

# Cloudy with a Chance of Short RTTs

## Analyzing Cloud Connectivity in the Internet

The Khang Dang<sup>‡†</sup> Nitinder Mohan<sup>‡†</sup> Lorenzo Corneo<sup>#</sup> Aleksandr Zavodovski<sup>#</sup>  
Jörg Ott<sup>‡</sup> Jussi Kangasharju<sup>b</sup>

<sup>‡</sup>Technical University of Munich <sup>#</sup>Uppsala University <sup>b</sup>University of Helsinki

<sup>†</sup>Equal contribution

### ABSTRACT

Cloud computing has seen continuous growth over the last decade. The recent rise in popularity of next-generation applications brings forth the question: “Can current cloud infrastructure support the low latency requirements of such apps?” Specifically, the interplay of wireless last-mile and investments of cloud operators in setting up direct peering agreements with ISPs globally to current cloud reachability and latency has remained largely unexplored.

This paper investigates the state of end-user to cloud connectivity over wireless media through extensive measurements over six months. We leverage 115,000 wireless probes on the Speedchecker platform and 195 cloud regions from 9 well-established cloud providers. We evaluate the suitability of current cloud infrastructure to meet the needs of emerging applications and highlight various hindering pressure points. We also compare our results to a previous study over RIPE Atlas. Our key findings are: (i) the most impact on latency comes from the geographical distance to the datacenter; (ii) the choice of a measurement platform can significantly influence the results; (iii) wireless last-mile access contributes significantly to the overall latency, almost surpassing the impact of the geographical distance in many cases. We also observe that cloud providers with their own private network backbone and direct peering agreements with serving ISPs offer noticeable improvements in latency, especially in its consistency over longer distances.

### CCS CONCEPTS

• **Networks** → **Public Internet**; **Network measurement**.

### KEYWORDS

Cloud connectivity, Last-mile latency, Peering, Edge computing

#### ACM Reference Format:

The Khang Dang, Nitinder Mohan, Lorenzo Corneo, Aleksandr Zavodovski, Jörg Ott and Jussi Kangasharju. 2021. Cloudy with a Chance of Short RTTs: Analyzing Cloud Connectivity in the Internet. In *ACM Internet Measurement Conference (IMC '21)*, November 2–4, 2021, Virtual Event, USA. ACM, New York, NY, USA, 18 pages. <https://doi.org/10.1145/3487552.3487854>

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).

IMC '21, November 2–4, 2021, Virtual Event, USA

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-9129-0/21/11...\$15.00

<https://doi.org/10.1145/3487552.3487854>

### 1 INTRODUCTION

Cloud computing has become the core enabler for an ever-increasing growth of networked services on the Internet over the past decade [21]. Cloud providers have made significant investments to expand their global footprint, not just by deploying datacenters in new locations [9] but also installing private backbones interconnecting vast geographical regions [8, 29, 90], deploying Point-of-Presence (PoPs) at Internet eXchange Points (IXPs) [2] and colocation facilities [47] closer to their customers [67]. Due to these advancements in the backbone, the cloud infrastructure was able to handle the sudden rise in user traffic as the majority population moved to work-from-home model around the globe in 2020 [31].

Beyond improving cloud computing infrastructure, interest has recently grown in “edge computing”, a paradigm deploying compute servers closer to the users and outside the managed cloud infrastructure, e.g., on ISP premises [28] or in city-owned buildings [53]. The trend of edge computing is primarily driven by a widespread *belief* that the current cloud infrastructure is too sparsely deployed to support the latency requirements of next-generation mission-critical applications [23], such as AR/VR [56], autonomous vehicles [49], etc. However, the cloud infrastructure has improved dramatically since the inception of edge computing in 2009 [72]. Along with advances in the backbone, cloud hypergiants have also invested heavily in installing new datacenters in previously under-provisioned locations [76]. Furthermore, many new small-to-medium-sized cloud providers, such as Vultr, Linode, DigitalOcean, etc., have entered the market and focus their services on specific geographical regions.

However, the growth in the cloud ecosystem has remained largely unnoticed by researchers. This can be primarily attributed to a shortage of impartial studies that investigate the state of cloud reachability and the factors that impact it globally. Few previous works in this space are either out-of-date since they do not capture the recent expansion of cloud infrastructure [48], cover only a limited set of cloud providers [8], or use vantage points that do not consider users in home environments using wireless connectivity [22]. In this paper, we plug this gap in research by providing a well-rounded, comprehensive analysis of cloud connectivity representative of the majority of real Internet users across the globe. Specifically, we make the following key contributions in this paper:

(1) We conduct a large-scale measurement study spanning over *six months* targeting the compute cloud regions of *nine* major cloud providers with a global presence - totalling 195 datacenters deployed in 28 countries (§3.1). We use 115,000 probes in 140 countries from the commercial measurement platform Speedchecker [52] as our vantage points. Speedchecker probes are end-user mobile devices

deployed in thousands of networks across the globe (§3.2). Our vantage point selection allows us to assess cloud connectivity from ASes that are estimated to host 95.6% of the world’s Internet users. We measure user-to-cloud latency (ping) and path (traceroute) over TCP and ICMP, respectively. We find that the geographical location of the datacenter has the most impact on cloud access latency as users in under-provisioned continents (like Africa or South America) get significantly worse performance than North America or Europe. For large parts of Africa and South America, traversing long undersea cables to reach datacenters in neighbouring better-provisioned continents can result in lower overall latency compared to relying on limited in-continent options.

(2) We compare our Speedchecker measurements to the previous reachability experiments conducted over 8000+ RIPE Atlas probes deployed in 184 countries targeting the same cloud regions (§4.2). We find that the Atlas probes achieve significantly lower latency in all continents (except South America, due to skewed probe distribution in countries hosting datacenters) almost consistently. Further investigation reveals the primary contributing factors to be (a) the wired nature of last-mile access of Atlas hardware probes; and (b) often managed (and non-residential) deployment locations of the probes. As a result, we find that the choice of measurement platform significantly affects the measurement results and analyses outcomes as RIPE Atlas may not accurately represent the connectivity of typical Internet users. On the other hand, the results over RIPE Atlas are a good yardstick for estimating cloud reachability for enterprise (non-residential) customers of cloud providers.

(3) As the Speedchecker probes use WiFi or cellular connections to access the Internet, we also isolate the impact of a wireless last-mile on overall cloud access latency (§5). We find that for a large majority of the population, wireless last-mile still acts as the primary bottleneck in user’s path to cloud - taking almost 40-50% of the total median latency globally. Compared to measurements from RIPE Atlas probes using wired connections, wireless can account for 2-3× additional latency. Since future applications will continue to rely on wireless medium irrespective of computing being handled by cloud or edge, the last-mile will make support for latency-critical applications quite problematic. Interestingly, we find that the type of wireless access (WiFi vs. cellular) does not have a significant impact on end-to-end latency as both connection types show similar variations in last-mile.

(4) We identify different types of interconnections that exist between ISPs and cloud providers and quantify the performance differences caused by them (§6). Our client-facing peering analysis reveals that the inbound traffic towards big-3 *hypergiant* cloud providers (Amazon, Microsoft, and Google) avoids the public Internet paths altogether, thanks to direct peering agreements between these providers and the majority of serving ISPs globally. However, our findings show that latency performance benefits of setting up direct peering are limited in developed continents like Europe as public Internet is well-provisioned and offers minimal overhead. On the other hand, in developing regions such as Asia, direct (or private) peering, along with the use of private WAN, results in significant improvement in latency variations – allowing connections to achieve *consistent* latencies even while traversing large geographical distances. As a result, the approach seems to be the best fit in

continents where a cloud provider intends to deliver a consistent quality-of-service to its clients despite limited motivation to deploy new datacenters.

To foster reproducibility, we publish our collected dataset of 3.8M ping and 7+M traceroute measurements at [60] and scripts at [25]. Additionally, readers can find other supporting datasets related to our study at <https://cloudreachability.github.io/>.

## 2 BACKGROUND & RELATED WORK

### 2.1 Cloud Access over the Internet

Significant efforts have been made over the years to understand the connectivity and latencies within the Internet at different levels. Researchers have focused on mapping an accurate representation of the Internet topology at router level [10, 11], AS-level [35, 57], and PoP-level [77]. Based on these works, several studies have shined a light on how recent advancements in cloud expansion - with the rise of IXPs [2, 46] and cloud-owned private WANs [8, 29] - have resulted in the “flattening” of traditionally hierarchical Internet topology [9]. The endeavours to reduce overheads of the transit Internet backbone have also been fuelled by significant competition within new and existing cloud providers, all contending to control the multi-billion-dollar cloud services market [36].

However, despite these advancements assisting cloud infrastructure, efforts to evaluate global cloud access latency have remained fairly limited. Related works on the subject were either conducted before the growth of cloud networks [48] or focused on a single cloud provider [45]. Others have concentrated on either analyzing the impact of private WAN from within the cloud network to client ISP [9] or for providing multi-cloud inter-connectivity [92]. ThousandEyes annual report in 2019 [86] compared latency for five different cloud providers, but only utilized 98 vantage points - all hosted in datacenters. Corneo et al. [22] conducted a global cloud reachability study targeting nine different cloud providers globally (same as this study) but over RIPE Atlas platform [81]. However, RIPE Atlas is known to be influenced by deployment biases as many vantage points are hosted within managed infrastructures, e.g., premises of network service providers, educational institutes, etc. [12, 14, 78] – hence not accurately representing the connectivity of real Internet users globally.

The study by Arnold et al. [8] is most noteworthy to us. The focus of author’s work was to isolate (possible) latency gains when using cloud provider’s private WAN compared to the public Internet. The authors used Speedchecker probes [52] as vantage points (same as this study) and targeted their 22 VM-based endpoints (11 using private WAN and 11 using public Internet) deployed in two hypergiant cloud networks - Amazon and Google. In contrast to [8], the focus of this study is to analyze the reachability and impact of cloud expansion for Internet users across the globe. As a result, we use 195 compute cloud regions operated by *nine* different providers (with a mix of hypergiants and small providers) as endpoints. As such, our study presents a broader overview and gives us an accurate insight into real Internet user metrics when they connect to the cloud for accessing a myriad of networked services.

Since one purpose of our study is to understand if the growth in current cloud infrastructure is feasible for supporting the latency requirements of mission-critical applications for Internet users

globally, we use the following quality-of-experience directives [59] when discussing the latency aspects of this study (§4). **Motion-to-Photon (MTP)** is the delay between user input, and it’s reflecting on the display, which is estimated to be  $\approx 20$  ms. Keeping below this threshold is a strict requirement for immersive applications like AR and VR to avoid motion sickness and dizziness. **Human Perceivable Latency (HPL)** of  $\approx 100$  ms is the threshold when a user starts to experience lags - and is influential for applications such as cloud gaming. **Human Reaction Time (HRT)** denotes the delay difference between a visual stimulus and the associated motor response and is estimated to be  $\approx 250$  ms. The threshold guides the operation of applications involving human-controlled tasks like remote surgery.

## 2.2 Last-Mile Latencies

The “last-mile” is generally regarded as the segment connecting the end-user to its ISP, either via wired or wireless access technology. Previous efforts have focused on studying the characteristics of fixed broadband at a large scale [18, 33, 83, 87]. In [13], the authors investigated last-mile latency from residential probes in Europe and the United States, not including latencies within the home network. Despite the fixed connection to the managed backhaul, previous studies on the topic have revealed last-mile to be the primary congestion and latency bottleneck [33].

While significant efforts have been made to analyze isolated characteristics of wireless technology [75, 84], there is a significant lack of visibility in understanding the impact of wireless on Internet connectivity at a large scale. The reason for this is primarily two-fold. Firstly, studies on this topic rely heavily on specialized monitoring methods, such as deploying custom hardware [70], using third-party datasets [83], designing trusted toolchains [68], or setting up large-scale operational networks [82]. Secondly, there is a lack of publicly-accessible global measurement platforms that allow researchers to conduct network experiments over wireless-equipped probes. For example, vantage points of publicly-accessible large-scale measurement platforms (such as RIPE Atlas, PlanetLab) are majorly deployed in fixed managed networks [14]. As a result, most research relies on setting up custom, short-lived, measurement platform targeting limited geographical regions [65, 89]. However, understanding the impact of the wireless last-mile from a global perspective has become an increasingly important factor in assessing cloud access latency; specially because cloud providers do not have much influence on this part of the connection as it is mainly controlled by regional Internet service providers (ISPs).

In this work, we leverage Speedchecker’s extensive network of wireless-enabled vantage points to plug this gap in this field of research. Being a commercial measurement platform, only limited studies have utilized Speedchecker in the past [7, 8, 34]. To the best of our knowledge, we are the first to estimate the impact and consistency of the WiFi and cellular link while accessing the compute infrastructure of popular cloud providers globally (§5).

## 2.3 Cloud Peering Interconnections

Despite significant investment in infrastructure and private WAN deployment, part of traffic to cloud network is handled by tenant’s serving ISP. To gain more control of their client path into

their managed network, cloud operators leverage several inter-connection approaches to bypass the public Internet altogether – resulting in the Internet “flattening” [9]. One possibility is using cloud exchanges [91], IXPs [15], and colocation facilities [47] which facilitate dedicated peering interconnections within managed third-party datacenter environments. For such interconnections, cloud providers get into a contractual agreement with service ISPs globally by signing a Letter of Authority and Customer Facility Assignment (LOA-CFA) [39]. This allows cloud providers to bypass any inbound tenant traffic originating from those service providers the transit providers altogether [9]. Direct peering enables cloud providers to skip many ASes on a path, thus allowing them to control a significant portion of their tenant connection and achieve high reliability and reduced latency [40].

If the tenant ISP prefers to not directly peer with the cloud provider, they can privately peer at the premises of a third-party transit provider hosting an edge *point-of-presence* (PoP) for that cloud [38]. These entities offer secure and private layer-3 connectivity and can be used by several cloud providers. Such interconnections are commonly referred to as Private Network Interconnects (PNI) [2, 38] and allows cloud providers to circumvent tenant-side regional transit connections, offering a much shorter path to their private WANs. For example, Arnold et al. [9] found networks of hypergiant cloud providers, i.e., Google, Amazon, and Microsoft, to have significantly high reachability with ASes globally, allowing them to be increasingly independent of Tier-1 and Tier-2 ISPs for transporting their traffic in the Internet. Paths without any special peering setup traverse the regular hierarchical public Internet.

Previous works have studied the impact of peering relationships by triggering active measurements from within cloud provider networks [9, 90], colocation facilities [63] and edge PoPs [74]. Arnold et al. [8] investigated the private WAN offerings within Amazon and Google networks. While the focus of [8] was to isolate the impact of cloud private WANs on routing in the Internet, we concentrate on uncovering possible QoS advantages enjoyed by cloud providers with (to without) private WANs from an end-user’s perspective. Our study extends their efforts by making a wider endpoint selection that includes datacenters from *nine* different cloud providers – both with and without a private WAN deployment (see §6).

## 3 MEASUREMENT METHODOLOGY

### 3.1 End-Points Selection

Since the aim of our analysis is to provide a comprehensive overview of cloud reachability and analyze factors affecting it across the globe, our provider selection was influenced by factors such as geographical presence, private WAN deployment, etc. We chose 195 cloud regions operated by *nine* different cloud providers as end-points, namely, Amazon, Google, Microsoft Azure, IBM, Oracle, Alibaba, DigitalOcean, Linode, and Vultr. Table 1 shows the distribution of our endpoints across continents, and Figure 1a shows their deployment density across the globe. For every provider, we filtered cloud regions that support *compute services* (e.g., Amazon ec2, Google compute engine, etc.) and retrieved the hostname of a public VM hosted in that region by CloudHarmony [20]. Some providers, such as Amazon, Google, Microsoft, etc., have built massive *private* WANs to shield tenant traffic from public Internet [76].

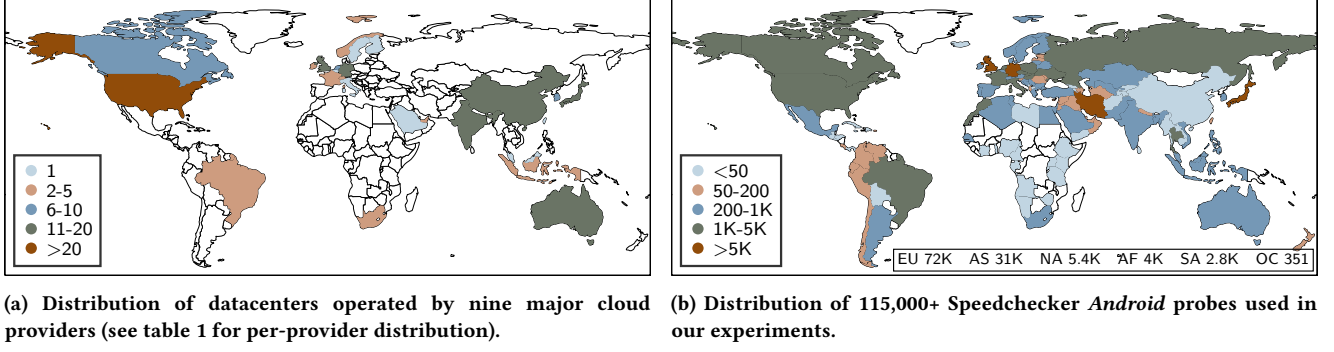


Figure 1: Global coverage of our measurement setup. Cloud datacenters in (a) represent our endpoints, and Speedchecker probes in (b) are the vantage points for our measurements.

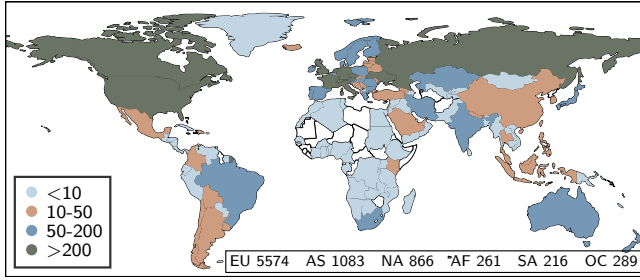


Figure 2: Distribution of 8500+ RIPE Atlas probes used by Corneo et al. [22].

Furthermore, as discussed in §2.3, these providers also deploy edge PoPs and sign contractual agreements with ISPs across the globe to *directly peer* user traffic into their WAN, avoiding transit paths altogether (we analyze the impact of such peering agreements on cloud reachability and path ownership in §6) [8, 19]. On the other hand, small-to-medium-sized providers such as Vultr and Linode, rely heavily on the *public Internet* for transporting their traffic horizontally (between datacenters) and vertically (between users and datacenters) [22]. Providers, such as DigitalOcean and IBM, only establish private backbones in certain geographical regions

Table 1: Global density of cloud provider endpoints, and their backbone network infrastructure.

	Datacenters per continent						Backbone N/W
	EU	NA	SA	AS	AF	OC	
Amazon EC2 (AMZN)	6	6	1	6	1	1	Private
Google (GCP)	6	10	1	8	-	1	Private
Microsoft (MSFT)	14	10	1	15	2	4	Private
Digital Ocean (DO)	4	6	-	1	-	-	Semi
Alibaba (BABA)	2	2	-	16	-	1	Semi
Vultr (VLTR)	4	9	-	1	-	1	Public
Linode (LIN)	2	5	-	3	-	1	Public
Amazon Lightsail (LTSL)	4	4	-	4	-	1	Private
Oracle (ORCL)	4	4	1	7	-	2	Private
IBM (IBM)	6	6	-	1	-	-	Semi
<b>Total</b>	<b>52</b>	<b>62</b>	<b>4</b>	<b>62</b>	<b>3</b>	<b>12</b>	

and use private interconnects to avoid public paths [44]. Table 1 also shows whether a cloud provider has a fully-private (*Private*), private within a continent (*Semi*), or a public Internet-based (*Public*) network backbone. We also include datacenters operated by Alibaba Cloud due to their immense presence in Asia, especially concentrated in China [3].

### 3.2 Vantage Points Selection

**Speedchecker vantage points.** The primary source of data collection in this study comes from our experiments over vantage points (VP) from Speedchecker platform [52]. Speedchecker is a global measurement platform that hosts several hundred thousand softwareized probes in over 170 countries; almost all of which are deployed exclusively in user devices and closely reflect *true* end-user experience. The platform allows researchers to trigger and record active network measurements, e.g., ping, traceroute, HTTP GET, etc., using an API [51]. Probes on Speedchecker are divided into three broad categories based on their operating platform - *router*, *PC*, and *Android*. For our experiments, we only utilized *Android* probes due to two reasons. First, we found that throughout our measurement campaign, Android VPs had the largest share of the total probes on the platform ( $\approx 89\%$ ) - around 470,000 probes deployed globally, of which at least 29,000 were available at any given time. Second, we verified from Speedchecker management team that the majority of Android probes are deployed in real user mobile phones and thus rely on wireless last-mile (WiFi/cellular) to connect to the Internet. Throughout our measurement period, we used upwards of 115,000 Speedchecker probes distributed in over 140 countries worldwide, the country-wise distribution of which is shown in Figure 1b (refer to Appendix A.1 for deployment density based distribution of VPs used in this study).

**RIPE Atlas Dataset:** We correlate and compare our active measurements over Speedchecker to cloud reachability study over RIPE Atlas platform [66] conducted by Corneo et al. [22]. RIPE Atlas is a global Internet measurement platform, driven by network enthusiasts, that includes thousands of small hardware and software probes deployed across the globe. The dataset includes ICMP pings and TCP traceroutes collected from over 8500+ Atlas probes to the same set of cloud regions shown in Table 1. Corneo et al. conducted the study between September 2019 to September 2020, and the 60GB+

dataset includes  $\approx 4\text{M}$  unique probe-to-cloud paths and  $\approx 2.3\text{M}$  ping data points. The dataset is available publicly at [30].

**Speedchecker VP coverage is significantly higher than RIPE Atlas.** Throughout our measurement study, we found that the number of connected Speedchecker probes consistently surpassed those of RIPE Atlas. For example, while RIPE Atlas offers  $\approx 9\text{K}+$  active probes, Speedchecker allows researchers to utilize  $\approx 29\text{K}+$  probes at any given time out of its 115K total available probes. More importantly, the geographical deployment density and availability of Speedchecker is much more comprehensive. Figure 1b and 2 shows the geo-distribution of Speedchecker and RIPE Atlas probes used in our study. While VPs from both platforms are highly concentrated within Europe and North America, Speedchecker’s probe density per geographical distance (a.k.a. `geoDensity`) is almost  $12\times$  in EU and  $6\times$  in NA compared to Atlas. Additionally, unlike Atlas, Speedchecker has at least 200 VPs in almost all countries within these two continents. Germany, Great Britain, Iran, and Japan have the densest VP coverage in Speedchecker, boasting of 5,000+ available probes. The platform’s coverage advantage is especially evident in developing regions of the globe, i.e., countries in SA, Africa, and Asia, where probe `geoDensity` is  $30\text{--}40\times$  higher than RIPE Atlas. However, despite significant availability, VP coverage of both platforms within these regions is relatively sparse, with most probes concentrated in only a few countries (for geographical distribution based on “closeness” of the probes, please refer to Appendix A.1).

From a networking perspective, Speedchecker coverage also overshadows RIPE Atlas quite significantly. Compared to RIPE Atlas, Speedchecker offers  $\approx 14\times$  probes that are hosted in  $\approx 12\text{K}$  ASes (compared to 8K for RIPE Atlas VPs as reported in [22]). To quantify the reach of both platforms for real Internet users, we utilize the user population per ASN dataset from Asia-Pacific Network Information Centre (APNIC) [5]. The dataset estimates the Internet user population coverage of ASes using ad-based measurements. We find that our Speedchecker VPs reside in ASes that cover 95.6% of the Internet user population compared to 69.2% for RIPE Atlas from [22]. It is also important to point out that the platform is also growing at an impressive pace as its population reach has increased by  $\approx 5\%$  since 2019 (as reported by Arnold et al. [8]). Furthermore, unlike the often privileged deployment of RIPE Atlas VPs within managed (mostly wired) network environments (that captures the state of connectivity of non-residential cloud customers) [12, 14, 78], Speedchecker probes are hosted on end-user devices, and the resulting measurements traverse ISP paths reflecting real end-user connectivity towards datacenters. Together, both Speedchecker and RIPE Atlas provides us with a complementary yet most complete picture of global user reachability to cloud till date.

### 3.3 Experiments

We are particularly interested in this work to analyze *two* key aspects of cloud connectivity: (i) state of real user latency to current cloud deployment, especially over wireless paths, and (ii) understanding the impact of cloud provider’s investments in shortening tenant paths to their infrastructure on end-user connectivity. To achieve this, we ran TCP pings and ICMP traceroutes from Speedchecker VPs to cloud region endpoints. Both experiments were conducted in parallel for *six-months*, i.e., from October 2020

to April 2021. Our collected dataset is available at [60], and the reproducibility (+ helper) scripts can be found at [25].

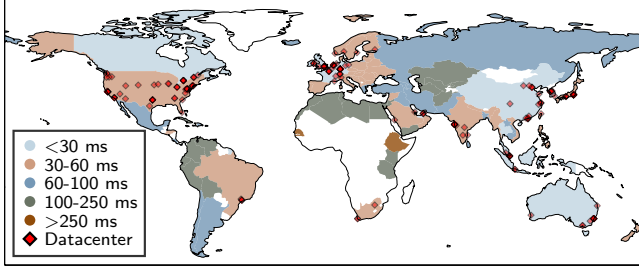
**Statistics and Confidence.** Unless otherwise specified in the rest of the paper, we opt for median round trip latency as our primary metric to assess user connectivity performance across multiple cloud providers. Unlike mean, the median is resilient to outliers that can occur due to bad performing probes, last-mile inconsistencies, and other analysis artifacts [32]. For assessing last-mile access variations (§5), we utilize all recorded measurements.

To make a statistically confident assessment in our analysis, we calculate the minimum measurement sample size required for each country. We define the required confidence interval for the measurement as  $n = \frac{z^2 \times \hat{p}(1-\hat{p})}{\epsilon^2}$ , where  $z$  is the  $z$ -score,  $\hat{p}$  is the population proportion,  $n$  is the target sample size, and  $\epsilon$  is the margin of error. Therefore, for a 95% confidence interval and an error of  $\epsilon = 2\%$ , we collect at least 2400 measurements per country.

**Probe Selection and Experiment Configuration.** Despite its significant reach and probe density, we encountered several challenges while using the Speedchecker platform that influenced our experiment setup. Firstly, we found that the majority of Android probes on the platform were *transient* across days and only became available for use unexpectedly. As a result, we were unable to explicitly trigger experiments over the same set of probes throughout our measurement period and instead had to rely on the platform’s in-built probe selection per geographical region. Secondly, we were provided access to the platform with a limited measurement budget that refreshed at the end of each day. To allow for global coverage with reliable results, we took inspiration from the experimental study of Arnold et al. [9]. Of our total per-day quota, we reserved a few API calls for collecting information about connected VPs, which we triggered at every four-hour interval. We logged all connected probe IDs, their IP addresses, connection type (router, PC, or Android), city-specific geolocation, and ASN - which allowed us to track consistently connected probes on the platform worldwide. We then configured our active network experiments to cycle through every country of each continent with at least 100 probes and targeted all cloud regions within the same continent. For VPs in continents with low datacenter density, e.g., Africa and South America, we also targeted datacenters in neighbouring continents, i.e., Europe and North America (see §4.3). To not overload the platform with our measurement requests, we employed a self-imposed rate limit of *one* measurement request/minute. It took us approximately two weeks to trigger experiments from all countries on the platform, at the end of which we restarted the cycle.

We use TCP ping and ICMP traceroute to estimate end-to-end latencies and distance between users and cloud datacenters, respectively. Overall, we collected over 3.8M ping data points and 7+M unique traceroutes within our study period. The majority of the data points are collected from probes in Europe (around 50%), followed by Asia ( $\approx 20\%$ ) and North America ( $\approx 10\%$ ). Both Africa and South America have almost similar overall contributions in our dataset, with intra-continental taking the larger share over inter-continental measurements ( $\approx 70\text{--}30$  ratio).

We compared the end-to-end latencies from ICMP and TCP measurements over Speedchecker for each `<country, datacenter>` pair and found little-to-no difference between the two protocols.



**Figure 3: Median latency from Speedchecker VPs to the closest datacenter worldwide. Geographical “closeness” is still the primary driving factor for better QoS as countries with in-house cloud deployment achieve much lower latencies than countries without. Africa shows the most uneven performance due to sparse and concentrated datacenter availability favoring southern countries.**

Latencies over TCP tend to be slightly lower than ICMP (within 2% range), which we attribute as possible outliers. The trend departs significantly in RIPE Atlas where ICMP latencies are consistently (and extensively) larger than TCP - especially in Asia, EU, SA, and NA [22]. In both platforms, TCP has lower variance than ICMP, although the median values of the two are comparable in Speedchecker (see Figure 15 in Appendix A.2 for details). Therefore, throughout the rest of the paper, we only use TCP latencies for RIPE Atlas but use both TCP and ICMP interchangeably when analyzing Speedchecker experiments. We solely rely on latencies from traceroutes when investigating the impact of wireless last-mile (§5) and cloud-ISP peering agreements (§6) on cloud access.

**Processing Traceroutes:** We use PyASN [41] to resolve IP-level traceroutes to AS-level paths. For any unresolved router hops (excluding those with private IP addresses) we use Team Cymru IP-to-ASN mapping tool [24]. We further query PeeringDB [1] and enrich our AS-level topology with additional information, such as organization name, location, network type, etc. This phase allows us to accurately identify the serving and transit ISPs on the path responsible for managing VP traffic. Furthermore, we specifically identify the presence of Internet eXchange Points (IXPs) on user paths to cloud using CAIDA IXP dataset [17]. We use GeoIPLookup [37] to geolocate all on-path router hops. However, since such geolocation databases are known to be quite inaccurate [50, 73], we refrain from making any geographical ISP-to-cloud traffic routing assessments in this study and leave that analysis for future work.

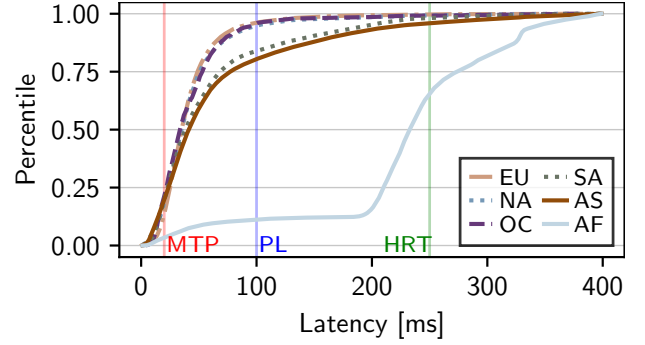
## 4 CLOUD ACCESS LATENCY

We now present our results on cloud access latency from 115,000 Speedchecker wireless VPs and also compare them against the study over RIPE Atlas [22].

### 4.1 Intra-Continental Latency

We start by providing an overview of the access latencies at the global scale shown in Figure 3. The world map represents the median RTT from ping measurements towards the *closest*<sup>1</sup> cloud datacenter (within the same continent) for each country with at least

<sup>1</sup>Datacenter with lowest mean latency over time is estimated to be closest to a probe.



**Figure 4: Distribution of all RTT values by all probes to the nearest datacenter grouped by continent. The vertical lines denote the strict latency thresholds desired by next-generation applications (see §2.1 for details).**

100 Speedchecker probes. The color of the country denotes the latency group (corresponding to latency requirements in §2.1) its median latency lies in. Since only China is able to achieve median RTT below MTP (i.e.s 20 ms), we keep the first latency group as 0-30 ms, all the way until HRT (250 ms). The red diamonds show the approximate locations of cloud regions targeted in this study.

We observe that geographical deployment locations of datacenters have a significant impact on overall cloud performance as countries with in-land datacenters exhibit the best median latency. Among these, China achieves the lowest latency (within MTP bounds), followed by central and northern Europe, North America and South America, India, South Africa, Oceania, and some Asian countries, e.g., Singapore, Indonesia, Thailand, etc. To gain further insight, we plot the distribution of (all) latency measurements recorded by the probe to the nearest datacenter grouped by continents in Figure 4. The results show a very clear trend. Continents well-provisioned with datacenters, i.e., Europe, North America, and Oceania, exhibit very similar latency distributions. Users in these continents can achieve the 100 ms HPL threshold with high probability (as evident by the 90% of the samples from these continents). Keep in mind that the plot includes latency due to the wireless last mile, which is known to be the primary bottleneck in an end-user’s connection [84] (we investigate this in §5). However, achieving MTP in these regions is difficult in the current state of cloud deployments. We investigate the cause of this gap later in this paper.

Countries in continents with sparse datacenter deployments, e.g., South America, Africa, and the Middle East, show significant latency overheads. Within this group, Asia and South America have similar distributions and meet the HPL threshold for roughly 80% of the latency samples, albeit the long tails. We believe the primary contributor to be the significantly lower ratio of available datacenters to total landmass area within these regions. Probes deployed close to a datacenter enjoy quite low latency (see Brazil in Figure 3), which degrades with increasing distances. This phenomenon is more prevalent in South America (Brazil, Argentina, and Chile), Africa (Morocco, Egypt, Algeria), and Asia (India, Pakistan, and Afghanistan) – resulting in significantly long tails in latency distributions within these continents.



The worst performance-hit continent is Africa, where only  $< 10\%$  of latency samples are below the HPL threshold. Closer inspection reveals that these samples belong to VPs near South Africa that also host the only three datacenters endpoints within the continent (correlation between DC deployment in Table 1 and probe availability in A.1). However, 65% of the latency samples satisfy the HRT threshold, the remaining 35% do not. Interestingly, latency distributions from Africa also differ the most out of all continents when compared between Speedchecker and RIPE Atlas platforms [22].

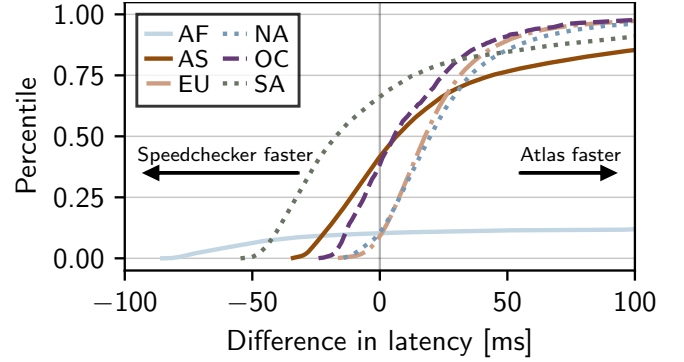
*Takeaway* — Achieving a consistent MTP threshold ( $\leq 20$  ms) is near impossible for Internet users around the globe. 96 out of a total 120 countries can support application requirements governed by HPL threshold ( $< 100$  ms), and all countries, except two in Africa, comply with the HRT threshold (250 ms).

## 4.2 Speedchecker vs. RIPE Atlas

We now compare our Speedchecker measurements to the RIPE Atlas dataset from Corneo et al. [30]. Figure 5 shows the cumulative distribution of differences in latencies recorded from all probes on the two platforms to the nearest datacenter. For clarity, we crop the long tails in this plot. Readers can refer to Figure 4 and [22] for the distributions of all measurements from Speedchecker and RIPE Atlas, respectively. Distributions leaning towards the left indicates faster connectivity over Speedchecker; conversely, distributions towards right implies RIPE Atlas is faster.

The result shows that Atlas probes enjoy slightly better connectivity in Europe and North America compared to Speedchecker. The chasm between the platforms is greatest in Africa, where measurements over RIPE Atlas are significantly faster. The results make sense when one considers the differences in probe deployment location and connectivity type between the two platforms. Within Africa, almost all Atlas probes are situated near the south – physically closer to the in-continent DCs. On the other hand, a large portion of Speedchecker African probes is in the north (see 1b), which takes significantly longer to access the DCs in south. Furthermore, the majority of Atlas probes are hosted by network enthusiasts in managed network environments and connect to the Internet via wired access [12, 14, 78]. On the other hand, Speedchecker probes are exclusively placed on the wireless last-mile since they are deployed as an Android application of end-user’s mobile devices [52]. We investigate the overheads due to wireless last-mile in §5.

As noted, most of the Atlas probes are concentrated close to the datacenter locations (e.g., see Africa in Figure 2), which drives latencies from these countries towards the lower end. Performance differences between the two platforms is similar in Oceania and Asia ( $\sim 60\%$ ) due to similar probe distributions and large geographic distances between VPs and closest datacenters. On the other hand, nearly 70% of the Speedchecker samples from South America are faster than RIPE Atlas. Our explanation for this is as follows. The South American RIPE Atlas dataset contains measurements from probes which are more evenly spread throughout the countries,  $\approx 40\%$  are located in Brazil (where the SA datacenters are). Conversely, more than 80% of the Speedchecker probes are from Brazil, hence delivering lower latency samples than RIPE.



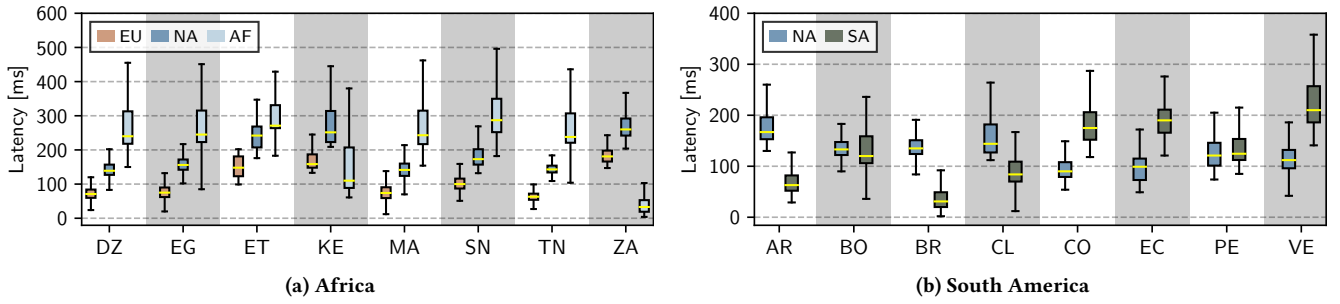
**Figure 5: Latency differences between all measurements from Speedchecker and RIPE Atlas VPs towards the nearest datacenter. The left side denotes samples where Speedchecker is faster, while the right side shows Atlas to be faster. The long tails of the plot have been clipped to maintain legibility. Atlas probes achieve significantly lower latencies than Speedchecker due to their largely wired connectivity.**

To achieve an apples-to-apples comparison, we filtered probes from both platforms with the same  $\langle \text{city}, \text{ASN} \rangle$  of the first hop targeting the same datacenter endpoint. Figure 16 in Appendix A.3 shows the distribution of latency differences between measurements conducted over these probes. Since we did not find enough probe intersection from the same  $\langle \text{city}, \text{ASN} \rangle$  in Africa, South America, and Oceania, we exclude the results from these continents. The result strengthens our arguments above as only a fraction of latency samples in North America are faster in Speedchecker, while for the rest RIPE Atlas achieves significantly lower latencies. While our results highlight the influence of measurement platform on derived conclusions, we are not criticizing the use of RIPE Atlas for cloud measurements. Thanks to its largely wired and managed deployment, RIPE Atlas probes are a good representation of enterprise customers of cloud providers. On the other hand, platforms like Speedchecker provide an accurate reflection of end-user connectivity in home and mobile environments.

*Takeaway* — Measurements over RIPE Atlas generally deliver lower latency compared to Speedchecker. This occurs because the majority of Atlas probes are dedicated hardware devices that connect to the backbone via a wired last-mile. On the other hand, Speedchecker measurements provide a more accurate representation of real Internet user connectivity as the used probes are end-user Android devices connected via a wireless access medium.

## 4.3 Inter-Continental Latency

Previous results revealed that cloud connectivity can be significantly longer in continents with limited datacenter deployment [22]. Therefore we now analyze if the latencies within these regions improve if the users connect to datacenters in neighbouring (better-provisioned) continents. The aim of the analysis is to investigate if shortcomings of sparse geographical datacenter deployments can be overcome by private and faster network backbones. We consider two target regions for this analysis - Africa and South America -



**Figure 6: Cloud access latency from probes in countries within (a) Africa and (b) South America to nearest cloud datacenters within the same and in neighbouring continents.**

since both host only a few datacenters, but are physically close to well-served continents (North America and Europe, respectively).

Figure 6a shows the latency distributions of *all* measurements recorded from African countries to nearest DCs within Africa, Europe, and North America. North African countries like Egypt (EG) and Morocco (MA) have a relatively fast track to Europe due to their physical proximity. Conversely, the path from these countries to datacenters in South Africa is much longer, which manifests as significantly higher access latency. Interestingly, we find that it is faster for these countries to access North American datacenters via undersea cables than in-land ones [16, 85]. Unsurprisingly, probes from South Africa (ZA) have the quickest access to in-land cloud since all three datacenters in the continents used for our measurements are colocated in nearing regions. The most interesting results are shown by Kenya (KE) - a country in the central east side of Africa, which is (almost) equidistant from Europe and South Africa. Here we observe that the lowest median latency is achieved when accessing ZA datacenters, albeit with significant variation. On the other hand, it takes longer to access datacenters in Europe from Kenya, but the distribution appears to be a lot more stable.

Figure 6b plots similar results for VPs in South American countries connecting to DCs in Brazil and NA. Notice that the lowest latencies from the continent are measured from probes in Brazil (BR) and Argentina (AR) when accessing the in-continent datacenters in BR. AR, being the furthest away, clocks the highest latency towards NA datacenters. Results from Bolivia (BO) and Peru (PE) are particularly interesting. Despite being geographically closer to BR than NA, both countries have almost identical latency distributions for the two endpoint regions. This is a likely result of high bandwidth submarine fiber cables connecting both countries directly to North America [85]. Countries located in the north of the continent, e.g., Colombia (CO), Ecuador (EC), and Venezuela (VE), reach NA datacenters quicker than the SA ones. Once again, we verify that cloud access latency is highly influenced by datacenter distance (see BR and AR). However, our analysis also shows that strong networking infrastructure can greatly help in case of local datacenter scarcity (see BO and PE).

**Takeaway** – Networking infrastructure can play an instrumental role in bringing down latencies for regions with sparse datacenter deployments. Remote countries (such as Bolivia, Peru, and Kenya)

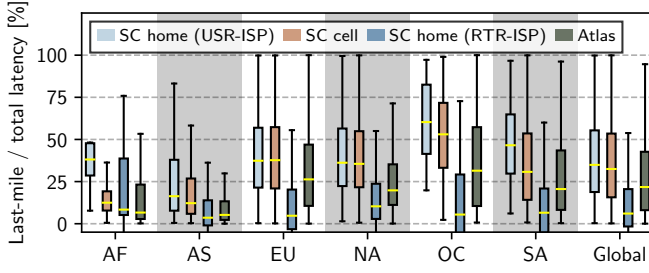
can achieve similar performance connecting to datacenters in-land or within neighbouring continents due to a well-provisioned networking backbone. However, for most countries within SA, Africa, and Asia, physical proximity to datacenters is the driving factor affecting overall access latencies.

## 5 INFLUENCE OF WIRELESS LAST-MILE

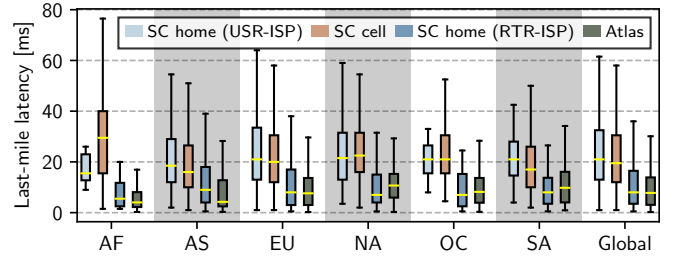
We leverage our traceroute measurements to analyze the impact of the wireless last-mile on cloud access. We infer the last-mile as the link segment between probe IP address and first hop within ISP AS. Since Speedchecker Android probes can either use WiFi or cellular links for connecting to the Internet, we divide them into two broad categories - *home* and *cell*. As the name suggests, *home* VPs are user devices deployed in home networks that use WiFi as wireless connectivity. We identify such VPs through their network paths which traverse a private first-hop (home router) before ingressing the ISP AS. Within this set, we breakup last-mile latency into 1) *wireless inclusive*, i.e., from the probe to ISP (SC home [USR-ISP]) and 2) *wireless exclusive*, i.e., home router to ISP (SC home [RTR-ISP]). The SC *cell* category includes measurements from VPs that have a direct one-hop link to ISP ASN. These probes are, with high likelihood, user devices using cellular wireless medium to access the Internet, and the RTT of the last-mile reflects latency between the device and the cellular tower. Keep in mind that there are several caveats associated with our categorization approach, which may impact the accuracy of our inferences. Firstly, the first hop responding to our traceroutes might not be the basestation itself (home or cellular). As a result, our inferred last-mile may include part of ISP internal network in addition to the wireless media. Similarly, for connections to the Internet via a VPN or carrier-grade NATs (CGN) [71], private addresses will be translated to public IPs; which would directly impact our home-cell probe classification. Secondly, previous research has shown that latency estimates from traceroutes can be inflated due to path inconsistencies, probe processing from underpowered networking devices, and so on. [32, 55, 80]. Such delays are hard to accurately detect post-measurement, and thus may unduly impact our study.

**Last-mile share of user path to cloud.** Figure 7a shows the percentage share of wireless last-mile to the overall cloud access latency for home probes (SC home [USR-ISP]) and cellular probes (SC cell). Firstly, we find that the distribution of the latency share is





(a) Share of wireless last-mile to total cloud access latency for Speedchecker probes.

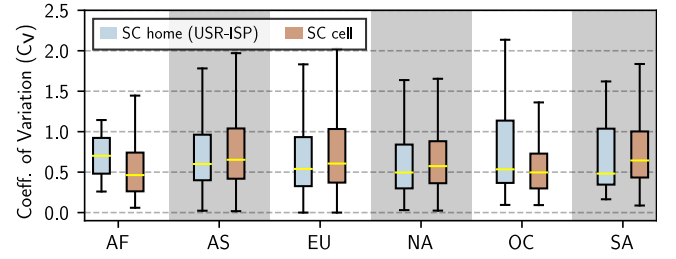


(b) Absolute latency at the wireless last-mile for Speedchecker and RIPE Atlas probes.

**Figure 7: Impact of the wireless last-mile on the cloud access latency grouped by continents.** SC home (USR-ISP) is the latency between the VP and the ISP (via a home router), SC home (RTR-ISP) is the latency between the home router and the ISP, SC cell is the latency between the VP and the first hop of cellular network, and Atlas is last-mile latency of RIPE Atlas probes.

quite similar for both access technologies, irrespective of the probe location. Secondly, the wireless last-mile accounts for a significant share of the total cloud access latency and is higher in continents with more provisioned cloud deployment (i.e., NA and EU.). The result is somewhat expected as not only the overall latency to reach the nearest datacenter is significantly lower within these regions (see §4.1), the latency due to transit is also quite low due to significant deployment of cloud-owned WANs. As a result, the effect of last-mile to overall cloud access latency is more pronounced within these continents. In developing continents, such as Africa and Asia, the percentage share of latency due to last-mile is much smaller as paths to cloud traverse large geographical distances due to relatively sparse datacenter deployment. We also observe that the impact of the last-mile is higher in home probes than cellular probes within developing regions compared to the rest of the globe. Figure 19 in Appendix A.5 illustrates a similar percentage share of last-mile access to end-to-end latency per probe, but only for measurements towards the nearest cloud datacenter. The distribution trend remains fairly unchanged from Figure 7a and further strengthens our inferences drawn above. However, we now find that the latency due to the last-mile is more likely to be the primary bottleneck – as it exceeds the 50% share almost globally.

To understand the behaviour of the last-mile further, we compare the absolute latency at the last-mile for both home and cellular connections in Figure 7b. The plot also compares latency due to the wired part of the home connection (SC home [RTR-ISP]) and last-mile of probes in RIPE Atlas dataset (Atlas) (§4.2). The result indicates that the nature of last-mile (cellular or WiFi) has little influence on the overall cloud access latency across the globe. The latency distribution of path between probe and ISP is similar across continents as the median value hovers around 20–25 ms for both home and cellular connection types. Interestingly, last-mile (irrespective of the access technology) borders close to the MTP threshold ( $< 20$  ms) worldwide. This indicates that even if a compute edge server is deployed directly at the last-mile hop, the latency due to wireless would make MTP almost unachievable for next-generation applications. We also find that the percentage share of the last-mile latency is significantly lower for RIPE Atlas ( $\approx 20\%$ ) (not shown in Figure 7a for brevity) compared to Speedchecker ( $\approx 40\%$ ). Considering that the absolute latency due to last-mile in

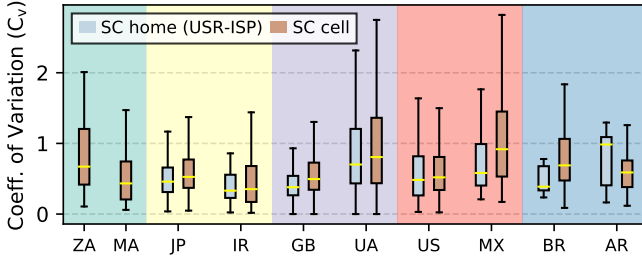


**Figure 8: Coefficient of variation ( $C_v$ ) of wireless last-mile latencies across all measurements per Speedchecker probe grouped by continents.** Higher  $C_v$  indicates higher variation in last-mile latencies of the probe. Similar  $C_v$  for both home and cellular connections hints that the latencies at the last-mile are consistent (and similar) across geographical regions and connectivity types.

Atlas is  $\approx 10$  ms (Figure 7b) provides further validation to the wired nature of Atlas VP access connectivity. As a matter of fact, the latencies from Atlas probes closely resemble the wired part of the Speedchecker home-to-ISP path (SC home [RTR-ISP]).

The large discrepancy between relative and absolute numbers for Africa in Figure 7 can be explained by the probe distribution within the continent. Speedchecker only hosts a limited number of home-based probes in Africa, almost all of which are deployed in the southern part of the continent (coincidentally close to the datacenters). The majority of the remaining (cellular-based) probes ( $\approx 75\%$ ) are concentrated around the north of Africa (in countries like Egypt, Algeria, etc.). This very skewed probe to datacenter deployment distribution results in lower last-mile percentage shares but higher absolute latencies for cellular probes – as their overall path is longer while the last-mile latency stays relatively similar.

**Consistency of the wireless access.** Many networked services (such as content delivery, live video analytics, etc.) care more about consistent than absolute latencies to deliver optimal quality-of-experience to their users [6]. Most of these applications usually employ device buffers to handle long delays, which can react negatively to sudden latency peaks [54]. To understand the feasibility



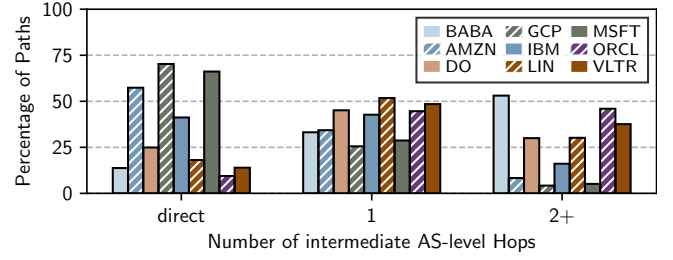
**Figure 9: Coefficient of variation ( $C_v$ ) of wireless last-mile latencies across all measurements recorded from VPs in two representative countries in Africa, Asia, Europe, North America and South America (denoted by different background colors in that particular order). We exclude  $C_v$  for home probes in Africa (ZA & MA) due to insufficient measurement samples from home probes in these two countries.**

of such applications over the cloud, we now analyze the latency variation at the last-mile. Figure 8 plots the coefficient of variation ( $C_v$ ) of the last-mile access (both home and cellular) delays across all measurements per probe across continents. The metric has been used effectively in previous research to identify the quality and stability of the wireless connection [84]. We calculate last-mile  $C_v$  as  $\sigma/\mu$ , where  $\sigma$  is the standard deviation of multiple measurements per probe and  $\mu$  is the mean. In a nutshell,  $C_v$  quantifies the extent of variability with respect to the mean value at the last-mile of each probe - higher values indicating higher variations in latency. We calculate  $C_v$  for all <probe, datacenter> pairs with at least 10 samples. The result shows that both WiFi-based home probes and cellular probes show similar variation across time, with the median  $C_v$  hovering around 0.5. Correlating the results with the absolute latency achieved by home and cellular probes (Figure 7b) confirms that all currently deployed wireless access technologies have similar behaviors and account for a significant portion of the latency to the cloud. Figure 9 sheds more light on our results and shows  $C_v$  of probes in two representative countries in each continent. Even though the plot illustrates subtle stability differences in last-mile delays across different countries, the state, and latencies due to the wireless media is comparable (and significant) throughout the globe. While new technologies like 5G promise to improve the last-mile connectivity, preliminary studies measuring its current deployment in-the-wild show minimal improvements over existing technologies [64, 65]. However, since 5G deployment is still in its nascent stages, its performance is expected to improve in the future [69].

**Takeaway** — Despite significant efforts to improve network connectivity, the last-mile link continues to be the primary bottleneck for cloud providers. As the coveted hop remains out of cloud operator’s influence, latencies due to wireless will make support for latency-critical applications difficult - unless the wireless media improves significantly.

## 6 CLOUD & ISP INTERCONNECTIONS

Our study till now has focused on understanding the state of cloud access across the globe and how user-side of the network impacts



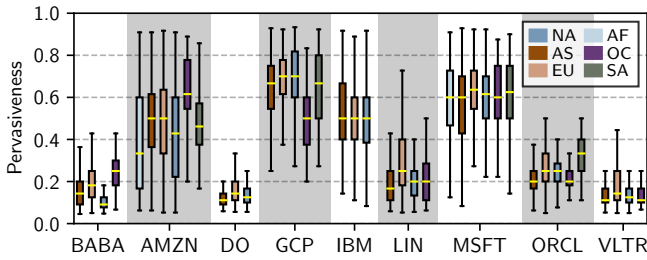
**Figure 10: Different ISP-cloud interconnections globally. *Direct* indicates direct peering between probe ISP and cloud WAN. Paths traversing one AS indicates likely presence of a private transit provider while paths with more on-route ASes possibly showcases cloud providers relying on the public Internet for connectivity.**

overall latencies. As noted in §2.3, cloud providers have made significant investments for shortening the path between tenant ISP and their private WAN by setting up direct and private peering agreements globally. In this section, we isolate the occurrences of such undertakings by the cloud providers in our measurements and analyze their possible impact on reducing user latencies.

### 6.1 Identifying ISP-Cloud Peering.

To accurately identify interconnections between VP ISPs and cloud providers, we remove any unresponsive IP addresses and map the remaining to their respective ASes using the methodology described in §3.3. We identify and tag any IXPs on a path using CAIDA [17] and PeeringDB [1] datasets, and remove them from AS-level topology as they only act as points of traffic exchange. Further, we classify paths where the cloud and probe ISP AS are directly connected neighbours as *direct peering*. Paths where an intermediate AS acts as transit between cloud and VP ISP are tagged as *private peering*. Finally, paths with more than one transit ASes are categorized as *public Internet*. Please note that our peering relationship identification may include several artifacts. Firstly, it is not guaranteed that IXP hops will show up in traceroutes, and therefore we might miss classify routes that traverse via IXPs as direct. Secondly, since we conduct our measurements from probes outside of cloud and ISP networks, our resulting traceroutes may not include router hops within these WANs, thus resulting in mis-identification of interconnections. A more complete approach for accurately identifying cloud peering relationships would be to simultaneously measure from both client-side (like this study) and from within cloud networks (like [8, 74]). Finally, different ISPs globally may have different peering relationships with the cloud providers, and by grouping them together we may miss out on regional-specific routing trends. While we do shed some light on country-specific peering case studies in §6.2 and Appendix A.4, a thorough examination of routing relationships between ISPs and cloud providers is required (similar to [9]), which we plan to undertake in future.

Figure 10 shows the percentage breakup of paths belonging to the three interconnection categories for all cloud providers in our target list. Our results verify the advertised backbone network type of cloud providers shown in Table 1. Majority of the connections



**Figure 11: Degree of pervasiveness of different cloud providers globally. High pervasiveness in Google, Microsoft and Amazon routes shows that majority of routers on end-user’s path to the nearest DC are within ASes owned and operated by the providers themselves.**

bound to networks of the three hypergiants – Amazon, Google, and Microsoft – bypass transit providers altogether, as tenant traffic from serving ISPs directly peers into the provider’s private WANs. For client ISP without direct peering, we find that cloud providers increasingly employ carrier peering via private Tier-1 ISPs (e.g., Telia carrier - AS1299, GTT comm. - AS3257, etc.). Private peering interconnections are used by almost all cloud providers as the peering providers host edge PoPs for multiple operators [2]. Medium-sized cloud providers, e.g., IBM and DigitalOcean, benefit greatly from private peering as their private WANs are still localized, and they can divert their investments into expanding their infrastructure by deploying more datacenters [27]. We find that IBM follows a hybrid interconnection approach as it relies on private peering for shorter paths (concentrated mainly within Europe and North America) but public transit for longer paths (mostly in Asia). Lastly, we find that paths destined to small-sized cloud providers, such as Linode, Vultr, and Oracle, often include two or more on-path ASes, likely hinting routing via the public Internet. Interestingly, Alibaba, despite its massive datacenter and private WAN deployment [4], also uses public Internet paths to interconnect users to its cloud regions. We attribute this behavior to the low availability of Speedchecker probes in China (see Fig. 1b), which does not provide us visibility into Alibaba’s primary operational region. Outside of China, Alibaba operates its datacenters as independent “islands”, only allowing ingress into their WAN via public transit providers.

We also analyze the router-level traceroute data and calculate *pervasiveness* in Figure 11. We define pervasiveness as the ratio between the number of routers owned by the cloud providers to the overall path length to the cloud. High pervasiveness degree hints at most of the end-user route to the cloud to be owned, controlled and operated by the provider themselves – highlighting the reach of their private WAN. We find that pervasiveness of the cloud providers follows a similar trend to the AS-level hop distribution; with Google, Microsoft, and Amazon owning more than 60% of the path in almost every continent. Similarly, providers with two or more ASes only own  $\approx 20\%$  of routers on a path, further strengthening the correctness of our methodology to identify types of ISP-cloud interconnections.

**Takeaway** – Hypergiant cloud providers (Amazon, Google, Microsoft) usually have direct peering with clients’ ISPs ( $> 50\%$ )

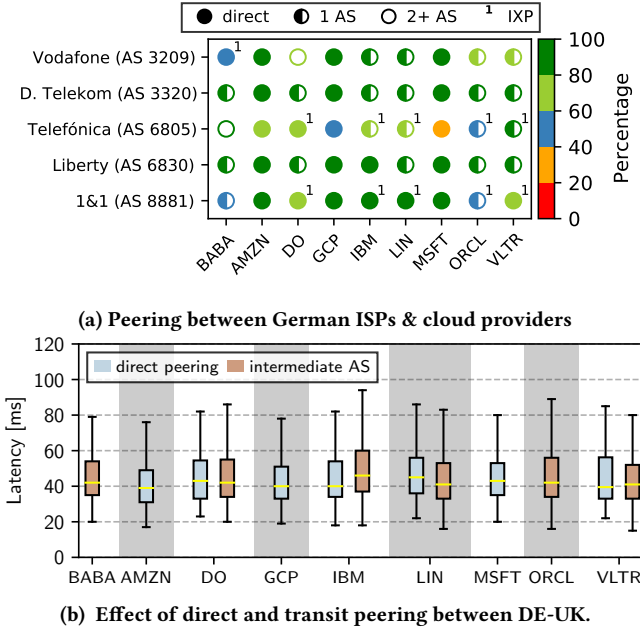
across the globe. When direct peering is not possible, cloud providers prefer to use private interconnects via Tier-1 ISPs like Telia carrier. Smaller providers like Linode, Vultr, Oracle mostly rely on the public Internet for routing their tenant traffic.

## 6.2 Impact of ISP Peering on Latency

We now turn our attention towards understanding the impact of direct ISP-cloud peering interconnections on user cloud access latency. For a thorough analysis, we choose to focus on measurements from Europe (VPs in Germany to DCs in the UK) and Asia (VPs in Japan to DCs in India). Cloud providers have been known to focus their infrastructure investments within Europe and North America, to maximize their profits from the existing user base [67, 92]. However, these continents are already well-provisioned with a reliable Internet backbone and have remained within the limelight of networking innovation for decades [19, 42, 88]. Previous studies have shown that the benefits of private cloud WANs decrease with decreasing geographical distance between user and datacenter [8], and as evident from Table 1, users in both EU and NA have several options for accessing the nearest DC. Comparatively, DC deployment in continents such as Asia, SA, and Africa, is highly scattered – favoring *only* a few select countries. As a result, the impact of using privately managed WANs operated by the cloud providers should be more noticeable within these regions. To keep the analysis comparative across these two continents, we select Germany and Japan as originating countries since both have a dense availability of Speedchecker VPs (see Figure 1b). Similarly, UK and India are selected as endpoints since both have DC deployment from almost all providers in our target list. With this analysis, we aim to understand the continent-specific routing policies set up by cloud providers to transport tenant traffic. We provide more case studies within these continents, specifically Bahrain VPs to India DCs (for Asia) and Ukraine VPs to UK DCs (for Europe) in Appendix A.4 to strengthen the inferences we draw in this section.

Figures 12 and 13 highlight the impact of using different cloud-ISP interconnections in Europe and Asia, respectively. Let’s first focus our attention on Europe. Figure 12a shows the different peering types used by German ISPs<sup>2</sup> while transporting traffic bound to cloud providers. The color denotes the percentage of paths belonging to the majority interconnection type between the ISP and the cloud provider. The result validates our findings in Figure 10. The three hypergiants – Amazon, Google, and Microsoft – exclusively peer directly with almost all serving ISPs in Germany. As a result, the majority of traffic originating from Germany towards DCs of these providers traverses a very “flat” Internet – avoiding even the transit Tier-1 [9]. For other cloud operators, except for traffic originating from Telefonica (AS 6805) towards Alibaba and Vodafone (AS 3209) towards DigitalOcean, almost all German ISPs route their traffic via private interconnection facilities that support the PoP of that provider. We also find that as a medium-sized operator, IBM uses a combination of direct and private interconnects to support tenant traffic. However, it also exchanges traffic at public IXPs more than any of its contemporaries. The trend of setting up direct peering agreements by hypergiants repeats for

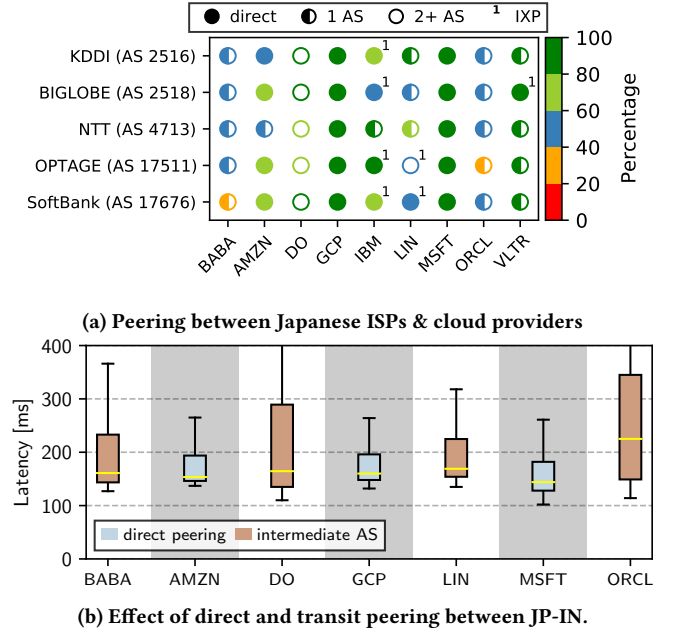
<sup>2</sup>We only show top-5 ISPs ordered by number of recorded measurements.



**Figure 12: Case study of ISP-cloud peering in Europe.** Figure 12a identifies peering interconnections from German ISPs to DCs in the UK. The color denotes the percentage of paths in  $\langle \text{ISP}, \text{cloud} \rangle$  that used the same interconnection type. Figure 12b compares the impact of different interconnection on cloud access latency.

Japanese ISPs as well (Figure 13a) – with the exception of Amazon traffic originating from NTT (AS 4713). Interestingly, we find that DigitalOcean strictly relies on the public Internet for transporting tenant traffic in Asia. We attribute this behaviour to the possible lack of PoP deployment for the cloud provider within the continent. It must be noted here that the transit fabric differs based on the region of the targeted DC. For example, in the case of Japan, we find that for traffic ingress into cloud provider’s WAN (that does not support direct peering) is transited over NTT (AS2914) when both the VP and the DC are co-located within Japan. On the other hand, traffic from VPs in Japan to DCs in India is handled by TATA Comm. (AS6453) for transit. Keep in mind that this behavior is missing in paths over direct peering links where the cloud provider directly handles all of the ingress traffic from tenant’s ISP.

Figure 12b compares the impact of direct peering agreements versus other interconnections on user-to-cloud latency within Europe. For increased confidence, we only show latency values for  $\langle \text{peering type}, \text{cloud provider} \rangle$  pairs – if at least 100 measurements were made for that group. We observe that direct peering between ISPs and cloud operators has minimal effect on cloud access latency between Germany and UK. Latency distributions of the two interconnection categories are quite similar across all cloud providers – indicating that the user latency to the cloud is affected more by geographical distances (§4.3) than routing. We also find minimal latency differences between paths using private peering and public Internet connecting Germany and the UK. The inference



**Figure 13: Case study of ISP-cloud peering in Asia.** Figure 13a identifies peering interconnections from Japanese ISP to DCs in India. The color denotes the percentage of paths in  $\langle \text{ISP}, \text{cloud} \rangle$  that used same interconnection type. Figure 13b compares the impact of different interconnection on cloud access latency.

remains true for paths between Ukrainian ISPs and DCs in UK (see Figure 17 in Appendix A.4), showcasing the after-effects of significant innovations in the backhaul within the continent that has left little margin for overheads due to network management.

The trend, however, differs significantly in paths from Japanese VPs to Indian DCs (see Figure 13b). While the median cloud access latencies are comparable across providers, we find that direct peering significantly reduces the latency variations in the connection (notice shorter box heights in the plot). This result could be a possible outcome of cloud provider’s investment into undersea cables [58, 79]. To investigate this further, we plot the interconnections between VPs in Bahrain to same DCs in India (see Figure 18b in Appendix A.4). Since these two countries are connected by land, routing within these regions is not dependent on common undersea cables. Here we observe a significant latency difference between direct and indirect interconnections, with direct peering achieving consistently shorter latencies.

Similar to the study by Arnold et al. [8], we observe that privately interconnecting paths in Europe can either ingress cloud WAN close to the VP or the server. While different ingress location can affect user path length differently ( $\approx 30\%$  reduction when ingress is close to VP), we found little to no impact of this behaviour on overall access latency. On the other hand, direct peering paths from Japan almost always ingresses cloud WAN within the country itself.



*Takeaway* — Direct peering between cloud and ISP has almost negligible impact on cloud access latency in Europe; showcasing the already well-provisioned public backbone within the region. In Asia, direct peering significantly reduces the long latency tails, which can be especially useful for cloud-backed immersive applications. For in-land interconnections within the continent, direct peering also improves the median latency by a significant margin.

## 7 DISCUSSION

Although our experiments cover a wide range of scenarios, we are inherently limited by the measurement platforms and the nature of network connections. Analysis based on traceroutes, such as ours, are susceptible to inconsistencies from asymmetric forwarding and reverse paths [26, 32]. Likewise, traceroutes only provide us the base network latency between the measurement points. The actual user-observed delay at the application level can be higher due to processing and internal queueing. In this respect, our reported latencies represent the best-case scenario and can be considered as lower bounds on achievable performance. Our final limitation comes from the Speedchecker platform. Our experiments over Speedchecker do not include the last-mile access type (WiFi/cellular) throughout the duration of the measurement. As a result, our analysis inferring the type of wireless access through traceroutes can contain several false positives – including possible switches between WiFi and cellular within the measurement test duration.

In light of the factors affecting cloud access latencies for Internet uses on a global scale, we now discuss the utility of deploying compute edge servers outside of the cloud domain.

**Which networks can live without the edge?** Our results indicate that the latencies to the cloud in regions with dense datacenter deployments are quite stable, regardless of the wired or wireless last-mile connectivity. Developing regions, with poorer connections to cloud datacenters show much more promise for bringing services closer to the users, such as via edge computing [23]. Many of these regions would see considerable improvements in connectivity even with a sparser edge deployment, e.g., via a regional edge or a small datacenter [59, 62]. Developed regions with many cloud datacenters can only see benefits when the deployment of edge is very dense and widely spread, hence their capabilities would not be much improved by edge. Considering that the investments required to set up an edge infrastructure would exceed cloud investments in peering and private WANs, the final preference is likely going to be dictated by the responsible entity, e.g., ISPs would prefer to use the edge while cloud would likely extend their existing WAN.

**Which applications can live without edge?** As our results show, out of the three key latency thresholds (§2.1), cloud is able to satisfy HRT in almost all of the measured cases, and HPL is easily achievable in regions with denser datacenter deployments. However, when it comes to MTP-constrained applications, the picture becomes muddier. As showcased in §5, the absolute wireless last-mile latencies are already on the order of 20+ ms, which makes MTP-constrained applications infeasible, unless all of the processing happens on-board the mobile device. The results hold true for all of the regions and are independent of the density of cloud datacenters. While wireless last-mile latencies can be expected to decrease

(e.g., 5G promising latencies down to 1 ms), it is far from certain whether the reduction would be substantial enough to enable edge deployments since the latency overhead due to transit is minimal (at least in developed regions). But already now, cloud can fully support both HRT and HPL-based applications to an increasing extent. MTP-constrained applications are not really feasible, especially with wireless last-mile, and barring dramatic improvements in wireless technology, are likely to remain infeasible.

While peering agreements between operators help ensure lower latency variations, our results do not indicate that they would markedly reduce the base latencies (especially in regions with a well-provisioned public backhaul). More consistent latencies will aid applications as they make the network more predictable [61]. For example, a video streaming service can make more accurate estimates about the need for buffering and optimize video quality better when the network is stable. Deploying edge servers would help reduce the base latency, which would be beneficial to many applications. However, as discussed above, for many applications, the base latencies are already short enough, so it is not clear what benefits an extensive investment in edge deployment would bring.

It should be stressed that our analysis focuses on network performance issues, and other considerations for edge, such as locality, privacy, etc., are beyond our current scope. As our results show little compelling technical reasons aiding large edge deployments, these non-technical factors should be the focus of future investigations.

## 8 CONCLUSION

Over the past decade, cloud providers have made significant investments for widening their global infrastructure by deploying new datacenters and expanding their private WANs to become much closer to their clients. Furthermore, cloud providers also employ services of colocation facilities and private interconnects hosting their edge PoPs that allows tenant traffic to completely bypass public Internet paths. In this work, we investigate the impact of such cloud advances on overall access latencies for real Internet users globally. Our results show that cloud performance is almost consistent and comparable across providers in continents hosting developed countries due to significant datacenter availability. In developing regions, user latency to the cloud is largely sub-optimal and highly influenced by geographical distance. It is also within these regions, the effects of investments in private WAN and direct ISP peering are more pronounced, as the cloud providers can deliver consistent (and in some areas lower) latencies. Finally, we find that the wireless last-mile is still the primary bottleneck in user cloud access irrespective of the geographical region, hinting at the significant room for research for improving wireless performance.

## ACKNOWLEDGMENTS

We would like to acknowledge the Speedchecker team, especially Janusz Jezowicz, for providing us access to their platform. We also thank our shepherd Alexander Marder and the anonymous IMC reviewers who provided us useful feedback. This work was supported by the Swedish Foundation for Strategic Research with grant number GMT-14-0032 (Future Factories in the Cloud), the Academy of Finland in the BCDC (314167), AIDA (317086), WMD (313477) projects and Celtic project Piccolo (C2019/2-2).



## REFERENCES

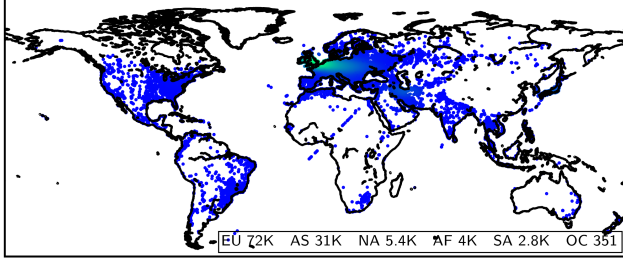
- [1] 2021. PeeringDB. <https://www.peeringdb.com/>.
- [2] Bernhard Ager, Nikolaos Chatzis, Anja Feldmann, Nadi Sarrar, Steve Uhlig, and Walter Willinger. 2012. Anatomy of a large European IXP. In *Proceedings of the ACM SIGCOMM 2012 conference on Applications, technologies, architectures, and protocols for computer communication*. 163–174.
- [3] Alibaba. 2021. Alibaba Cloud's Global Infrastructure. <https://www.alibabacloud.com/global-locations>.
- [4] Alibaba. 2021. Express Connect. <https://www.alibabacloud.com/help/doc-detail/44848.htm>.
- [5] APNIC. 2021. Visible ASNs: Estimated Customer Populations. <https://stats.labs.apnic.net/aspop>. Accessed: 2021-05-24.
- [6] Atakan Aral, Ivona Brandic, Rafael Brundo Uriarte, Rocco De Nicola, and Vincenzo Scoca. 2019. Addressing application latency requirements through edge scheduling. *Journal of Grid Computing* 17, 4 (2019), 677–698.
- [7] Todd Arnold, Matt Calder, Italo Cunha, Arpit Gupta, Harsha V. Madhyastha, Michael Schapira, and Ethan Katz-Bassett. 2019. Beating BGP is Harder than We Thought. In *Proceedings of the 18th ACM Workshop on Hot Topics in Networks* (Princeton, NJ, USA) (*HotNets '19*). Association for Computing Machinery, New York, NY, USA, 9–16. <https://doi.org/10.1145/3365609.3365865>
- [8] Todd Arnold, Ege Gürmerçililer, Georgia Essig, Arpit Gupta, Matt Calder, Vasileios Giotsas, and Ethan Katz-Bassett. 2020. (How Much) Does a Private WAN Improve Cloud Performance?. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*. IEEE, 79–88.
- [9] Todd Arnold, Jia He, Weifan Jiang, Matt Calder, Italo Cunha, Vasileios Giotsas, and Ethan Katz-Bassett. 2020. Cloud Provider Connectivity in the Flat Internet. In *Proceedings of the ACM Internet Measurement Conference*. Association for Computing Machinery, New York, NY, USA.
- [10] Brice Augustin, Xavier Cuvelier, Benjamin Orgogozo, Fabien Viger, Timur Friedman, Matthieu Latapy, Clémence Magnien, and Renata Teixeira. 2006. Avoiding traceroute anomalies with Paris traceroute. In *Proceedings of the 6th ACM SIGCOMM conference on Internet measurement*. 153–158.
- [11] Brice Augustin, Balachander Krishnamurthy, and Walter Willinger. 2009. IXPs: mapped?. In *Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement*. 336–349.
- [12] Vaibhav Bajpai, Steffie Jacob Eravuchira, and Jürgen Schönwälder. 2015. Lessons learned from using the ripe atlas platform for measurement research. *ACM SIGCOMM Computer Communication Review* (2015).
- [13] Vaibhav Bajpai, Steffie Jacob Eravuchira, and Jürgen Schönwälder. 2017. Dissecting Last-Mile Latency Characteristics. *SIGCOMM Comput. Commun. Rev.* (10 2017), 10 pages.
- [14] Vaibhav Bajpai, Steffie Jacob Eravuchira, Jürgen Schönwälder, Robert Kistelevi, and Emile Aben. 2017. Vantage point selection for IPv6 measurements: Benefits and limitations of RIPE Atlas tags. In *2017 IFIP/IEEE Symposium on Integrated Network and Service Management (IM)*. <https://doi.org/10.23919/INM.2017.7987262>
- [15] Timm Böttger, Gianni Antichi, Eder Leao Fernandes, Roberto di Lallo, Marc Bruyere, Steve Uhlig, and Ignacio Castro. 2018. Shaping the Internet: 10 Years of IXP Growth. (10 2018).
- [16] Doug Brake. 2019. Submarine Cables: Critical Infrastructure for Global Communications. *Information Technology & Innovation Foundation: Washington, DC, USA* (2019).
- [17] CAIDA. 2021. CAIDA IXP Dataset. <https://www.caida.org/data/ixps/>. Accessed: 2021-03-23.
- [18] Marshini Chetty, Srikanth Sundaresan, Sachit Muckaden, Nick Feamster, and Enrico Calandro. 2013. Measuring Broadband Performance in South Africa. In *Proceedings of the 4th Annual Symposium on Computing for Development* (Cape Town, South Africa) (*ACM DEV-4 '13*). Association for Computing Machinery, New York, NY, USA, Article 1, 10 pages. <https://doi.org/10.1145/2537052.2537053>
- [19] Yi-Ching Chiu, Brandon Schlinker, Abhishek Balaji Radhakrishnan, Ethan Katz-Bassett, and Ramesh Govindan. 2015. Are We One Hop Away from a Better Internet?. In *Proceedings of the 2015 Internet Measurement Conference* (Tokyo, Japan) (*IMC '15*). Association for Computing Machinery, New York, NY, USA, 7 pages. <https://doi.org/10.1145/2815675.2815719>
- [20] CloudHarmony. 2020. Transparency for the cloud. <https://cloudharmony.com/>.
- [21] CNBC. 2019. Apple is Amazon's biggest customer. <https://www.cnbc.com/2019/04/22/apple-spends-more-than-30-million-on-amazon-web-services-a-month.html>.
- [22] Lorenzo Corneo, Maximilian Eder, Nitinder Mohan, Aleksandr Zavodovski, Suzan Bayhan, Walter Wong, Per Gunningberg, Jussi Kangasharju, and Jörg Ott. 2021. Surrounded by the Clouds. In *Proceedings of The Web Conference 2021* (*WWW '21*). Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3442381.3449854>
- [23] Lorenzo Corneo, Nitinder Mohan, Aleksandr Zavodovski, Walter Wong, Christian Rohner, Per Gunningberg, and Jussi Kangasharju. 2021. (How Much) Can Edge Computing Change Network Latency?. In *2021 IFIP Networking Conference (IFIP Networking)*. IEEE, 1–9.
- [24] Team Cymru. 2021. IP address to ASN mapping service. <https://team-cymru.com/community-services/ip-asn-mapping/>.
- [25] The Khang Dang, Nitinder Mohan, Lorenzo Corneo, Aleksandr Zavodovski, Jörg Ott, and Jussi Kangasharju. 2021. Cloudy with a Chance of Short RTTs – Reproducibility. <https://github.com/tkdang97/Cloudy-with-a-Chance-of-Short-RTTs>
- [26] Amogh Dhamdhare, David D Clark, Alexander Gamero-Garrido, Matthew Luckie, Ricky KP Mok, Gautam Akiwate, Kabir Gogia, Vaibhav Bajpai, Alex C Snoeren, and Kc Claffy. 2018. Inferring persistent interdomain congestion. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*. DigitalOcean. 2017. What's New With DigitalOcean's Network. <https://blog.digitalocean.com/whats-new-with-the-digitalocean-network/>.
- [27] Utsav Drolia, Rolando Martins, Jiaqi Tan, Ankit Chheda, Monil Sanghavi, Rajeev Gandhi, and Priya Narasimhan. 2013. The case for mobile edge-clouds. In *IEEE 10th International Conference on Ubiquitous Intelligence and Computing*. IEEE.
- [28] Ramakrishnan Durairajan, Paul Barford, Joel Sommers, and Walter Willinger. 2015. InterTubes: A study of the US long-haul fiber-optic infrastructure. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*. 565–578.
- [29] Maximilian Eder, Lorenzo Corneo, Nitinder Mohan, Aleksandr Zavodovski, Suzan Bayhan, Walter Wong, Per Gunningberg, Jussi Kangasharju, and Jörg Ott. 2021. Surrounded by the Clouds. <https://doi.org/10.11459/2020mp1593899>
- [30] Anja Feldmann, Oliver Gasser, Franziska Lichtblau, Enric Pujol, Ingmar Posee, Christoph Dietzel, Daniel Wagner, Matthias Wichthuber, Juan Tapiador, Noese Vallina-Rodriguez, Oliver Hohlfeld, and Georgios Smaragdakis. 2020. The Lock-down Effect: Implications of the COVID-19 Pandemic on Internet Traffic. In *ACM Internet Measurement Conference 2020*.
- [31] Romain Fontugne, Cristel Pelsser, Emile Aben, and Randy Bush. 2017. Pinpointing delay and forwarding anomalies using large-scale traceroute measurements. In *Proceedings of the 2017 Internet Measurement Conference*. 15–28.
- [32] Romain Fontugne, Anant Shah, and Kenjiro Cho. 2020. Persistent Last-mile Congestion: Not so Uncommon. In *Proceedings of the ACM Internet Measurement Conference*. 420–427.
- [33] Agustín Formoso, Josiah Chavula, Amreesh Phokeer, Arjuna Sathiseelan, and Gareth Tyson. 2018. Deep diving into africa's inter-country latencies. In *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*. IEEE, 2231–2239.
- [34] Lixin Gao. 2001. On inferring autonomous system relationships in the Internet. *IEEE/ACM Transactions on networking* 9, 6 (2001), 733–745.
- [35] Gartner. 2020. Gartner Forecasts Worldwide Public Cloud Revenue to Grow 6.3% in 2020. <https://www.gartner.com/en/newsroom/press-releases/2020-07-23-gartner-forecasts-worldwide-public-cloud-revenue-to-grow-6point3-percent-in-2020>.
- [36] geoiplookup.net. 2021. GeoIP Lookup XML API. <http://geoiplookup.net/xml-api/>. Accessed: 2021-03-23.
- [37] Google. 2021. Carrier Interconnect overview. <https://cloud.google.com/network-connectivity/docs/carrier-peering>.
- [38] Google. 2021. Cloud Interconnect overview. <https://cloud.google.com/network-connectivity/docs/interconnect/concepts/overview>.
- [39] Google. 2021. Colocation facility locations with low latency. <https://cloud.google.com/network-connectivity/docs/interconnect/concepts/choosing-colocation-facilities-low-latency>.
- [40] Hadi Asghari and Arman Noroozian. 2021. PyASN. <https://pypi.org/project/pyasn/>. Accessed: 2021-03-23.
- [41] Osama Haq, Mamoona Raja, and Fahad R. Dogar. 2017. Measuring and Improving the Reliability of Wide-Area Cloud Paths. In *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*, Rick Barlett, Rick Cummings, Eugene Agichtein, and Evgeniy Gabrilovich (Eds.). ACM. <https://doi.org/10.1145/3038912.3052560>
- [42] Zi Hu, Liang Zhu, Calvin Ardi, Ethan Katz-Bassett, Harsha V. Madhyastha, John Heidemann, and Minlan Yu. 2014. The Need for End-to-End Evaluation of Cloud Availability. In *Passive and Active Measurement*, Michalis Faloutsos and Aleksandar Kuzmanovic (Eds.). Springer International Publishing, Cham.
- [43] IBM. 2021. IBM Cloud Direct Link (2.0). <https://cloud.ibm.com/docs/dl?topic=dl-about>.
- [44] Yuchen Jin, Sundararajan Renganathan, Ganesh Ananthanarayanan, Junchen Jiang, Venkata N. Padmanabhan, Manuel Schroder, Matt Calder, and Arvind Krishnamurthy. 2019. Zooming in on Wide-Area Latencies to a Global Cloud Provider. In *Proceedings of the ACM Special Interest Group on Data Communication* (Beijing, China) (*SIGCOMM '19*). Association for Computing Machinery, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3341302.3342073>
- [45] Rowan Klöti, Bernhard Ager, Vasileios Kotronis, George Nomikos, and Xenofontas Dimitropoulos. 2016. A comparative look into public IXP datasets. *ACM SIGCOMM Computer Communication Review* 46, 1 (2016), 21–29.
- [46] Vasileios Kotronis, George Nomikos, Lefteris Manassakis, Dimitris Mavrommatis, and Xenofontas Dimitropoulos. 2017. Shortcuts through colocation facilities. In *Proceedings of the 2017 Internet Measurement Conference*. 470–476.
- [47] Ang Li, Xiaowei Yang, Srikanth Kandula, and Ming Zhang. 2010. CloudCmp: Comparing Public Cloud Providers. In *Proceedings of the 10th ACM SIGCOMM*

- Conference on Internet Measurement* (Melbourne, Australia) (IMC '10). Association for Computing Machinery, New York, NY, USA, 14 pages. <https://doi.org/10.1145/1879141.1879143>
- [49] Shih-Chieh Lin, Yunqi Zhang, Chang-Hong Hsu, Matt Skach, Md E. Haque, Lingjia Tang, and Jason Mars. 2018. The Architectural Implications of Autonomous Driving: Constraints and Acceleration. *SIGPLAN Not.* (03 2018), 16 pages.
- [50] Ioana Livadariu, Thomas Dreiholz, Anas Saeed Al-Selwi, Haakon Bryhni, Olav Lysne, Steinar Bjørnstad, and Ahmed Elmokashfi. 2020. On the Accuracy of Country-Level IP Geolocation. In *Proceedings of the Applied Networking Research Workshop*. 67–73.
- [51] Speedchecker Ltd. 2021. Probe API Documentation. <https://www.speedcheckerdn.com/probe-api/documentation.html>
- [52] Speedchecker Ltd. 2021. Speedchecker Platform. <https://www.speedchecker.com/>
- [53] Ivan Lujic, Vincenzo De Maio, Klaus Pollhammer, Ivan Bodrozic, Josip Lasic, and Ivona Brandic. 2021. Increasing Traffic Safety with Real-Time Edge Analytics and 5G. In *Proceedings of the 4th International Workshop on Edge Systems, Analytics and Networking*. 19–24.
- [54] Kevin J Ma, Radim Bartos, Swapnil Bhatia, and Raj Nair. 2011. Mobile video delivery with HTTP. *IEEE Communications Magazine* 49, 4 (2011), 166–175.
- [55] Harsha V Madhyastha, Thomas Anderson, Arvind Krishnamurthy, Neil Spring, and Arun Venkataramani. 2006. A structural approach to latency prediction. In *Proceedings of the 6th ACM SIGCOMM conference on Internet measurement*.
- [56] Simone Mangiante, Guenter Klas, Amit Navon, Zhuang GuanHua, Ju Ran, and Marco Dias Silva. 2017. VR is on the edge: How to deliver 360 videos in mobile networks. In *Proceedings of the Workshop on Virtual Reality and Augmented Reality Network*. ACM.
- [57] Zhuoqing Morley Mao, Jennifer Rexford, Jia Wang, and Randy H Katz. 2003. Towards an accurate AS-level traceroute tool. In *Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications*. 365–378.
- [58] Microsoft. 2021. Marea: The future of subsea cables. <https://news.microsoft.com/marea/>.
- [59] Nitinder Mohan, Lorenzo Corneo, Aleksandr Zavodovski, Suzan Bayhan, Walter Wong, and Jussi Kangasharju. 2020. Pruning Edge Research with Latency Shears. In *Proceedings of the 19th ACM Workshop on Hot Topics in Networks*. 182–189.
- [60] Nitinder Mohan, The Khang Dang, Lorenzo Corneo, Aleksandr Zavodovski, Jörg Ott, and Jussi Kangasharju. 2021. Cloudy with a Chance of Short RTTs: Analyzing Cloud Connectivity in the Internet. <https://doi.org/10.1145/2021mp1624200>
- [61] Nitinder Mohan, Tanya Shreedhar, Aleksandr Zavodovski, Otto Waltari, Jussi Kangasharju, and Sanjit K. Kaul. 2018. Redesigning MPTCP for Edge Clouds. In *24th Annual International Conference on Mobile Computing and Networking (MobiCom '18)*. <https://doi.org/10.1145/3241539.3267738>
- [62] Nitinder Mohan, Aleksandr Zavodovski, Pengyuan Zhou, and Jussi Kangasharju. 2018. Anveshak: Placing edge servers in the wild. In *Proceedings of the 2018 Workshop on Mobile Edge Communications*. 7–12.
- [63] R. Motamedi, B. Yeganeh, B. Chandrasekaran, R. Rejaie, B. M. Maggs, and W. Willinger. 2019. On Mapping the Interconnections in Today's Internet. *IEEE/ACM Transactions on Networking* (2019).
- [64] Arvind Narayanan, Eman Ramadan, Jason Carpenter, Qingxu Liu, Yu Liu, Feng Qian, and Zhi-Li Zhang. 2020. A first look at commercial 5G performance on smartphones. In *Proceedings of The Web Conference 2020*. 894–905.
- [65] Arvind Narayanan, Xumiao Zhang, Ruiyang Zhu, Ahmad Hassan, Shuowei Jin, Xiao Zhu, Xiaoxuan Zhang, Denis Rybkin, Zhengxuan Yang, Zhuoqing Morley Mao, Feng Qian, and Zhi-Li Zhang. 2021. A Variegated Look at 5G in the Wild: Performance, Power, and QoE Implications. In *Proceedings of the 2021 ACM SIGCOMM 2021 Conference* (Virtual Event, USA) (SIGCOMM '21). Association for Computing Machinery, New York, NY, USA, 610–625. <https://doi.org/10.1145/3452296.3472923>
- [66] RIPE NCC. 2021. RIPE Atlas. <https://atlas.ripe.net/>.
- [67] Enric Pujol, Ingmar Poese, Johannes Zerwas, Georgios Smaragdakis, and Anja Feldmann. 2019. Steering hyper-giants' traffic at scale. In *Proceedings of the 15th International Conference on Emerging Networking Experiments And Technologies*.
- [68] Mohammad Rajiullah, Andra Lutu, Ali Safari Khatouni, Mah-Rukh Fida, Marco Mellia, Anna Brunstrom, Ozgu Alay, Stefan Alfreðsson, and Vincenzo Mancuso. 2019. Web Experience in Mobile Networks: Lessons from Two Million Page Visits. In *The World Wide Web Conference* (San Francisco, CA, USA) (WWW '19). Association for Computing Machinery, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3308558.3313606>
- [69] Eman Ramadan, Arvind Narayanan, Udhaya Kumar Dayalan, Rostand AK Fezeu, Feng Qian, and Zhi-Li Zhang. 2021. Case for 5G-aware video streaming applications. In *Proceedings of the 1st Workshop on 5G Measurements, Modeling, and Use Cases*. 27–34.
- [70] Shravan Rayanchu, Ashish Patro, and Suman Banerjee. 2012. Catching whales and minnows using wifinet: Deconstructing non-wifi interference using wifi hardware. In *9th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 12)*. 57–70.
- [71] Philipp Richter, Florian Wohlfart, Narseo Vallina-Rodriguez, Mark Allman, Randy Bush, Anja Feldmann, Christian Kreibich, Nicholas Weaver, and Vern Paxson. 2016. A multi-perspective analysis of carrier-grade NAT deployment. In *Proceedings of the 2016 Internet Measurement Conference*. 215–229.
- [72] Mahadev Satyanarayanan, Paramvir Bahl, Ramón Caceres, and Nigel Davies. 2009. The case for vm-based cloudlets in mobile computing. *IEEE pervasive Computing* 8, 4 (2009), 14–23.
- [73] Quirin Scheitle, Oliver Gasser, Patrick Sattler, and Georg Carle. 2017. HLOC: Hints-based geolocation leveraging multiple measurement frameworks. In *2017 Network Traffic Measurement and Analysis Conference (TMA)*. IEEE, 1–9.
- [74] Brandon Schlinder, Italo Cunha, Yi-Ching Chiu, Srikanth Sundaresan, and Ethan Katz-Bassett. 2019. Internet Performance from Facebook's Edge. In *Proceedings of the Internet Measurement Conference* (Amsterdam, Netherlands) (IMC '19). Association for Computing Machinery, New York, NY, USA, 179–194. <https://doi.org/10.1145/3355369.3355567>
- [75] Philipp Schulz, Maximilian Matthe, Henrik Klessig, Meryem Simsek, Gerhard Fettweis, Junaid Ansari, Shehzad Ali Ashraf, Bjoern Almeroth, Jens Voigt, Ines Riedel, Andre Puschmann, Andreas Mitschele-Thiel, Michael Muller, Thomas Elste, and Marcus Windisch. 2017. Latency Critical IoT Applications in 5G: Perspective on the Design of Radio Interface and Network Architecture. *IEEE Communications Magazine* (2017).
- [76] Amazon Web Services. 2019. AWS Global Infrastructure Map. "<https://aws.amazon.com/about-aws/global-infrastructure/>".
- [77] Rob Sherwood, Adam Bender, and Neil Spring. 2008. Discarte: a disjunctive internet cartographer. In *Proceedings of the ACM SIGCOMM 2008 conference on Data communication*. 303–314.
- [78] Rachee Singh, Arun Dunna, and Phillipa Gill. 2018. Characterizing the Deployment and Performance of Multi-CDNs. In *Proceedings of the Internet Measurement Conference 2018* (Boston, MA, USA) (IMC '18). Association for Computing Machinery, New York, NY, USA, 168–174. <https://doi.org/10.1145/3278532.3278548>
- [79] Ben Treynor Sloss. 2018. Expanding our global infrastructure with new regions and subsea cables. *Google Cloud* (2018).
- [80] Neil Spring, Ratul Mahajan, and Thomas Anderson. 2003. The causes of path inflation. In *Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications*. 113–124.
- [81] RN Staff. 2015. RIPE Atlas: A global internet measurement network. *Internet Protocol Journal* (2015).
- [82] Kaixin Sui, Mengyu Zhou, Dapeng Liu, Minghua Ma, Dan Pei, Youjian Zhao, Zimu Li, and Thomas Moscibroda. 2016. Characterizing and improving wifi latency in large-scale operational networks. In *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services*. 347–360.
- [83] Srikanth Sundaresan, Walter de Donato, Nick Feamster, Renata Teixeira, Sam Crawford, and Antonio Pescapè. 2012. Measuring Home Broadband Performance. *Commun. ACM* (Nov. 2012), 10 pages. <https://doi.org/10.1145/2366316.2366337>
- [84] Srikanth Sundaresan, Nick Feamster, and Renata Teixeira. 2016. Home network or access link? locating last-mile downstream throughput bottlenecks. In *International Conference on Passive and Active Network Measurement*. Springer.
- [85] TeleGeography. 2019. Submarine Cable Map. "<https://www.submarinecablemap.com/>".
- [86] ThousandEyes. 2019. *Cloud Performance Benchmark 2019–2020 Edition*. Technical Report. ThousandEyes.
- [87] Martino Trevisan, Danilo Giordano, Idilio Drago, Marco Mellia, and Maurizio Munafo. 2018. Five Years at the Edge: Watching Internet from the ISP Network. In *Proceedings of the 14th International Conference on Emerging Networking Experiments and Technologies* (Heraklion, Greece) (CoNEXT '18). ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3281411.3281433>
- [88] Florian Wohlfart, Nikolaos Chatzis, Caglar Dabanoglu, Georg Carle, and Walter Willinger. 2018. Leveraging interconnections for performance: the serving infrastructure of a large CDN. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*. 206–220.
- [89] Dongzhu Xu, Anfu Zhou, Xinyu Zhang, Guixian Wang, Xi Liu, Congkai An, Yiming Shi, Liang Liu, and Huadong Ma. 2020. Understanding operational 5g: A first measurement study on its coverage, performance and energy consumption. In *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*. 479–494.
- [90] Bahador Yeganeh, Ramakrishnan Durairajan, Reza Rejaie, and Walter Willinger. 2019. How Cloud Traffic Goes Hiding: A Study of Amazon's Peering Fabric. In *Proceedings of the Internet Measurement Conference* (Amsterdam, Netherlands) (IMC '19). 15 pages.
- [91] Bahador Yeganeh, Ramakrishnan Durairajan, Reza Rejaie, and Walter Willinger. 2019. How Cloud Traffic Goes Hiding: A Study of Amazon's Peering Fabric. In *Proceedings of the Internet Measurement Conference*. 202–216.
- [92] Bahador Yeganeh, Ramakrishnan Durairajan, Reza Rejaie, and Walter Willinger. 2020. A First Comparative Characterization of Multi-cloud Connectivity in Today's Internet. In *International Conference on Passive and Active Network Measurement*. Springer, 193–210.

## A APPENDICES

Appendices include supporting material that has not been peer-reviewed.

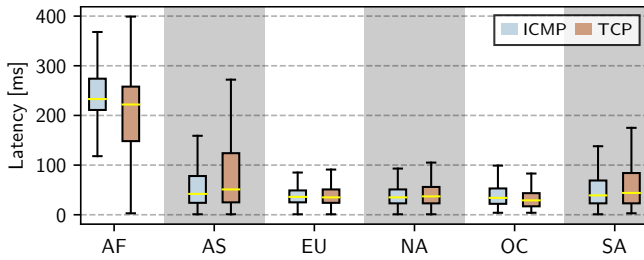
### A.1 Speedchecker VP Infrastructure



**Figure 14: Distribution of 115000+ Speedchecker probes used in this study grouped by their geographical “closeness”.**

Figure 14 showcases the geographical distribution of Speedchecker probes used in this study grouped based on their “closeness” density. Denser deployment of probes is denoted by greener hues. The illustration provides further granularity to the Speedchecker probe deployment shown in Figure 1b. The most noteworthy is the scattered availability of probes in both north and south of Africa – which drives up latencies towards in-continent datacenter deployment within the continent (see Figure 4). This deployment trend differs significantly from RIPE Atlas [22], where both the probes and the targeted datacenters are within close geographical proximity. The figure highlights how geographical deployment density and availability of vantage points belonging to different measurement platforms can affect the outcome of the resulting analysis.

### A.2 ICMP vs. TCP Probe-to-Cloud Latencies

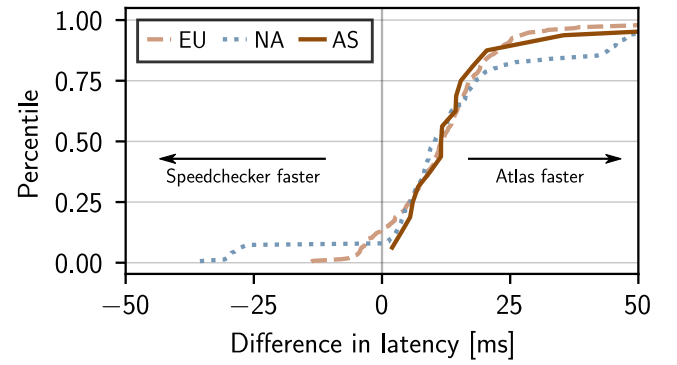


**Figure 15: Difference between end-to-end latencies over ICMP and TCP in Speedchecker grouped by continents.**

Figure 15 plots the end-to-end latencies recorded over ICMP traceroute and TCP pings over Speedchecker VPs for each <country, datacenter> pairing grouped by continents. For regions with dense and highly managed network backhaul (i.e., Europe, North America and Oceania), our result shows little-to-no difference between the latencies over two protocols. In Asia and South America, TCP and ICMP latencies show minor differences (especially at 75th percentiles), however, the median latencies over both protocols are

still comparable. The chasm between the two is largest in Africa, where latencies over TCP tend to be lower compared to ICMP. Here, we must point out that latency differences between the two protocols can be a likely side-effect of the measurement tool itself. Previous research has shown that traceroute is susceptible to inconsistent latency inflations due to path inconsistencies [32, 55, 80]. On the other hand, measurements over ICMP can be affected by load balancers/firewalls in cloud WANs which can route packets over longer paths, put them in lower priority queues, or drop them altogether [43]. In this regard, measurements over TCP are guaranteed to be end-to-end and provide a close estimate of connection latencies encountered by real applications operating in the cloud.

### A.3 Speedchecker vs. RIPE Atlas – <city, AS>



**Figure 16: Latency differences between measurements from Speedchecker and RIPE Atlas in same <city, ASN> towards the nearest DC. The left side denotes samples where Speedchecker is faster, while the right side shows Atlas to be faster.**

To further investigate the impact of different measurement platforms on global cloud accessibility and reachability (see §4.2), we plot the cumulative differences in latencies from probes located in the same city with the first hop within the same ASN. Within this analysis we filter measurements over probe that are handled by the same serving ISP in similar locations – thereby providing an apples-to-apples comparison between the two platforms. Figure 16 shows our results. Since we were unable to find enough probe intersections across the two platforms in Africa, South America and Oceania (largely due to sparse probe availability in RIPE Atlas in these regions), we exclude the results from these continents.

The result strengthens our analysis in §4.2, highlighting the significant connectivity and deployment differences between the two protocols. Only a fraction of Speedchecker probes in North America and Europe achieve better latencies than RIPE Atlas, but for the large majority RIPE Atlas is significantly faster than Speedchecker. In Asia, Atlas is *always* faster – hinting at the impact of wired vs. wireless access differences on overall latency. We plan to conduct measurements over Speedchecker wired probes in future to thoroughly investigate the affect of deployment (managed vs. home) on end-to-end cloud latency.

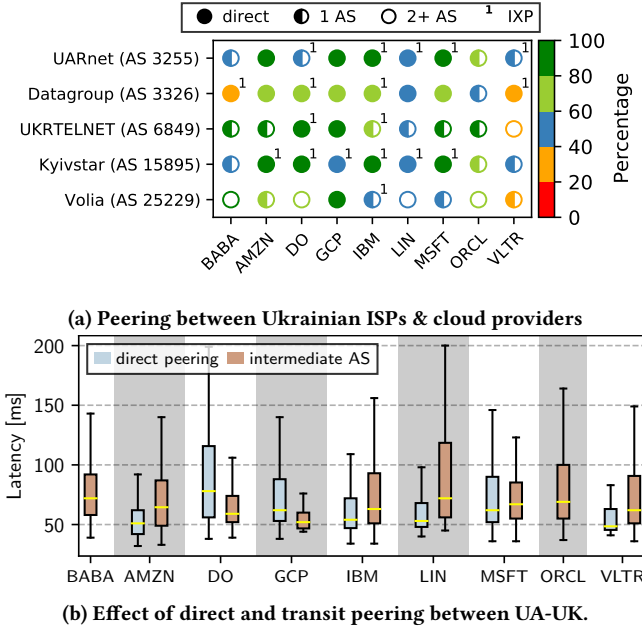


Figure 17: Another case study of ISP-cloud peering in Europe. (a) identifies peering interconnections from Ukrainian ISPs to DCs in the United Kingdom while (b) shows impact of those interconnections on cloud access latency.

#### A.4 Cloud-ISP Peering Additional Case Studies

In §6.2 we analyzed the impact of different cloud-ISP peering interconnections in Europe (from German ISPs to UK cloud DCs) and Asia (from Japanese ISPs to UK cloud DCs). However, our analysis can be considered incomplete due to the following reasons. Firstly, the affect of peering within Europe cannot be accurately assessed by only focusing on interconnections between the two countries (DE and UK) that are largely known for their well-provisioned network backhaul. Secondly, Japan-to-India connectivity does not pose itself as the most compelling use-case for peering in Asia since Japan also has a well-provisioned Internet backhaul (barring its dependence on submarine cables for global connectivity) along with a dense deployment of datacenters within the country itself.

To generalize our analysis on ISP-cloud peering within these two continents, we present additional connectivity case studies. For Europe, we examine the connections from VPs in Ukraine (UA) to DCs in the UK (see Figure 17). For Asia, we analyze peering between serving ISPs from Bahrain (BH) to DCs in India (see Figure 18). The colors in Figure 17a and Figure 18a denote the percentage of paths in  $\langle \text{ISP}, \text{cloud} \rangle$  pairings that used the same interconnection type. We refer the reader to §6.1 for our methodology on identifying ISP-cloud interconnections. Figure 17b and Figure 18b compares the impact of different interconnection types on cloud access latency. It must be noted that both UA and BH have no local DC availability and must rely on deployment in other countries within the continent via in-land backhaul for cloud connectivity.

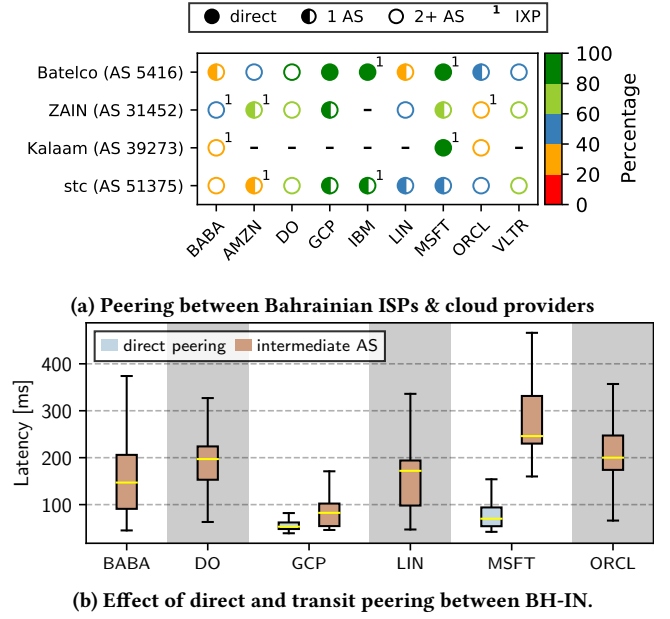


Figure 18: Another case study of ISP-cloud peering in Asia. (a) identifies peering interconnections from Bahrainian ISPs to DCs in India while (b) shows impact of those interconnections on cloud access latency.

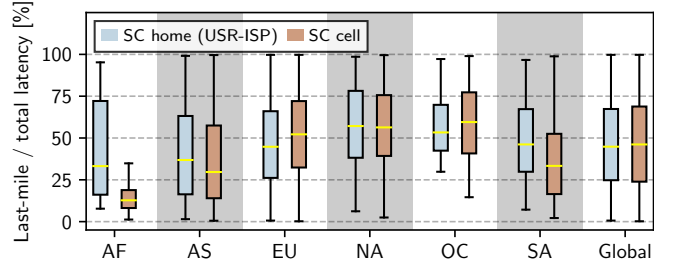


Figure 19: Share of wireless last-mile to the end-to-end latency towards nearest cloud from Speedchecker probes.

We observe that the peering trend in Europe is largely repeatable as the hypergiant cloud providers (Google, Amazon, and Microsoft) have set up direct peering interconnections with most of the Ukrainian ISPs. Similar to our observations in §6.2, we do not see a significant latency advantage for measurements over direct vs. indirect peering links as both achieve comparable median end-to-end latencies. The trend, however, departs significantly for Bahrain-India connections. We find that direct interconnections between these two countries are less common – other than Microsoft and Google peering directly with a handful of serving ISPs. The rest of the cloud providers either provide connectivity via private interconnects or via a public backhaul. Interestingly, here we see a clear latency advantage of direct peering over other interconnection types – with direct peering links achieving significantly lower latencies than their counterparts.

### A.5 Last-mile Latency Share to Closest Cloud

Figure 19 shows the percentage share of wireless last-mile to end-to-end latency towards the nearest cloud datacenter per probe. The result accompanies our last-mile share analysis in §5. Here, SC<sub>home</sub> denotes the last-mile access share of home probes in Speedchecker that likely use WiFi for connectivity. On the other hand, SC<sub>cell</sub> are Speedchecker probes that are likely accessing the cloud via cellular connectivity. The result is similar to those shown in Figure 7a. Firstly, we find that the type of last-mile access (cellular vs. WiFi) does not impact much differently as both technologies exhibit almost similar latency shares. Secondly, the latency at the

last-mile is more pronounced for measurements towards the closest datacenter (understandably so since the overall latency is now much shorter). Here we find that the last-mile can account for most of the connection latency (almost 50% on the global scale) – showcasing it as the primary bottleneck affecting cloud connectivity. Finally, at first glance, only within Africa does cellular connectivity seem to outperform home WiFi connections consistently. However, as noted in §5, this trend is a likely artifact of geographic probe availability within the continent, as most of the home probes are in the south closer to the in-continent DCs while cellular probes are mostly located in the north of the continent (see Appendix A.1).