocate IP addresses in China. Tian *et al.* [23] provided a clustering heuristic for improving Chinese commercial GeoIP databases. Our approach differs from the above three studies in that we use the hop count information from traceroute measurements and our geo-localization accuracy is within a last-mile level, substantially outperforming existing province-level geo-localization approaches [14, 20, 23].

### 2.4 Remark

A common thread across these approaches is that they rely on independent or even complementary approaches to geo-localize IP addresses, none of which alone has sufficient accuracy to perform reliable IP geo-localization in China's Internet. Consequently,

researchers are currently in the dark regarding this important domain of the last-mile level of China’s Internet topology. This has important implications for services, security, and performance. In our work, we leverage the insight that within a province or city, the general Ethernet is limited to several kilometers over optical fiber, so there is an opportunity to incorporate hop count to optimize IP geo-localization accuracy as a whole. The next section describes how we validate this insight by large-scale measurements.

### 3 MEASURING CHINA’S INTERNET

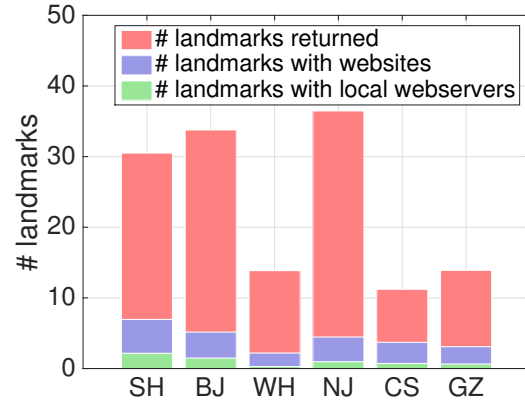
In this section, we first describe the landmark datasets used in this paper. Then we study the general characteristics of China’s Internet, *i.e.*, the correlation between RTT and physical distance as well as the correlation between traceroute hop count and physical distance. Our study shows that the RTT latency has little correlation with its physical distance in China’s Internet, which challenges the fundamental assumption of the latency-based IP geo-localization approach. After performing a large-scale traceroute measurement, we demonstrate that hop count is a more accurate and reliable physical distance estimator than RTT.

#### 3.1 Landmark Dataset

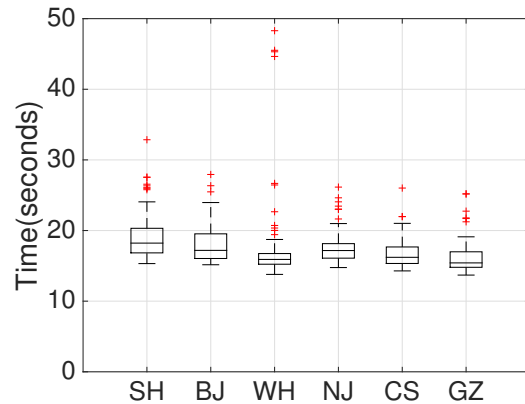
**3.1.1 Attempt to Acquire Landmarks.** The RTT-based approach always exhibits a large estimation error. For example, Octant, the state-of-the-art approach, has a median estimation error of 22 miles [26]. More early works, such as TBG [16] (median estimation error is 67 km) and CBG [13] (median estimation error is 228 km), are less accurate than Octant. Therefore, to achieve street-level geo-localization accuracy, it is considered imperative to increase the landmark density by acquiring closer and more accurate landmarks around the target IP address [24], where Wang *et al.* [24] proposed a map-based mining approach to identify landmarks.

However, it is much harder to collect street-level IP geo-localization landmarks with the map-based mining approach nowadays than a decade ago. The reason is that most of the businesses choose to host their web services remotely (*e.g.*, on clouds or CDNs) rather than hosting them locally. Therefore, the number of landmarks discovered by the map-based mining approach is also expected to be far less than a decade ago. Note that the trend that more and more web services are hosted on clouds and CDNs does not necessarily suggest that the number of web servers hosted locally is shrinking, because it is possible that the numbers of locally and remotely hosted servers have both increased. By reproducing the approach proposed in [24], we illustrate that the number of landmarks discovered by the map-based mining approach is less than two on average in a 20 km-radius geo-fence area for the most populated cities in China. Furthermore, the overhead (in terms of overall response time) and the dollar cost to collect more landmarks are also infeasible for us when using the map-based mining approach in China’s Internet.

We re-implemented the map-based mining approach proposed in [24] and try to find more landmarks in half a dozen top cities in China, such as Shanghai, Beijing, Guangzhou, Wuhan, Changsha, and Nanjing. There is no field containing web information in the geo-query responses with the top two map applications in China (Gaode Map[3] and Baidu Map[1]). We landed with Google Map,



(a) Number of landmark returned by Google nearby queries



(b) The overall time to execute each query

**Figure 1: Results of map-based mining approach [24] for China’s top cities**

which is the only map app that contains the web information for each location. We query the Google Map Place API [5] with a coordinate randomly generated within the range of 20 km of the city center. In doing so, we first use the nearby query to list all the places nearby. The radius of the issued nearby query is 20 km. Next, we use the place detail query to retrieve the web URL and zip code of the place. We then check whether the place satisfies the following two conditions based on the original approach: (1) the zip code of the place must match the zip code of the randomly generated coordinate, and (2) the web server must be in the queried 20 km radius geofence area, *i.e.*, the web server must be local. To geo-localize where the web server is running, we first resolve the IP address through DNS. Then we either query whois database and use the address of the entity that owns the IP prefix as the web server’s location or we directly query a commercial GeoIP database [4] to resolve the IP address to a city level. Different from the original approach that actually issues an HTTP request to figure out whether the web server is running on a CDN, we run the script

on a server on the US East Coast and DNS will automatically redirect us to a US CDN server rather than a China's; thus we could easily rule out the server running on CDNs. Our implementation can be found on Github [7]. We did 100 runs for each city and show the average number of places returned in each run, the average number of places returned with a web URL in each run, and the average number of places returned with a local webserver.

Figure 1(a) shows that the average number of places returned in each nearby query is in a range of 10 (Changsha) to 35 (Nanjing) across different cities. However, the average number of returned landmarks with a local web server is less than two for all the cities except Shanghai. The average number of returned landmarks with a local web server is slightly above two, *i.e.*, 2.16, in this measurement. There are two reasons for the small number of landmarks with a local web server: (1) the total number of landmarks with a web URL (websites) is small (as shown in the purple bar), and (2) the fraction of the places that host the web server locally is also small. In this study, we did not verify whether the found web servers are pingable. If some of the discovered web servers are not pingable, it would further reduce the effectiveness of this approach.

Figure 1(b) shows the time that takes to find the landmarks in this study across 100 runs. It roughly takes 15 – 20 seconds to retrieve the nearby search results and go through the information of all the places. Note, we did not heavily optimize the process to send detail queries. If we send all detail queries in parallel, the time is expected to reduce to 10 seconds. A back-of-envelope calculation reveals that it may take a whole year to retrieve the data for China with this approach (there are 1.8 million zip codes in China [2] and each search takes 15 seconds):  $1.8M \times 15 / 3600 / 24 = 320$  days. In addition to the time cost, we estimate that it may cost us 1.2 million US dollars to collect all the data based on Google's Place API Pricing policy [6], which is also unaffordable for us.

**3.1.2 RTBAsia Landmark Dataset.** Therefore, we directly acquire our landmark dataset from RTBAsia [8], a dominant IP geo-location company in China's market. The whole dataset consists of millions of IPs across China (Figure 2). For our experiments, we choose IPs located in 10 out of 31 China's provinces, which consist of 244344 pingable IP addresses associated with their geographic coordinates in latitudes and longitudes (see the highlighted dots in Figure 2). Based on the dataset description, it is as accurate as GPS geolocation, *i.e.*, within an error of 50 m. The dataset also covers three dominant ISPs in China, *i.e.*, Telecom, Unicom, and Mobile. Table 1 shows the breakdown of landmarks across 10 provinces mostly located in southeastern china.

### 3.2 Network RTT

To investigate the correlation between RTT and geographical distance in China's Internet, we collect the ping latency from one vantage point in Shanghai to all pingable IP addresses in our dataset (over 200K). We ping each target IP address with 5 requests and use the minimum ping latency as the RTT between the target and our vantage point. The scatter plot of RTT versus geographical distance across ten provinces is shown in Figure 3(a). Interestingly, Figure 3(a) shows a linear upper-bound of physical distance with respect to RTT latency for all these provinces.

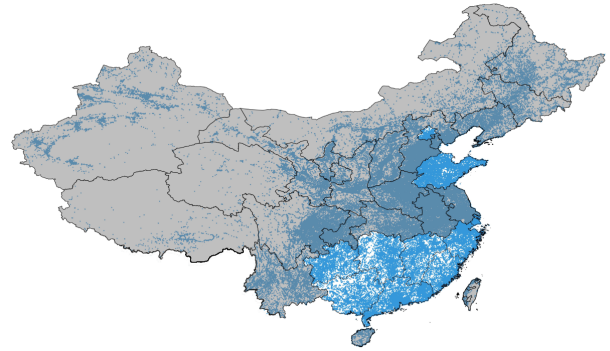


Figure 2: Landmark locations across China

Although RTT latency seems to give us a good approximation ( $1/7 \cdot c$ ) to the *upper bound* of line-of-sight delay (rather than the delay in itself), the geo-localization still heavily relies on multilateration, which needs massive vantage points to do active measurement and costs tremendous complexity to converge to a relatively small overlapping area albeit faithfully [17, 24–26]. Moreover, when we zoom in each province by the order of ten kilometers for physical distance, for example, Shanghai, the data points almost randomly scatter. The correlation coefficient of RTT and physical distance in Shanghai is only 0.098 (see Figure 3(b)). This result aligns with the prior work [20] which suggests that using RTT to approximate physical distance can be highly inaccurate in China's Internet. Empirically, we claim such a weak correlation is potentially due to the high density of the population that incurs unstable communication and computational overheads (such as queuing and processing delays) for routers.

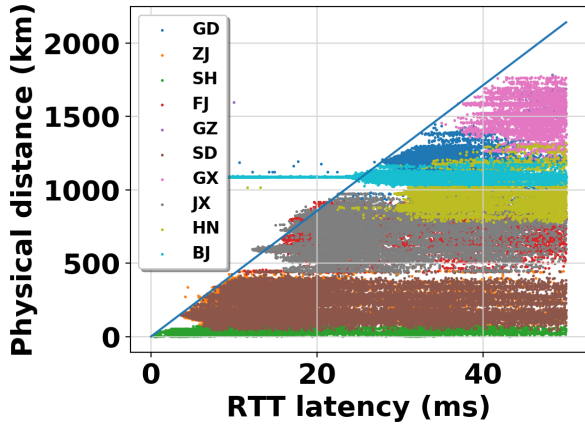
*In summary, we validate that conventional latency based physical distance estimation approaches do not work well for China's Internet. The inaccurately estimated distance hinders further IP geo-localization. Hence, we sidestep latency-based approach and examine alternative measures.*

### 3.3 Network Hop

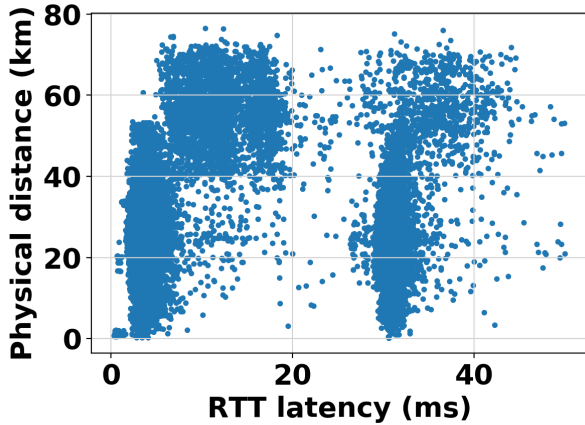
In this part, inspired by the fact that the general Ethernet is limited to several kilometers over optical fiber for efficiency [18], we then take a look at how hop count is correlated to physical distance by conducting a traceroute measurement with the same vantage point

Table 1: The breakdown of landmarks

Province	# Landmarks	Area (km <sup>2</sup> )	GDP per capita (USD)
Beijing (BJ)	34286	16410	18418
Shanghai (SH)	22283	8240	17795
Zhejiang (ZJ)	42421	92057	18418
Fujian (FJ)	17202	124020	11853
Guangdong (GD)	34479	179810	11584
Shandong (SD)	43497	157130	10407
Hunan (HN)	13487	211850	7223
Jiangxi (JX)	15503	166890	6455
Guangxi (GX)	12688	237560	5993
Guizhou (GZ)	8498	176150	5422



(a) Across 10 provinces



(b) In Shanghai

Figure 3: Physical distance versus RTT latency

as used in Section 3.2 and collecting the traceroute results for all IP addresses in our dataset.

**Analysis of the last hop router.** The last hop router is identified as the last visible router interface shown in the traceroute result towards a target IP address. Unlike intermediate hops which could travel across provinces and span hundreds of kilometers, last hops usually physically present in the same block as the target IP, or sometimes even in the same building. Therefore, the IPs sharing the same last hop router tend to be close to each other. To validate this heuristic, we compute the pairwise physical distances among IPs bounded by the same last hop router and plot the CDF of the 291239 distances in Figure 4(a). It shows that the median is about 5 km and 80% of these physical distances are within 14 km. This result implies that, if a target and a landmark happen to share the same last hop router, the estimated geo-localization error tends to be small with a high probability by simply adopting the location

of the landmark for the target IP. Unfortunately, in most cases, we cannot find a landmark that shares the same last hop with the target IP. To address this problem, we extend the concept of the last hop router to the last common router and propose to use hop count for physical distance estimation.

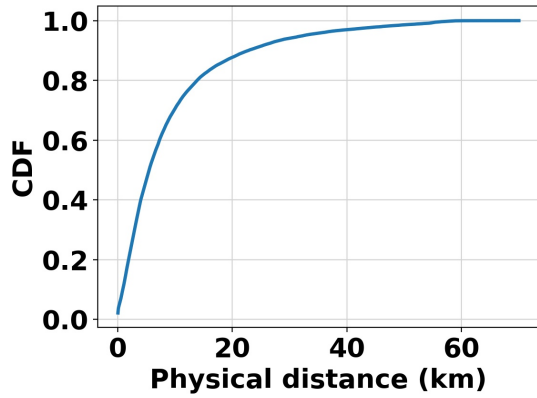
**Analysis of the last common router.** The traceroute paths from one vantage point to different IPs usually share the same route trace until they diverge at the last common routers. Similar to the last hop, the last common router and the target IP are usually within a city, or a province, and serve limited coverage radii. Therefore, we can use the hop count along these routers, which we call provincial routers, as well as their coverage radii to estimate physical distances even when a target IP is not bounded by the same last hop with any landmark. For this measurement, we try to find the last common router along the paths from the vantage point to two different landmarks. We denote  $n_a$  the sum of hop counts from the last common router to the landmark  $a$ . We denote  $n_{hop}(a, b) = n_a + n_b$  as hop count distance between two landmarks  $a$  and  $b$  that are bounded by the same common router. To abuse the notation, we average router coverage radius by  $r = d/n_{hop}$ , once we obtain physical distance  $d$  for a landmark pair. We use the landmark pairs with hop counts less than nine, since large hop counts tend to introduce estimation noise. Figure 4(b) shows the coverage radius of all routers used for analysis. Surprisingly, we find that, for almost 80% of the data, the estimated router radius coverage is within 5 km. Therefore, assigning a fixed value to each router radius can provide a reasonable distance estimation between landmarks, thereby validating the assumption that each router only serves a limited area (*i.e.*, within approximately 5 km from our dataset).

**Analysis of router radius and packet traveling speed.** To further motivate our hop count based geo-localization approach, we take a direct comparison between router radius and packet traveling speed. From the traceroute information, we can find the last common routers and the shortest paths. Similar to hop count distance, the RTT distance for the landmark pair is  $t = t_a + t_b - 2 \cdot t_c$ , where we denote  $t_a$  and  $t_b$  the RTTs from the vantage point to two different end point landmarks  $a$  and  $b$ , and  $t_c$  the RTT from the vantage point to the last common router. We plot the distribution of packet traveling speed  $s = 2 \cdot d/t$  in Figure 4(c). As can be seen, the packet traveling speed is less stable within a certain region. The coefficient of variation (CV) of traveling speed (0.88) is larger than that of router radius (0.75), where CV is the ratio of the standard deviation to the mean and is a measure of relative variability. Also note that the RTT-based distance estimation requires additional ping requests besides traceroute information.

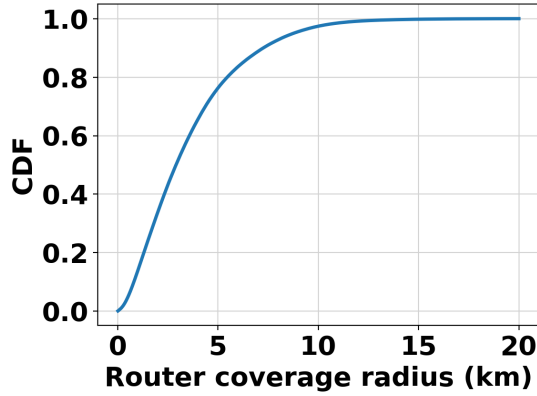
*In summary, we show that hop count with provincial router radius is a more reliable and cost-efficient distance measure than RTT with packet traveling speed.*

## 4 HOP COUNT BASED IP GEO-LOCALIZATION APPROACH

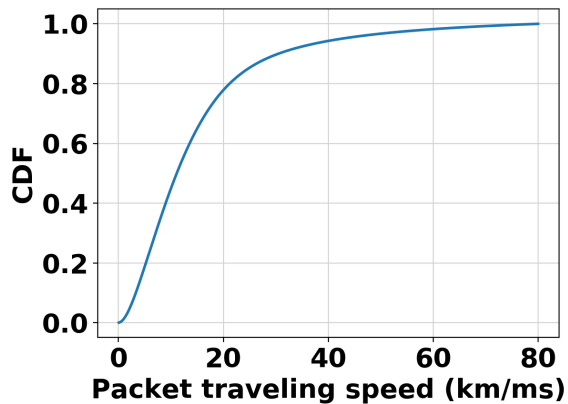
In this section, we incorporate our insights obtained from the above analysis into our final IP geo-localization approach. Note that we only geo-localize IP within a province.



(a) CDF analysis of the physical distances between the landmarks bounded to the same last hop router



(b) CDF analysis of the router coverage radius



(c) CDF analysis of the packet traveling speed

**Figure 4: Traceroute measurements****ALGORITHM 1: Hop-count based IP geo-localization****Input:**

- A ping server  $v_s$ ;
- A set of landmarks  $\mathcal{V}_l$  with their IPs and geo-locations;
- A target IP  $v_t^k$  to be geo-localized.

**Output:**

The geo-location of the target  $v_t^k$ .

# Tier 1.

- 1: obtain the traceroute between ping server  $v_s$  and IP node  $v \in \mathcal{V}_l \cup \{v_t^k\}$ ;

# Tier 2.

- 2: For each IP node pair  $v^i, v^j \in \mathcal{V}_l \cup \{v_t^k\}$ , identify the shortest paths and hop count distances;
- 3: For each landmark  $v_l^i \in \mathcal{V}_l$ , estimate the coverage radii for its nearby provincial routers  $r^i$ ;

# Tier 3.

- 4: For target  $v_t^k$  and each landmark  $v_l^i \in \mathcal{V}_l$ , calculate the physical distance  $\hat{D}_{phy}(v_t^k, v_l^i)$ .

- 5: Find the nearest landmark  $v_l^* = \arg \min_{v_l^i \in \mathcal{T}^k} \hat{D}_{phy}(v_t^k, v_l^i)$ ;

- 6: **return** the geo-location of  $v_l^*$ .

**4.1 Approach Overview**

We model our network topology as an undirected graph,  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where each node  $v \in \mathcal{V}$  is a node representing a router, a landmark, a target IP or the vantage point and each edge  $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$  is a hop connecting two nodes. Our three-tier IP geo-localization approach is listed as follows:

- We obtain traceroute paths from a single vantage point to target IPs and landmarks through active probing, and formulate an network topology of our dataset  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ .
- We find the last common router  $v_c$  along the traceroute paths from the vantage point to the target  $v_t$  and the landmark  $v_l$  and identify the shortest traceroute paths.
- We apply radius estimation for routers around each landmark and subsequently assign an estimated radius to each router in graph  $\mathcal{G}$ .
- We calculate the estimated physical distances between the targets and landmarks through the shortest traceroute paths, and geo-localize the targets to the positions of the nearest landmarks.

The pseudocode is detailed in Algorithm 1.

**4.2 Tier 1. Deriving the Network Topology**

To obtain the topology of the network, we first need to set a vantage point (*i.e.*, where we take measurements). We take traceroute measurements from the vantage point to all the landmarks and targets. Based on each traceroute path, we incrementally add nodes and edges to the topology graph  $\mathcal{G}$ . Note that  $\mathcal{G}$  is only a subgraph of the real network topology since we only use one active probe server with one-time traceroute requests.

**4.3 Tier 2(a). Identifying the Shortest Paths**

After we obtain the network topology  $\mathcal{G}$  from the measurement, we search along the shortest traceroute paths through the last common routers between targets and landmarks. Let  $v_t \in \mathcal{V}_t \subset \mathcal{V}$ ,

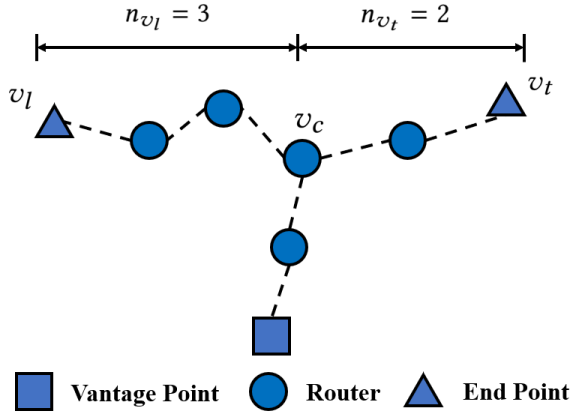


Figure 5: Illustration of the last common router and the shortest path

$v_l \in \mathcal{V}_l \subset \mathcal{V}$ , and  $v_c \in \mathcal{V}_c \subset \mathcal{V}$ , be the target, the landmark, and the last common router, respectively. Then the hop count distance between  $v_l$  and  $v_t$  is  $n_{hop}(v_l, v_t) = n_{v_t} + n_{v_l}$  (see Figure 5). It should be highlighted that since the landmarks and targets are within the same province, the shortest traceroute paths are also at the provincial level, regardless of the locations of a vantage point. Therefore, our assumption on provincial router radius can be used for physical distance estimation.

#### 4.4 Tier 2(b). Estimating Coverage Radius for Routers

In Section 3.3, we show that it is reasonable to assign a limited coverage radius to each provincial router. We further assume that the provincial routers within a small region (e.g., a block) tend to have similar service radii while routers from separated regions may have different radii; hence, we assign a separate coverage radius for routers around each landmark (within a small region) instead of using a universal value. Specifically, for each landmark  $v_l^i \in \mathcal{V}_l$ , we estimate the radius  $r^i$  of nearby routers by using any other landmark  $v_l^j \in \mathcal{V}_l$  such that  $n_{hop}(v_l^i, v_l^j) < T$ , given a certain threshold  $T$ :

$$r^i = \frac{1}{N} \sum_{v_l^j \in \mathcal{L}^i} \frac{D_{phy}(v_l^i, v_l^j)}{n_{hop}(v_l^i, v_l^j)}$$

$$\mathcal{L}^i = \{v_l^j \in \mathcal{V}_l \mid n_{hop}(v_l^i, v_l^j) < T\}, N = \#\mathcal{L}^i, \text{ for any } i \neq j,$$

where  $D_{phy}(v_l^i, v_l^j)$  denotes physical distance between landmarks  $v_l^i$  and  $v_l^j$ . Note that we only use nearby landmarks for radius estimation because the traceroute path between two close IPs is usually less circuitous and better encodes the physical distance information.

#### 4.5 Tier 3. Calculating Physical Distances and Geo-localizing Target IPs

We estimate the physical distance between the target and landmark by minimizing the dot product of the radius  $r^i$  and corresponding hop count distance  $n_{hop}$ . We ultimately geo-localize the target IP to the position of its nearest landmark  $v_l^*$ , as shown below:

$$\hat{D}_{phy}(v_l^k, v_l^i) = r^i \cdot n_{hop}(v_l^k, v_l^i), \text{ given } k,$$

$$v_l^* = \arg \min_{v_l^i \in \mathcal{T}^k} \hat{D}_{phy}(v_l^k, v_l^i),$$

$$\mathcal{T}^k = \{v_l^j \in \mathcal{V}_l \mid n_{hop}(v_l^k, v_l^j) < T\}, \text{ for any } j \neq k.$$

### 5 REAL-WORLD EVALUATION

In this section, we first describe our experiment setup and evaluation metric. Then we take a detailed evaluation of our proposed hop count based geo-localization approach.

#### 5.1 Experiment Setup

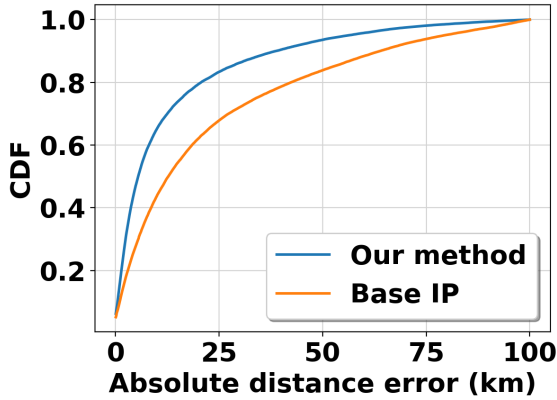
In our experiment, we use one vantage point for geo-localization to simulate the limited ping capacity. We also study the effect by using multiple vantage points in the later part of this section. We split the ground truth dataset into landmarks and targets with a portion of 8 to 2. When geo-localizing target IPs, we only send one single traceroute request to collect essential information. Note that both targets and landmarks are within the same province so that we can apply the assumption on provincial router coverage radius. **Evaluation metric.** We evaluate our experimental results in terms of *absolute distance error* between the predicted IP address location and the actual IP address location based on the ground-truth data.

#### 5.2 Evaluation Results

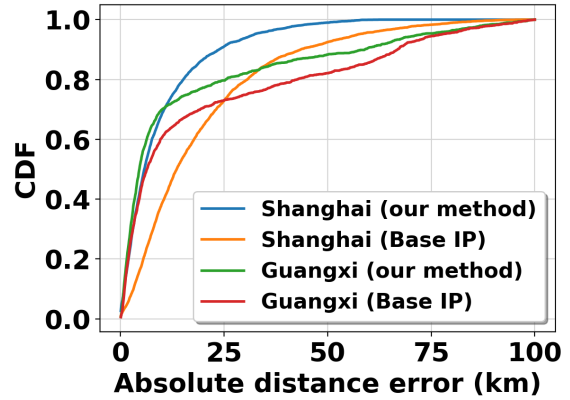
In this subsection, we first evaluate our approach over the whole dataset in comparison with the street-level IP geo-localization approach [24] as a baseline, which we denote as Base-IP. Next, we take a detailed study about how different factors influence the accuracy of our hop count based IP geo-localization.

**In comparison with Base-IP [24].** Since we already have a large number of landmarks with high-confidence ground truth locations, we discard the searching process for available landmarks proposed in the Base-IP approach. The difference between Base-IP and our approach is the measurement used for physical distance estimation: Base-IP uses RTT latency while ours uses hop count. To simulate Base-IP, we send five repeated ping requests from a single server located in Shanghai to targets, landmarks, and last common routers, respectively, and take minimum RTTs. We then geo-localize the target to the position of the landmark with the least RTT distance. For our method, we use the same ping server to collect traceroute information. The comparison result is presented in Figure 6(a). As we can see, our approach completely outperforms Base-IP in terms of estimation error. We can see that for over 60% IPs, the distance error is larger than 20 km for Base-IP while the error is within 10 km for our approach.

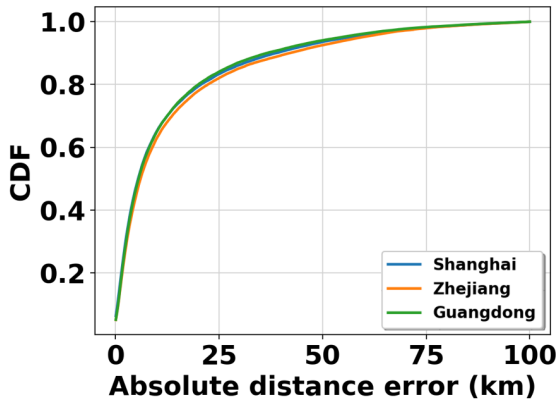
**Geo-localization performance across provinces.** We further look into how the geo-localization performance varies across different provinces. Our single ping server for this experiment is located



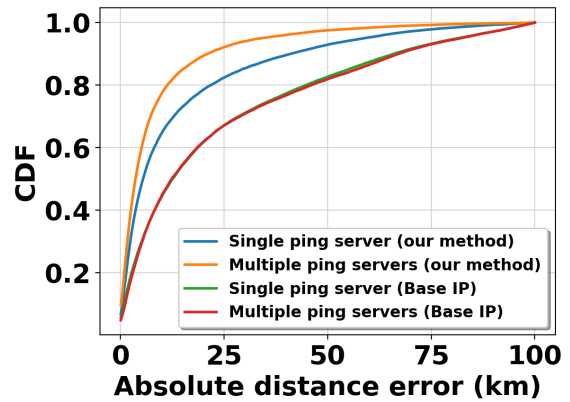
(a) CDF analysis of the absolute distance error in ten provinces from our method and Base-IP



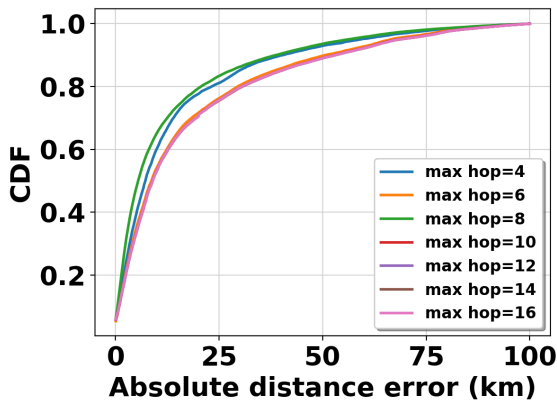
(b) CDF analysis of the absolute distance error in Shanghai and Guangxi



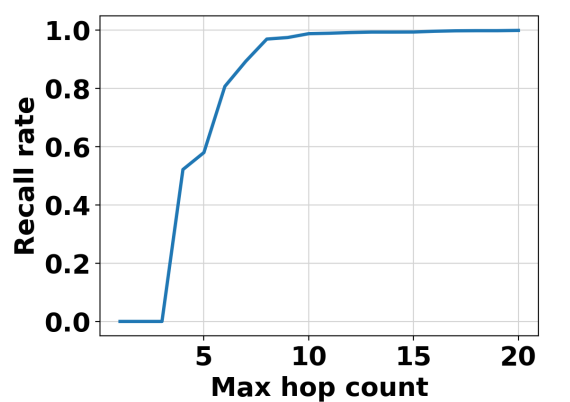
(c) CDF analysis of the influence of the vantage point location



(d) CDF analysis of the incorporation of multiple ping servers



(e) CDF analysis of the absolute distance error across max hop count distance



(f) Recall rate versus max hop count distance

Figure 6: Evaluation of our IP geo-localization approach



**Table 2: The median absolute error across ten provinces**

Province	SH	BJ	GZ	GD	HN	ZJ	GX	FJ	SD	JX
Median absolute error (km)	5.7	1.0	14.1	7.4	6.1	7.2	6.5	4.7	7.6	8.1
Landmark density (per km <sup>2</sup> )	3.5	2.1	0.1	0.2	0.1	0.4	0.1	0.1	0.3	0.1
Population density (per km <sup>2</sup> )	3800	1300	200	600	320	550	190	300	630	270

in Shanghai. The median absolute error for different provinces is presented in Table 2. From Table 2, we see that the median error differs across provinces — the error diminishes especially in well-developed metropolises, such as Shanghai and Beijing, which are administratively equal to provinces, and the error inflates dramatically in less-developed provinces with lower population density.

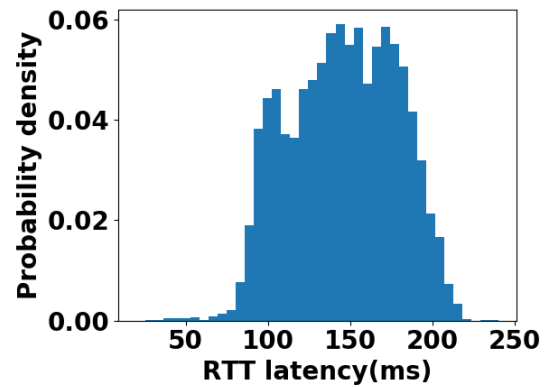
**Influence of landmark density and population density.** To further interpret Table 2, we pick two particular provinces, which are Shanghai (the highest population density in Table 2) and Guangxi (the lowest population density in Table 2) for geo-localization performance comparison. As we can see from Figure 6(b), our approach outperforms Base-IP for both provinces. The subtle inconsistency of performance in different locations is perhaps due to the population and router density. For provinces with high population density, we can imagine that the router density also tends to be high in order to serve more people. The high router density compresses space of provincial routers and results in small router radii; thus enables more precise physical distance estimation. We further tabulate the landmark density and population density in Table 2. We can see that the general trend is that provinces with higher landmark and population density are likely to have smaller estimation errors.

**Influence of the vantage point location.** We use three servers located in Shanghai, Zhejiang, and Guangdong, respectively, to study the influence of the vantage point location on the geo-localization performance. We plot the CDF of estimation errors in Figure 6(c). As we can see from the figure, the estimation errors from different ping servers are similar. This is because traceroute paths and the used provincial routers are all located within the target province separately. The results further validate our reasoning that the geo-localization accuracy is mostly correlated to the property of the target province (*e.g.*, population and landmark density). We have shown that our method is reliable, regardless of the location of a ping server, and this promising method enables more scalable and reliable location-based services.

**Incorporation of multiple ping servers.** Prior RTT-based methods usually use multiple ping servers to find the path or RTT with least inflation. We further check whether multiple ping servers contribute positively to our hop count based approach. To do this, we incorporated the traceroute information from three servers in Shanghai, Zhejiang, and Guangdong to geo-localize the same set of target IPs. Specifically, we took the minimum physical distances estimated by three different ping servers as the final distance estimation and localize the target to its nearest landmarks. We also include the result for Base-IP from multiple servers for comparison. As shown in Figure 6(d), our approach still substantially outperforms the RTT-based approach. Compared to Figure 6(d), we can see a boost in estimation accuracy. We conclude that an entity with



(a) IPs under the same BRAS (in Beijing)



(b) RTT frequency (from Shanghai)

**Figure 7: Broadband remote access server (BRAS)**

large ping capacity can greatly benefit from our hop count based geo-localization approach.

**Influence of max hop count distance.** One crucial point to our approach is to choose a proper hop count threshold  $T$ . A large  $T$  will increase computational overheads and trigger much data noise while a small  $T$  will lead to fewer landmarks used for distance estimation. We analyze the influence of max hop count in Figure 6(e). As we can see, under our experiment setting, the max hop count constraint does not greatly influence the geo-localization, which shows the robustness of our approach. Figure 6(e) also shows that the smaller max hop count tends to have slightly better geo-localization performance. However, we have to mention that a too small max hop count may fail geo-localization, since such landmarks with hop count distance to the target IP within the threshold  $T$  may be unavailable. Hence, a careful trade-off between geo-localization accuracy and recall rate is needed. We set max hop count to 8 in the experiment, since Figure 6(f) shows that max hop count distance equal to 8 has an almost 100% recall rate in geo-localization performance.

### 5.3 Discussion

In this subsection, we discuss the impact of the ubiquitous broadband remote access server (BRAS) on the effectiveness of the IP geo-localization approach. In access networks, BRAS serves as one of the key components that provide important functionalities, such as enforcing operator policies (e.g., traffic shaping, firewalling) and IP Quality of Service (QoS). BRAS is often deployed at the edge of an ISP's core network, routing traffic through broadband remote access devices, such as digital subscriber line access multiplexers (DSLAM). In the real-world deployment, since a BRAS router can cover a certain area, the IPs associated with the same BRAS router will be dynamically assigned regardless of their hosts' physical locations. Furthermore, the data forwarding and propagation within the BRAS can take multiple hops, which are invisible to external observers and thus are treated with a single hop by using traceroute. All these make the correlation between RTT and physical distance quite weak.

To support the argument, Figure 7(a) shows a geographic region of 50 km<sup>2</sup> in Beijing, covered by some of the red dots representing the locations of IPs bounded by the same BRAS (partial data from RTBAsia [8] are shown.). Figure 7(b) is plotted vertically as a histogram with bars representing probability density with respect to the RTT measured by using a vantage point located in Shanghai. From Figure 7(a), we see that IPs bounded by the same BRAS are spatially scattered in a narrow area; while from Figure 7(b), although all the IPs are bounded by the same last hop BRAS, their RTTs physically differ much from [100, 200] ms. Further recall that Figure 3(a) shows a roughly bounded coefficient of  $1/7 \cdot c$ , and the distance estimated from RTT could be vertically fluctuated  $\pm 4285$  km! These are all threats to validity for RTT based IP geo-localization in China's Internet.

By contrast, as the routing paths through the BRAS are usually circuitous and invisible for Internet users due to the black-box property of BRAS, the last hop count to the destination is compressed to the same, making 19 hops stable enough to reach all destination IPs plotted in Figure 7(a). This surreptitiously justifies the validity of our hop count based IP geo-localization approach.

## 6 CONCLUSION

In this paper, we approach the problem of cost-efficient and reliable IP geo-localization in China. We show that physical distance has almost no correlation with RTT latency for China's Internet. We conclude that the ping-measured distance is sensitive in China's Internet shaped by the complex and uneven internal structure of ISPs. We have overcome the problem of the latency-based measurement by proposing a three-tier hop count based IP geo-localization approach. Through our extensive experiments on the dataset covering three dominant ISPs in China, we show that the proposed approach can geo-localize IP addresses within last-mile level with one traceroute request from a single server, and our approach substantially outperforms any previous studies in terms of accuracy and cost-efficiency. We believe the new measure described in this paper is particularly amenable to IP geo-localization in China's Internet, and hope the methodology developed in this paper will help researchers and practitioners understand IP topology in our increasingly dynamic and fluid world.

## ACKNOWLEDGMENTS

This work was supported by National Science Foundation of China, under Grant No. 61672350.

## REFERENCES

- [1] [n.d.]. Baidu Map. <https://map.baidu.com/>.
- [2] [n.d.]. China Zip Code Database. <https://shuju.youbianku.com/cn-postcode>.
- [3] [n.d.]. Gaode Map. <https://gaode.com/>.
- [4] [n.d.]. GeoIP. <https://www.maxmind.com/en/geoip-demo>.
- [5] [n.d.]. Google Map Place API. <https://developers.google.com/places/web-service/search>.
- [6] [n.d.]. Google Places API Usage and Billing. <https://developers.google.com/places/web-service/usage-and-billing>.
- [7] [n.d.]. Re-implementation of map-based mining approach. <https://github.com/ZhaoyuUmass/GeoIP>.
- [8] [n.d.]. RTBAsia. <http://www.rtbasia.com>
- [9] Doxa Chatzopoulos and Marios Kokkodis. 2007. IP geolocation. *Computer Science and Engineering Dept, UC Riverside, Tech. Rep* (2007).
- [10] Michael J Freedman, Mythili Vutukuru, Nick Feamster, and Hari Balakrishnan. 2005. Geographic locality of IP prefixes. In *Proceedings of the 5th ACM SIGCOMM conference on Internet Measurement*. USENIX Association, 13–13.
- [11] Bamba Gueye, Steve Uhlig, and Serge Fdida. 2007. Investigating the Imprecision of IP Block-based Geolocation. In *International Conference on Passive and Active Network Measurement*. Springer, 237–240.
- [12] Bamba Gueye, Steve Uhlig, Artur Ziviani, and Serge Fdida. 2006. Leveraging Buffering Delay Estimation for Geolocation of Internet Hosts. In *International Conference on Research in Networking*. Springer, 319–330.
- [13] Bamba Gueye, Artur Ziviani, Mark Crovella, and Serge Fdida. 2006. Constraint-based Geolocation of Internet Hosts. *IEEE/ACM Transactions On Networking* 14, 6 (2006), 1219–1232.
- [14] Chuanxiong Guo, Yunxin Liu, Wenchao Shen, Helen J Wang, Qing Yu, and Yongguang Zhang. 2009. Mining the Web and the Internet for Accurate IP Address Geolocations. In *INFOCOM 2009, IEEE*. IEEE, 2841–2845.
- [15] Bradley Huffaker, Marina Fomenkov, et al. 2014. DRoP: DNS-based router positioning. *ACM SIGCOMM Computer Communication Review* 44, 3 (2014), 5–13.
- [16] Ethan Katz-Bassett, John P John, Arvind Krishnamurthy, David Wetherall, Thomas Anderson, and Yatin Chawathe. 2006. Towards IP Geolocation using Delay and Topology Measurements. In *Proceedings of the 6th ACM SIGCOMM Conference on Internet Measurement*. ACM, 71–84.
- [17] Mohammad Taha Khan, Joe DeBlasio, Geoffrey M Voelker, Alex C Snoeren, Chris Kanich, and Narseo Vallina-Rodriguez. 2018. An Empirical Analysis of the Commercial VPN Ecosystem. In *Proceedings of the Internet Measurement Conference 2018*. ACM, 443–456.
- [18] James F Kurose and Keith W Ross. 2017. *Computer Networking: A Top-Down Approach*. Pearson Education.
- [19] Sándor Laki, Péter Mátray, Péter Hágá, Tamás Sebők, István Csabai, and Gábor Vattay. 2011. Spotter: A Model based Active Geolocation Service. In *INFOCOM, 2011 Proceedings IEEE*. IEEE, 3173–3181.
- [20] Dan Li, Jiong Chen, Chuanxiong Guo, Yunxin Liu, Jinyu Zhang, Zhili Zhang, and Yongguang Zhang. 2013. IP-Geolocation Mapping for Moderately Connected Internet Regions. *IEEE Transactions on Parallel and Distributed Systems* 24, 2 (2013), 381–391.
- [21] James A Muir and Paul C Van Oorschot. 2009. Internet Geolocation: Evasion and Counterevasion. *Comput. Surveys* 42, 1 (2009), 4.
- [22] Venkata N Padmanabhan and Lakshminarayanan Subramanian. 2001. An Investigation of Geographic Mapping Techniques for Internet Hosts. In *ACM SIGCOMM Computer Communication Review*, Vol. 31. ACM, 173–185.
- [23] Ye Tian, Ratan Dey, Yong Liu, and Keith W Ross. 2013. Topology Mapping and Geolocating for China's Internet. *IEEE Transactions on Parallel and Distributed Systems* 24, 9 (2013), 1908–1917.
- [24] Yong Wang, Daniel Burgener, Marcel Flores, Aleksandar Kuzmanovic, and Cheng Huang. 2011. Towards Street-Level Client-Independent IP Geolocation. In *NSDI*, Vol. 11. 27–27.
- [25] Zachary Weinberg, Shinyoung Cho, Nicolas Christin, Vyas Sekar, and Phillipa Gill. 2018. How to Catch when Proxies Lie: Verifying the Physical Locations of Network Proxies with Active Geolocation. In *Proceedings of the Internet Measurement Conference 2018*. ACM, 203–217.
- [26] Bernard Wong, Ivan Stoyanov, and Emin Gün Sirer. 2007. Octant: A Comprehensive Framework for the Geolocalization of Internet Hosts. In *NSDI*, Vol. 7. 23–23.
- [27] Artur Ziviani, Serge Fdida, José F de Rezende, and Otto Carlos MB Duarte. 2005. Improving the Accuracy of Measurement-based Geographic Location of Internet Hosts. *Computer Networks* 47, 4 (2005), 503–523.