

Bias in Internet Measurement Platforms

Pavlos Sermpezis[†], Lars Prehn^{*}, Sofia Kostoglou[†], Marcel Flores[‡], Athena Vakali[†], Emile Aben[¶]

[†] Aristotle University of Thessaloniki; {sermpezis, sofikost, avakali}@csd.auth.gr

^{*}Max Planck Institute for Informatics; lprehn@mpi-inf.mpg.de

[‡]Edgio; mfloros@edg.io

[¶]RIPE NCC; emile.aben@ripe.net

Abstract—Network operators and researchers frequently use Internet measurement platforms (IMPs), such as RIPE Atlas, RIPE RIS, or RouteViews for, e.g., monitoring network performance, detecting routing events, topology discovery, or route optimization. To interpret the results of their measurements and avoid pitfalls or wrong generalizations, users must understand a platform’s limitations. To this end, this paper studies an important limitation of IMPs, the *bias*, which exists due to the non-uniform deployment of the vantage points. Specifically, we introduce a generic framework to systematically and comprehensively quantify the multi-dimensional (e.g., across location, topology, network types, etc.) biases of IMPs. Using the framework and open datasets, we perform a detailed analysis of biases in IMPs that confirms well-known (to the domain experts) biases and sheds light on less-known or unexplored biases. To facilitate IMP users to obtain awareness of and explore bias in their measurements, as well as further research and analyses (e.g., methods for mitigating bias), we publicly share our code and data, and provide online tools (API, Web app, etc.) that calculate and visualize the bias in measurement setups.

I. INTRODUCTION

Public Internet Measurement Platforms (IMPs) like RIPE Atlas [3], RIPE RIS [4], or RouteViews [5] are fundamental building blocks of networking research and operations. Network operators and researchers frequently use their measurement capabilities and publicly archived data to, e.g., detect routing events and malicious networks [25], [53], [57], analyze the Internet’s structure [11], [32], [41], understand and optimize (their own) routing policies [30], [51], [56], or detect outages and performance bottlenecks [28], [54], [58].

IMPs operate a broad range of globally distributed vantage points. While RIPE Atlas hosts around 11,000 measurement probes in 3,300 autonomous systems (ASes), RIPE RIS and RouteViews collect routing information from around 300 and 500 ASes, respectively. Despite their presence in thousands of ASes, IMPs only provide a limited view into the routing ecosystem. It is well-known that IMPs capture incomplete views of the Internet [7], [11], [29], [41], [46] and sometimes offer misleading or incomplete answers for seemingly simple questions [23], [34], [49], [59]. This incompleteness problem resulted in approaches for extending the observed AS topology via other data sources [11], [14], [20], [24], [29] or by adding new, favorably-positioned vantage points to IMPs [22], [32], [37], [48]. While deploying IMP infrastructure aiming to increase completeness (e.g., "hunting for the most AS links" by deploying route collectors to IXPs) has a clear

value, it frequently leads to *unequal* (or, "biased") visibility of different parts of the Internet. This bias can come along many dimensions, such as network types, geographic placement, etc.

Despite extensive studies on *incompleteness*, it remains unclear how *representative* our view of the entire Internet routing ecosystem is: *Do we have equal visibility to all types of networks? And, if not, how biased are the views we obtain from IMPs?* In this paper, we study this unexplored aspect and take first steps towards a comprehensive characterization of the bias in IMPs. Contrary to previous works that focus on specific aspects of bias, we argue that capturing representativeness is an inherently multi-dimensional problem. To this end, we make the following contributions:

- We formally define bias, and introduce a generic framework for quantifying the bias in a multi-dimensional context (§III-A). The framework receives as input information about characteristics of networks (connectivity, location, etc.), and quantifies how representative a set of vantage points is.
- We aggregate information from real-world datasets (§III-B), and apply our framework to quantify the bias of widely-used IMPs (§IV). Despite the numerous limitations that come with real-world data sets (e.g., abstractions, inaccuracies, or incompleteness), we observe that our framework is capable of replicating the findings of previous studies (e.g., [12], [13], [49]) conducted by domain experts. Besides these well-known issues, our framework can produce novel, and more nuanced, insights about the bias in IMPs. For example, our analysis confirms that RIPE RIS is heavily biased towards larger networks and IXPs [49], and it also reveals that while networks that peer at many IXPs are over-represented in RIPE RIS, their peering policies (PeeringDB) are representative of the Internet’s peering ecosystem.
- We extend our analysis to explore the improvement potential of IMPs (§IV-A), and study the biases involved in common measurement practices, such as RIPE Atlas probes selection or use of individual route collectors (§IV-B).
- We publicly share our code and data [1], and discuss how it can be parametrized to extend or adapt our analyses. Moreover, to further facilitate users to explore and quantify bias in IMPs or in custom measurement setups, we provide an API and a web portal with interactive visualizations (§V). We believe that having a framework to systematically quantify bias (and tools that automate it) can be valuable for Internet measurements: e.g., from *raising awareness* to users (and

lowering the bar for domain expertise) so that they avoid pitfalls in interpretation of measurement data, to generating a "bias assessment" for each measurement study.

We deem our work as a first step in this direction. In §VII, we provide a critical discussion about how our analysis can be extended or refined to overcome existing limitations, and future research directions that can build upon our framework.

II. INTERNET MEASUREMENT PLATFORMS AND BIAS: A PRIMER

In this section, we introduce the concept of bias on a general example (summarized in Table I). Afterwards, we introduce the three major IMPs that we analyze in this paper (§II-A), discuss some of their known biases, and motivate the research questions that our study aims to address (§II-B).

Let us assume a population consisting of 100 people, 50 of which are men and 50 women. If we run a survey with 10 people, of which 8 men and 2 women, our sample is biased towards men. We say that our sample is biased as there is a *difference in the distributions between the entire population and our sample*.

Table I: Bias example: population and sample statistics. The gender bias is 0.22 (KL-divergence metric; §III-A) and is higher than the country bias 0.03.

	Men	Women	Country A	Country B
Entire population	50%	50%	70%	30%
Survey sample	80%	20%	80%	20%

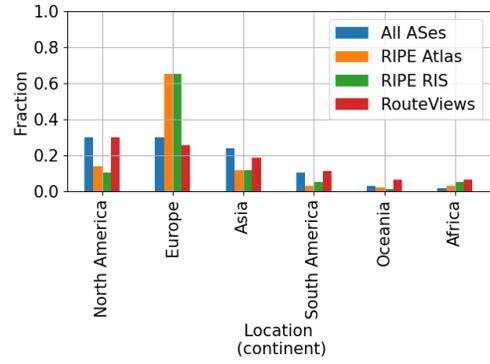
Measuring bias. To *identify* this bias, one could run statistical tests (e.g., Kolmogorov-Smirnov test) to compare the two distributions. To further *quantify* the bias, it is common to measure the *distribution distance* among the population and the sample distributions (e.g., with the Kullback-Leibler divergence metric)

Multi-dimensional bias. Let us consider that our survey focuses on the height of individuals. If we compare the distributions of height within our total population to that within our survey sample, we may find that they differ as men (who naturally tend to be around ~7% taller [43]) are over-represented. Now, let us consider that our survey further focuses on the native language of individuals. For this second case, the gender-bias in our sample would not affect our findings. In contrast, the country-bias (e.g., see the right side of the Table I) of our sample, may play a major role. In other words, *different bias dimensions (e.g., gender or country) may affect our measurements findings differently, depending on how they relate to the insights we want to gain*.

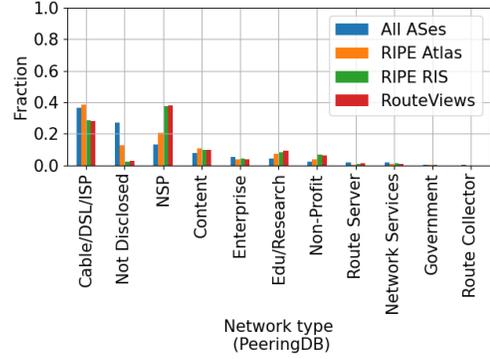
A. IMPs: RIPE Atlas, RIPE RIS, RouteViews

We briefly overview the 3 major IMPs on which we focus.

RIPE Atlas [3] is a platform that hosts more than 11,000 measurement "probes" in more than 3,000 ASes. Probes support a fixed set of measurement types (e.g., ping, traceroute, DNS). Users can select sets of probes and execute measurements (e.g., a traceroute towards a target IP), under some rate-limits.



(a) Location



(b) Network type

Figure 1: Distributions of (a) *Location (continent)* and (b) *Network type* for the entire population of ASes (blue bars) and the set of ASes that are part of the RIPE Atlas, RIPE RIS, and RouteViews platforms.

RIPE RIS [4] and RouteViews [5] are two global platforms that host "route collectors", which are dedicated devices that passively receive, dump, and publicly archive the routing information from their peering networks. Most route collectors are located at large IXPs such that they can quickly establish many sessions over the IXP's peering LAN. The "multi-hop"-enabled route collectors may establish indirect sessions with remote ASNs. In total, RIPE RIS and RouteViews host 27 (of which 3 multi-hop) and 36 (20 multi-hop) route collectors with more than 500 and 300 peer ASNs, respectively. A peering ASN may provide feeds for the entire routing table ("full feed") or only a part of it.

B. IMP biases: known and unknown aspects & user awareness

Location bias. A glance at the map with the locations of RIPE's infrastructure (see [40] for Atlas probes and [6] for RIS route collectors) reveals a higher density of the infrastructure in Europe, which is in imbalance with the spread of ASes around the world, i.e., there is location bias in RIPE Atlas and RIPE RIS. While this bias is well-known (or, easy to spot), a clear quantification is missing: *how far are we from an ideal scenario? Or, what is the room for improvement?*

On the contrary, Fig. 1(a) shows that RouteViews has route collectors deployed in more representative locations around

the world. However, *are users of RouteViews and RIS aware of this significant difference (and do they take it into account in their measurements)? And, can combining measurements from both platforms significantly reduce the location bias?*

Topological bias. Route collectors (RIPE RIS and RouteViews) are biased towards larger core networks and at Internet eXchange Points (IXPs) as reported in previous work [49]. *How could we enable novice users (without "10 years" of experience [49]) to easily identify such biases?*

Bias awareness. Neither the location nor the topological bias are new to expert users. However, not even expert-users might be able to accurately judge the extent of different biases on different IMPs. Similarly, other biases along (less prominent) dimensions, such as the network type (see Fig. 1(b)), might be even harder to judge. A questionnaire related to the topic of this paper that we ran (see details in Appendix A) supports the fact that not all users are aware of biases: out of the 50 questioned operators and researchers, only 26 (52%) consider IMPs to be biased, while 28% consider that there is no bias (or, probably not), and 20% "do not know". *This lack of (or, partial) awareness motivates our study to comprehensively quantify the bias in IMPs.*

III. QUANTIFYING BIAS: FRAMEWORK & METHODOLOGY

Similarly to the example of §II where people are characterized by two features (gender and origin country), the IMPs can also be characterized by a multitude of features, such as location, connectivity, traffic levels, etc.. Each characteristic/feature can be considered as a dimension, and the bias can be calculated over each dimension. Then, depending on the measurement use case, all or some of the dimensions can be taken into account, depending on their relevance (see §II).

In this section, we first introduce our framework, and formally define the bias and the metrics to quantify it (§III-A). The framework is generic: it takes as input a dataset of network characteristics, and returns in a unified way the bias along all characteristics (i.e., "bias dimensions"). Then, in §III-B, we present the dataset we compile and use in this paper as input to our framework for analyzing the bias of IMPs.

A. The bias quantification framework

Definitions. Let P be the distribution of a characteristic (e.g., network size) within a set of networks \mathcal{N} . If the characteristic takes K distinct values, its distribution is $P = [p_1, \dots, p_K]$, where p_i is the probability of a network having the i value (e.g., $p_{Europe}=0.32$ for the entire population of ASes; see Fig. 1(a)); formally, $p_i = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}} I_{j \rightarrow i}$, where $I_{j \rightarrow i}$ an indicator function that is 1 if the network j has the characteristic i , and $|\mathcal{N}|$ the size of the set \mathcal{N} .

Also, let a subset of networks $\mathcal{V} \subset \mathcal{N}$, and Q be the corresponding distribution within the set of networks \mathcal{V} . We define the bias (of the set \mathcal{V} wrt. to the set \mathcal{N}) as the distance between the distributions P and Q .

Identifying bias. If the distance between P and Q is statistically significant, then there is bias. There are several statistical

tests that could be applied. We use the Kolmogorov-Smirnov (or, KS-test), which is a nonparametric test that compares two distributions (two-sample KS-test). The KS-test answers the question "what is the probability that P and Q are drawn from the same distribution?". If this probability is small enough (e.g., less than 5%), then we can confidently state that the population (i.e., all ASes) and sample (i.e., VPs) follow different distributions, i.e., there is bias.

Bias Metrics. There exist several metrics to quantify the distance between two distributions. A metric that is commonly used (in particular, for concepts related to bias, e.g., [55]) is the Kullback–Leibler (KL) divergence:

$$B_{KL} = \sum_{i=1}^K p_i \cdot \log \left(\frac{p_i}{q_i} \right) \quad (1)$$

The KL-divergence takes on values in $[0, +\infty]$, where the higher the value the more the two distributions differ. In the paper, we use a bounded version of the KL-divergence that takes values in $[0, 1]$ [50], [55]¹, and we call it the *bias score*. We calculate a bias score per characteristic/dimension.

For example, in terms of the location distributions depicted in Fig. 1(a), the bias score for RIPE Atlas and RIPE RIS is $B_{KL} = 0.06$ and $B_{KL} = 0.07$, respectively, while for RouteViews, which follows a similar distribution to the entire population, the bias score is $B_{KL} = 0.01$. For the network type (Fig. 1(b)) the bias scores for RIPE Atlas, RIPE RIS, and RouteViews are 0.03, 0.12, and 0.11, respectively, clearly highlighting the higher bias in the route collector projects.

Remark: We tested other common metrics (e.g., Total Variation) for the bias score as well (see Appendix B). While the actual values of each metric are different, the qualitative findings of the paper remain the same.

The framework is generic with respect to \mathcal{N} , \mathcal{V} , bias dimensions, bias metrics, and input data. In this paper, we consider as

- \mathcal{N} : the entire population of ASes (i.e., more than 100,000 ASNs for which we have data)
- \mathcal{V} : the set of VPs of an IMP (e.g, the peers of RIPE RIS or RouteViews, or the probes of RIPE Atlas)

and in the next section we compile a dataset of several network characteristics/dimensions. Later, in §V, we discuss how other choices can be done for these parameters.

B. Data and bias dimensions

Data sources. We take into account the characteristics of the IMPs at an AS-level granularity (e.g., two RIPE Atlas probes in the same AS have the same AS-level characteristics). The reasons for this choice is twofold: data availability and scope. Specifically, at the AS-level there are several public datasets: at a finer granularity there is scarce information (which would limit our analysis to only a few dimensions) and compiling a rich dataset would need extensive measurements

¹We substitute $q_i \rightarrow (1-w) \cdot q_i + w \cdot p_i$, with $w = 0.01$ and normalize with its upper bound $\log \frac{1}{w}$, to get $B_{KL} = \frac{1}{\log \frac{1}{w}} \cdot \sum_{i \in \mathcal{K}} p_i \cdot \log \left(\frac{p_i}{(1-w) \cdot q_i + w \cdot p_i} \right)$.

per dimensions (which could be done only per use case, and thus would be beyond the scope—and space limitations—of this paper). Nevertheless, our framework is extensible to a more fine-grained level (e.g., per monitoring device, such as at a vantage point level or router level); we discuss these extensions and limitations of our analysis in §VII.

We compile a list of characteristics for each AS from the widely used CAIDA AS-rank [17] and AS-relationships [18] and PeeringDB [19], [44] datasets, as well as from public datasets that contain information for the network size/importance (Internet Health Report’s AS-hegemony metric [26], [35], and the Country-level Transit Influence index, or CTI, [27]) and network types (bgp.tools [15] and ASDB [61]).

"Vantage Points (VPs)". Since we study bias at an AS-level, in the remainder, we will not differentiate between different probes in RIPE Atlas that are hosted in the same AS, or between different peers of RIPE RIS and RouteViews with the same ASN. And, for brevity, we will refer to the ASes that host RIPE Atlas probes or provide feeds to RIPE RIS / RouteViews as "vantage points" or VPs.

Dimension categories. From the datasets we select all the characteristics that relate to the concept of bias, in order to make our analysis as general as possible. We end up to a set of 22 characteristics that relate to the concept of bias and group them in the following categories:

- **Location:** RIR region; Country; Continent
- **Network size:** Customer cone (#ASNs,#prefixes,#addresses); AS hegemony; CTI "origin" and "top" indices
- **Topology:** #neighbors (total, peers, customers, providers)
- **IXP-related:** #IXPs; #facilities; Peering policy
- **Network type:** Net. type; Traffic ratio; Traffic volume; Scope; Personal ASN; ASDB classification (level 1 and 2)

Remark: It is important to note that our methodology is generic and more characteristics can be included or grouped differently. We only use these groups to facilitate the discussion in the paper (i.e., to refer to multiple dimensions under a single term), but we present detailed results for all dimensions.

Figure 2 depicts an example of the compiled dataset (which is also available in [1]).

IV. ANALYZING IMP BIAS

In this section, we study the biases in RIPE Atlas, RIPE RIS, and RouteViews. Figure 3 shows a radar plot with bias scores for all dimensions. The colored lines—and their included area—correspond to the bias metric of a given IMP along a given dimension, e.g., the bias score for RIPE RIS (orange line) in the dimension “Location (country)” is 0.2. Larger bias scores (i.e., farther from the center) correspond to more bias, e.g., in the dimension “Location (country)” RIPE RIS is more biased than RIPE Atlas (blue line).

Remark: As knowing the entire distribution of a characteristic may help to better understand the bias along a certain dimension, we provide detailed distribution plots (i.e., similar to those in Fig. 1) for all characteristics in the extensive documentation of our code and data [1].

ASN	Location-related information			Network-size related information			Topology-related information			IXP-related information		Network type-related information	
	RIR Region	Country	Continent	Customer cone (in #ASNs)	AS hegemony	...	#neighbors (in #ASNs)	...	#IXPs connected to	...	Net. type (PeeringDB)	Net. type (ASDB)	...
174	ARIN	US	North America	32457	0.09	...	6614	...	0	...	NSP	ICT	...
1299	RIPE	SE	Europe	37162	0.10	...	2328	...	0	...	NSP	ICT	...
2497	APNIC	JP	Asia	507	0.01	...	338	...	16	...	NSP	NaN	...
3320	RIPE	DE	Europe	3015	0.01	...	667	...	5	...	NSP	ICT	...
3333	RIPE	NL	Europe	3	0.00	...	320	...	1	...	Non-profit	ICT	...
5470	RIPE	GR	Europe	1	0.00	...	1	...	NaN	...	NaN	Education & Research	...
15169	ARIN	US	North America	12	0.01	...	366	...	214	...	Content	ICT	...
...

Figure 2: An example depicting the compiled dataset with characteristics (columns) of ASes (rows).

Key findings. Based on Fig. 3, we can observe that:

- While the bias of IMPs differs significantly by dimension, RIPE Atlas is substantially less biased than RIPE RIS and RouteViews along most dimensions.
- RIPE RIS and RouteViews have significant topological bias (e.g., number of neighbors/peers) as most of their collectors are deployed at IXPs, where ASes establish many (peering) connections [46].
- RouteViews and RIPE RIS are also quite biased in terms of network size (“Customer cone” dimensions) because they peer with many large ISPs. Having feeds from large ISPs may be desired for visibility, however, users still should be aware of it since it may lead to biased measurements.
- In most IXP-related and network type dimensions (that correspond to data mainly from PeeringDB), all platforms have relatively low bias; with an exception of RIPE RIS and RouteViews that are biased in terms of number of IXPs/facilities the VPs are connected to.
- There are small differences between RIPE RIS and RouteViews. RIPE RIS is more biased in terms of topology (number of neighbors, total and peers), whereas RouteViews is more biased in terms of network sizes (“Customer cone” and “AS hegemony” dimensions).
- We applied the KS-test for all platforms and dimensions. In almost all cases, the KS-test rejected the null hypothesis that the IMPs vantage points follow the same distribution as the entire population of ASes. The only exceptions were the "Personal ASN" dimension for all IMPs, and the "RIR region" and "Location (continent)" for RouteViews (where bias scores are less than 0.01).

Table II shows the correlation between the network characteristics for the entire population of ASes. The characteristics are grouped in the categories of §III-B, and values correspond to averages among groups (i.e., values in the diagonal are not 1)². As expected, dimensions in the same category are correlated. Also, topology dimensions are significantly correlated with IXP-related dimensions. Nevertheless, comparing with Fig. 3, we see that correlated dimensions do not necessarily share similar bias scores. This highlights that a multi-dimensional bias exploration (Fig. 3) can give more insights.

²Since our dataset consists of both numerical and categorical data, we use (i) the Pearson correlation coefficient for pairs of numerical features, (ii) the correlation ratio for pairs of a numerical and a categorical feature, and (iii) Cramer’s V test for correlations between categorical features.

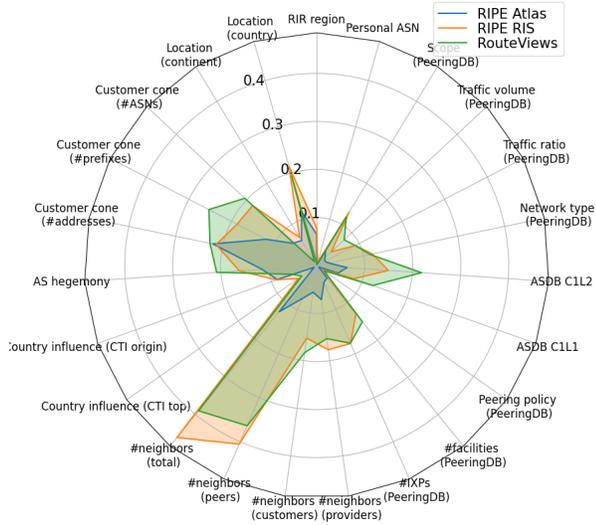


Figure 3: Radar plot depicting the bias score for RIPE Atlas (blue line/area), RIPE RIS (orange line), and RouteViews (green line) over the different dimensions (radius of the circle). Larger values of bias scores (i.e., far from the center) correspond to more bias.

Beyond this basic analysis, we conduct three similar analyses deepening our understanding of different IMP aspects.

Combining RIS and RouteViews. Using data from both RIPE RIS and RouteViews is common (e.g., via CAIDA BGPstream [42]); hence, we analyze the combined bias in Fig. 4(a). When considering vantage points from both projects, the bias slightly decreases in most dimensions. Interestingly, there are some exceptions, e.g., number of neighbors (total and peers), where it would be preferable—in terms of bias—to use only feeds from RouteViews.

Full vs. all feeds. Only 240 and 70 peers of the RIPE RIS and RouteViews peers provide feeds for the entire routing table ("full feeds"), respectively. Figure 4(b) compares the bias of only full feed peers against the entire IMPs. For RIPE RIS the increase in bias is small, whereas for RouteViews the set of full feeds is significantly more biased. In fact, while RIPE RIS is on average more biased than RouteViews, the opposite becomes true when considering only full feeds.

IPv4 vs IPv6 vantage points. Figure 4(c) compares the set of ASes hosting IPv4, IPv6, and all RIPE Atlas probes (i.e., both IPv4 and IPv6). The set of networks hosting IPv6 probes

Table II: Correlations between dimensions categories for the entire population of ASes.

	Location	Net. type	IXP-rel.	Topology	Net. size
Location	0.99	0.15	0.15	0.24	0.21
Net. type		0.31	0.33	0.26	0.17
IXP-rel.			0.75	0.48	0.26
Topology				0.69	0.38
Net. size					0.40

is slightly more biased than networks hosting IPv4 probes in most dimensions. The only exception is the #addresses in customer cone, which is mainly due to the differences in the IP space between the two versions. In RIPE RIS (not depicted in the plot), the differences between IPv4 and IPv6 peers is negligible. Due to the similarity in our analyses between IPv4 and IPv6 VPs, in the remainder we do not present separate results for each of these subgroups.

A. Analyzing improvement potential

Now that we have a basic understanding of the current biases in IMPs, we want to compare the current state to a (hypothetical) case, where vantage points are randomly deployed among all types of networks, locations, etc. This comparison (i) provides a better understanding of the potentially avoidable IMP bias, and consequently (ii) reveals room for improvement (under practical limitations).

Random sampling from the entire population is an unbiased process. A sufficiently large random sample would lead to zero bias. Yet, small samples tend to be biased especially for characteristics with large variance. We treat the bias score that can be achieved via random sampling as a non-biased baseline.

Table III compares the average bias over all dimensions³ of the IMPs against that of a random sample with the same number of vantage points (e.g., in the case of RIPE RIS we consider random samples of size $|\mathcal{V}|=539$). We repeat our random sampling 100 times and report the average bias. We observe that *with the same number of VPs as in the current IMPs, a random sample of ASes would have on average (almost) no bias*. This indicates that the "limited" number of vantage points is not the root cause of bias (which is mostly due to the deployment strategies; see §II-B). In other words, *we do not need more VPs, but more representative VPs*. This finding can be valuable for future extensions of IMPs (e.g. selection of deployments in under-represented parts of the Internet) and improvement of measurement techniques (e.g., carefully selecting representative subsets of existing VPs); we identify these two aspects as key future research directions that can stem from our framework.

Bias vs. number of vantage points. While the current set of VPs is clearly not optimal in terms of bias, we wonder how bias changes when we only use a smaller random set of VPs (e.g., measurements with few Atlas probes due to rate/credit limits, or collecting feeds from a subset of route collectors peers due to the large volumes of data [8]). Figure 5(a) shows

³There are infinite options of combining bias scores of different dimensions. Here, we consider averaging as an intuitive choice, however, our framework supports other options as well (e.g., weighted average, or bias for subsets of dimensions).

Table III: Bias of IMPs vs. random sample of vantage points.

Platform (#vantage points)	Atlas (3391)	RIS (539)	RV (340)	RIS & RV (762)
Platform bias	0.06	0.16	0.15	0.14
Random sample bias	0.00	0.01	0.01	0.01

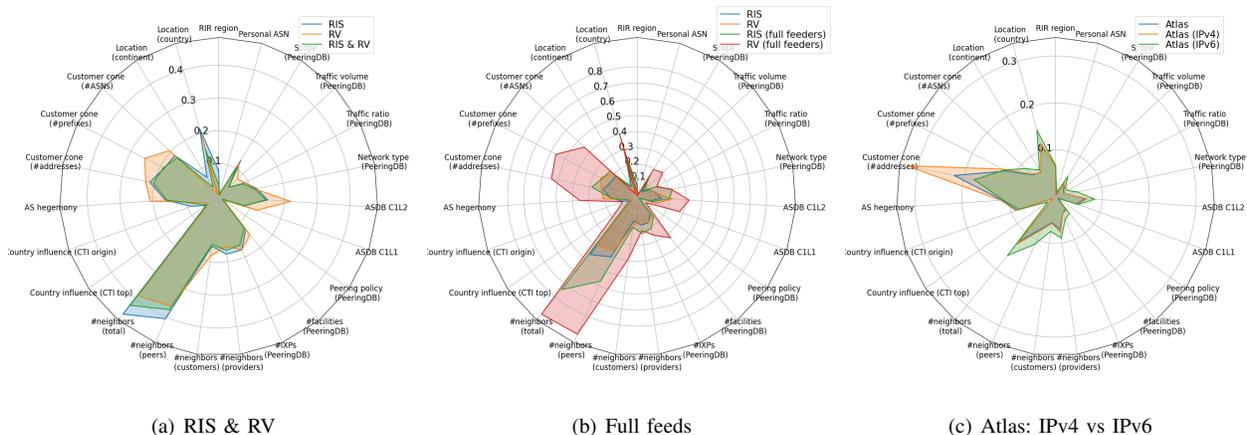


Figure 4: Radar plot of bias score for the cases of (a) the combined RIPE RIS and RouteViews vantage points, (b) full feeds (vs all feeds) for RIPE RIS and RouteViews, and (c) IPv4 vs IPv6 RIPE Atlas probes. Note the different bias ranges in each plot.

the average bias for different sample sizes drawn randomly from either the entire population of ASes ('all') or one of the three IMPs. Lines correspond to averages over 100 sampling iterations, and errorbars indicate 95% confidence intervals. For ease of comparison, dashed lines correspond to the bias values of using the entire infrastructure (i.e., the values in Table III). We observe that: (i) the bias decreases with the sample size (as expected), (ii) random sampling ("all") always has lower bias for the same number of VPs, (iii) even for very small samples, ≥ 20 VPs, random sampling ("all") has lower bias than the *entire* sets of RIPE RIS and RouteViews VPs (dashed lines), while the same holds for RIPE Atlas for ≥ 40 VPs.

For a deeper inspection of the bias in smaller sets of VPs, Fig. 5(b) presents the bias of random samples of RouteViews VPs of sizes 10, 20, and 100 (similar results hold also for RIPE RIS and Atlas). We can see how the bias decreases in all dimensions for larger subset sizes. Yet, the change in bias is not the same in all dimensions, e.g., in network type dimensions the relative increase in bias for small subsets is much larger than in topology. This type of analyses can help to design measurements, for example, to select the number of VPs (e.g., see [47] or [52]) based on the required level of bias per dimension.

B. Bias in common measurement practices

In this section, we briefly analyze the bias involved in common VP selection methods that users follow in practice.

RIPE Atlas probe selection algorithm. RIPE Atlas users can either select specific probes to use in their measurements or not specify them (which is the default choice; with parameters 10 probes from "worldwide locations"⁴). In the latter case, RIPE Atlas has an automated algorithm to assign probes to a measurement, which prioritises probes with less load over more loaded probes, which makes the probe selection procedure not equivalent to true random sampling.

In Fig. 6(a) we study how the RIPE Atlas selection algorithm, "Atlas (platform)", performs compared to random sampling from either all RIPE Atlas probes, "Atlas (random)", or from all ASes ("all"); the values for these latter cases are the same as in Fig. 5(a)). We considered the sets of probes that the RIPE Atlas platform returned when we initiated measurements with parameters `type="area"` and `value="WW"`. Lines correspond to averages over 100 sampling iterations, and errorbars indicate 95% confidence intervals. We observe that *when using the RIPE Atlas algorithm for selecting probes, "Atlas (platform)", then the bias is significantly higher compared to randomly selecting probes, "Atlas (random)".* In fact, the bias is almost two times higher. This indicates that even with the existing infrastructure, users could decrease bias by 50% by not depending on the built-in probe selection process, but select random probes themselves.

Feeds from a single route collector (RC) may be used in cases that there are processing limitations (e.g., in terms of real-timeness or storage) due to the large volume of data, see [8], [10], [31]. Figure 6(b) presents the average bias score per RC (i.e., the bias of the set of VPs that peer to a RC) in relation to its number of VPs. Overall, there is a clear (negative) correlation between the number of VPs and the bias score of a RC. Nevertheless, the size of a route collector does not predict its bias as (i) the three RCs of RIPE RIS (rrc01, rrc03, rrc12) that are significantly larger (>80 members) than the rest of RCs, are not less biased (in fact, there are several smaller RCs with lower bias) and (ii) there are several medium-size RCs (and even some with only 10-20 VPs) that have relatively low bias. For RIPE RIS, the three multihop RCs (rrc00, rrc24, rrc25) are less biased than most of the non-multihop RCs (which are deployed at IXPs). For further analyses, we provide through our API and online tools (§V; [2]) the detailed bias scores and radar-plots for each RC.

Summary of main takeaways: (i) *RIPE RIS and RouteViews are substantially more biased than RIPE Atlas; and their VPs are significantly more biased towards networks at IXPs (with*

⁴<https://atlas.ripe.net/docs/udm#probe-selection>

V. OPEN DATA, CODE, API, AND TOOLS

To facilitate users and further research and analyses, we provide data, code, and tools to calculate and visualize the bias in a set of networks [1], [2].

Data. The data we aggregated from different sources are provided as a table with rows corresponding to ASNs and columns to network characteristics (see §III-B) [1].

Code. We open-source the code for calculating the bias [1]. The methods receive as input (i) the data table, (ii) a set of ASNs that are considered the “population” \mathcal{N} , (iii) a subset \mathcal{V} of the population, whose bias we are interested in, (iv) the set of characteristics that will be taken into account. This enables users to apply the framework in a generic way; for example:

- use other datasets than the ones we compiled in §III-B
- perform a bias analysis with respect to a given region (e.g., setting as the “population” \mathcal{N} only the ASNs in the RIPE region, instead of all the ASes we considered in this paper)
- explore the bias of a custom set of VPs \mathcal{V} (e.g., a set of Atlas probes, or a set of route collector peers)
- consider only a subset of bias dimensions

Open API. To further facilitate access to data and methods, we provide an API [2] that provides (up-to-date) bias scores of the IMPs we analyzed, other IMPs (e.g., bpg.tools, CAIDA Periscope), individual route collectors (see §IV-B), or any custom set of VPs (or, ASNs in general) requested by a user.

Web portal. We provide a set of online tools for interactive visualizations of the bias data, namely, radar plots (as in Fig. 3) and the detailed distributions per characteristic (CDF plots or histograms) for all platforms [2].

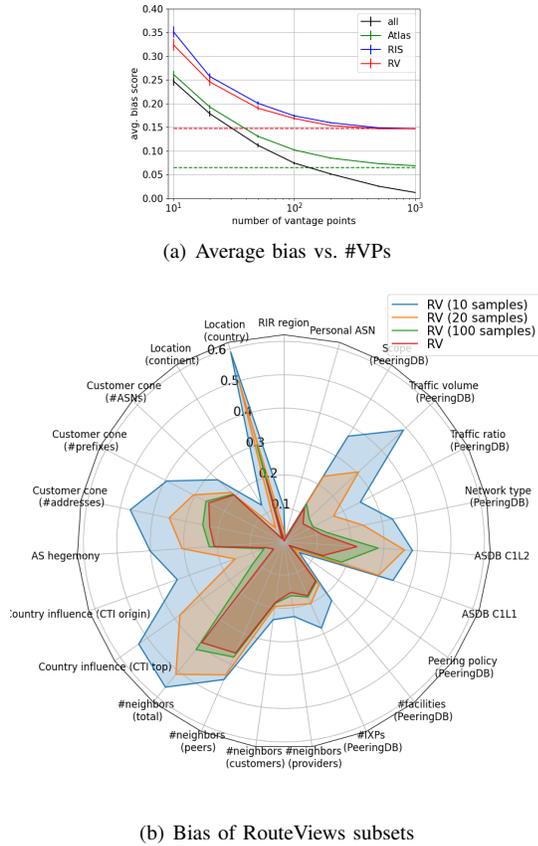


Figure 5: Bias vs. #VPs: (a) Average bias (y-axis) of random samples from the entire population of ASes (“all”) and from the IMP VPs vs. sample size (x-axis); (b) Bias of different sample sizes of RouteViews VPs in all dimensions.

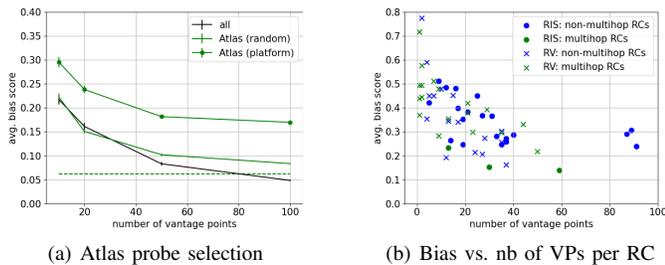


Figure 6: (a) Average bias score vs. sample size for the RIPE Atlas probe selection algorithm vs. random sampling. (b) Scatter plot of average bias (y-axis) vs. number of VPs (x-axis) per route collector of RIPE RIS and RouteViews.

many peering links) and larger networks. (ii) Considering only full-feeds further intensifies the bias of control-plane IMPs, while differences between IPv4 and IPv6 VPs are less important in all IMPs. (iii) If IMPs would choose VPs entirely random, their current set of VPs would be very close to an ideal sample; this indicates that bias is not due to the size of IMPs, but mainly due to VP deployment strategies. (iv) Common practices to limit the number of VPs yield higher bias than simple random samples from IMPs.

VI. RELATED WORK

Topological bias of route collectors. When analyzing the Internet’s topology, route collectors often miss many interconnections of CDNs [11], at IXPs [7], [29], [46], or due to complex routing setups [41]. While it is hard to remove these biases, many works tried to understand the importance of certain biases for their work by analyzing how their results would change when using only subsets of the available infrastructure, e.g., [36], [38], [39], [52], [47]. While we do not focus on the effect of biases in measurements, our framework enables to easily quantify biases, e.g., in these subsets of the infrastructure, and thus provide further insights on how it affects or correlates with the resulting measurement biases.

Biases and use cases. While Roughan et al. argued that route collectors are biased towards larger core networks and IXPs [49], Chung et al. [21] saw no substantial differences when comparing their view on the longitudinal deployment of route origin validation with that of Akamai gathered from an order of magnitude more monitors⁵. On the other hand, [45] shows that IMP data can lead to significant geographical and

⁵The study only analyzed prefix-origin pairs that were visible by the route collectors. It remains unclear whether this result would change when also considering Akamai’s privately received BGP announcements.

topological biases in AS relationships inference. This highlights that biases might be use-case dependant—a fact further supported by the work of Cittadini et al. from 2014 which showed that route collectors have different biases for topology analysis and iBGP policy inference [22]. In 2009, [33] argued that Internet measurements, in general, are biased by various (sometimes unknown) factors such as traffic volume, user populations, or topology, which is further supported by a series of exemplary experiments conducted by Bush et al. [16].

Bias in RIPE Atlas. In 2015, Bajpai et al. showed that the distribution of Atlas probes to ASes is heavy-tailed and also analyzed the network type distribution of probe hosting ASes (without comparing it to the overall type distribution) [12]. A later study by Bajpai et al. in 2017 further found that 91 % of RIPE Atlas probes are located in the RIPE and ARIN region and that the number of probes is not representative for the number of Internet users in countries such as Japan [13]. These types of bias explorations are facilitated (and extended to more bias dimensions) by the proposed framework.

Similarity of VPs. Two recent studies [8], [9] considered similarity of RIPE Atlas and RIPE RIS VPs, respectively, for subsampling VPs and thus avoiding redundant (i.e., over-represented) information. [9] calculates a similarity matrix between Atlas probes based on measurements, and proposes a method to select subsets of probes that are dissimilar. Similarly, [8] calculates VPs similarities based on topological characteristics, and applies a clustering algorithm to select a set of dissimilar of VPs aiming to achieve a good tradeoff between volume of information (i.e., less VPs) and observability of the AS topology. We consider these works complementary to ours; a difference is that [8], [9] are measurement-dependent, however, investigating relation between bias and VP similarities can lead to more efficient subsampling methods.

VII. CONCLUSION

This work aims to be the first effort for a systematic and comprehensive characterization of bias in IMPs, by providing a framework to quantify bias (metrics, data, code, etc.) and an analysis of popular IMPs. Being aware about the existence of bias and its "flavors" (e.g., how much and at what dimensions) can help the users of IMPs to carefully interpret the results of their measurements, and avoid pitfalls or wrong generalizations that may appear due to the bias.

Before our work, significant biases in IMP vantage point placements have been documented by experience papers from well-established scientists (e.g., [16], [49]) or via a few dedicated analyses [12], [13], [34]. Besides reproducing their original findings, the framework we introduced drastically facilitates finding new biases among diverse dimensions and tracking of the evolution of these biases over time.

Moreover, our findings and tools (data, code, API) can further help users to fine-tune their measurements (e.g., select a set of vantage points), and provide useful insights to IMP operators for extending their platforms. We see several promising messages in our results towards these directions.

We deem our work as an initial (but, necessary) step towards a complete understanding of bias in IMPs and its impact on user measurements. There are many research directions and improvements that would need a more extensive investigation and can be addressed in future work.

In the following we provide a critical discussion for some of these directions, in relation to our work:

AS-level granularity. We conducted our analysis at an AS-level, because the majority of data sources provide data at this granularity. It is straightforward to generalize our framework to a more fine-grained level (all methods, metrics, etc., directly apply). For example, if we have available data per prefix⁶, then we can consider as our "population" all the routed prefixes, and as "sample population" the prefixes that contain the IP addresses of the RIPE RIS / RouteViews peers or the RIPE Atlas probes. Our methods would then simply take as input a table as in Fig. 2 with rows the prefixes (instead of ASes) and columns the prefix characteristics (instead of AS characteristics).

Several use cases could benefit from such a more fine-grained granularity. However, the challenging part is the data availability. To extract even a single characteristic at this granularity, we may need extensive measurements and analyses. For example, a custom method is needed to infer per-prefix locations [60], while to infer customer cones per-prefix could lead to incomplete data since aggregating measurements from VPs in different prefixes would not be possible.

Dimensions of bias (per use case). Not all dimensions of bias may be relevant to a measurement study. For example, any bias in the "peering policy" dimension may not affect latency measurements (this is just a conjecture), whereas it is probable to affect BGP hijacking detection measurements. In another example, [52] has shown that estimating the impact of a hijack with RIPE RIS and RouteViews leads to a 10% higher error than custom measurements from random (i.e., unbiased) ASes; this error due to bias may be lower/higher in a different use case though. Identifying which dimensions are important per use case, could improve our understanding of bias and its role. However, this requires a per case analysis, since there are many different measurement use cases with a wide range of scopes and objectives. Our tools (§V) enable to exclude dimensions, thus, covering as many use cases as possible.

Accuracy, completeness, and bias in ground truth data. The input to our framework (i.e, the AS characteristics) is from public datasets. And, some of them are known to suffer from inaccuracies (e.g., country information per ASNs), incompleteness (e.g., only 25% of ASNs have records in PeeringDB), or even biases (e.g., data inferred based on measurements from the existing –biased– platforms, such as customer cones, topology, etc.). Improving the datasets would be beneficial, in general, and for the quantification of bias, in particular,

⁶Some ASes consist of many—sometimes globally distributed—routers that make independent decisions, which can be captured at a prefix level.

since they could reveal further insights⁷; nevertheless, this is an orthogonal task. As already discussed, changing the datasets does not change our framework, but only its input.

ACKNOWLEDGEMENTS

This research is co-financed by Greece and European Union through the Operational Program Competitiveness, Entrepreneurship and Innovation under the call RESEARCH-CREATE-INNOVATE (project T2EDK-04937) and RIPE NCC (AI4NetMon project).

REFERENCES

- [1] AI4NetMon - github repository. <https://github.com/sermpezis/ai4netmon>
- [2] The AI4NetMon project and tools. <https://ai4netmon.csd.auth.gr/>
- [3] RIPE Atlas. <https://atlas.ripe.net/> (2023)
- [4] RIPE RIS. <https://ris-live.ripe.net/> (2023)
- [5] RouteViews. <http://www.routeviews.org/> (2023)
- [6] Aben, E.: RIPE RIS route collectors map. <https://observablehq.com/ris-route-collectors-and-peer-locations> (2023)
- [7] Ager, B., Chatzis, N., Feldmann, A., Sarrar, N., Uhlig, S., Willinger, W.: Anatomy of a large european ixp. In: Proc. ACM SIGCOMM (2012)
- [8] Alfroy, T., Holterbach, T., Pelsser, C.: Mvp: measuring internet routing from the most valuable points. In: Proc. ACM IMC (2022)
- [9] Appel, M., Aben, E., Fontugne, R.: Metis: Better atlas vantage point selection for everyone. In: Proc. IEEE TMA (2022)
- [10] Ariemma, L., Liotta, S., Candela, M., Di Battista, G.: Long-lasting sequences of bgp updates. In: PAM conference. Springer (2021)
- [11] Arnold, T., He, J., Jiang, W., Calder, M., Cunha, I., Giotsas, V., Katz-Bassett, E.: Cloud provider connectivity in the flat internet. In: Proc. ACM IMC (2020)
- [12] Bajpai, V., Eravuchira, S.J., Schönwälder, J.: Lessons learned from using the ripe atlas platform for measurement research. ACM SIGCOMM Computer Communication Review **45**(3), 35–42 (2015)
- [13] Bajpai, V., Eravuchira, S.J., Schönwälder, J., Kistelegki, R., Aben, E.: Vantage point selection for ipv6 measurements: Benefits and limitations of ripe atlas tags. In: IFIP/IEEE IM (2017)
- [14] Battista, G.D., Refice, T., Rimondini, M.: How to extract bgp peering information from the internet routing registry. In: Proc. SIGCOMM workshop on Mining network data (2006)
- [15] BGP.Tools: Networks that have the following tag: Personal ASN. Available at <https://bgp.tools/tags/perso> (2022), last-accessed: Sunday, 8th May 2022
- [16] Bush, R., Maennel, O., Roughan, M., Uhlig, S.: Internet optometry: assessing the broken glasses in internet reachability. In: Proc. ACM IMC (2009)
- [17] CAIDA: AS-rank dataset. Available at <https://asrank.caida.org/> (2023)
- [18] CAIDA: AS-relationships dataset. Available at <https://publicdata.caida.org/datasets/as-relationships/> (2023)
- [19] CAIDA: PeeringDB Dataset. Available at <https://publicdata.caida.org/datasets/peeringdb/2022/04/> (2023)
- [20] Chen, K., Choffnes, D.R., Potharaju, R., Chen, Y., Bustamante, F.E., Pei, D., Zhao, Y.: Where the sidewalk ends: Extending the internet as graph using traceroutes from p2p users. In: Proc. ACM CoNEXT (2009)
- [21] Chung, T., Aben, E., Buijnzeels, T., Chandrasekaran, B., Choffnes, D., Levin, D., Maggs, B.M., Mislove, A., Rijswijk-Deij, R.v., Rula, J., et al.: Rpkis coming of age: A longitudinal study of rpkis deployment and invalid route origins. In: Proc. ACM IMC (2019)
- [22] Cittadini, L., Vissicchio, S., Donnet, B.: On the quality of bgp route collectors for ibgp policy inference. In: IEEE/IFIP Networking (2014)
- [23] Del Fiore, J.M., Merindol, P., Persico, V., Pelsser, C., Pescapé, A.: Filtering the noise to reveal inter-domain lies. In: 2019 Network Traffic Measurement and Analysis Conference (TMA), pp. 17–24. IEEE (2019)
- [24] Faggiani, A., Gregori, E., Improta, A., Lenzi, L., Luconi, V., Sani, L.: A study on traceroute potentiality in revealing the internet as-level topology. In: IEEE/IFIP Networking (2014)
- [25] Fontugne, R., Bautista, E., Petrie, C., Nomura, Y., Abry, P., Gonçalves, P., Fukuda, K., Aben, E.: Bgp zombies: An analysis of beacons stuck routes. In: PAM conference. Springer (2019)
- [26] Fontugne, R., Shah, A., Aben, E.: The (thin) bridges of as connectivity: Measuring dependency using as hegemony. In: PAM conference. Springer (2018)
- [27] Gamero-Garrido, A., Carisimo, E., Hao, S., Huffaker, B., Snoeren, A.C., Dainotti, A.: Quantifying nations’ exposure to traffic observation and selective tampering. In: PAM conference. Springer (2022)
- [28] Giotsas, V., Dietzel, C., Smaragdakis, G., Feldmann, A., Berger, A., Aben, E.: Detecting peering infrastructure outages in the wild. In: Proc. ACM SIGCOMM (2017)
- [29] Giotsas, V., Zhou, S., Luckie, M., Claffy, K.: Inferring multilateral peering. In: Proc. ACM CoNEXT (2013)
- [30] Gray, C., Mosig, C., Bush, R., Pelsser, C., Roughan, M., Schmidt, T.C., Wahlsch, M.: Bgp beacons, network tomography, and bayesian computation to locate route flap damping. In: Proc. ACM IMC (2020)
- [31] Green, T., Lambert, A., Pelsser, C., Rossi, D.: Leveraging inter-domain stability for bgp dynamics analysis. In: PAM conference. Springer (2018)
- [32] Gregori, E., Improta, A., Lenzi, L., Rossi, L., Sani, L.: On the incompleteness of the as-level graph: a novel methodology for bgp route collector placement. In: Proc. ACM IMC (2012)
- [33] Heidemann, J., Papadopoulos, C.: Uses and challenges for network datasets. In: Proc. IEEE CATCH (2009)
- [34] Holterbach, T., Pelsser, C., Bush, R., Vanbever, L.: Quantifying interference between measurements on the ripe atlas platform. In: Proc. ACM IMC (2015)
- [35] IIR: Internet Health Report. <https://ihr.ijlab.net/ihr/en-us/> (2023)
- [36] Jin, Z., Shi, X., Yang, Y., Yin, X., Wang, Z., Wu, J.: Toposcope: Recover as relationships from fragmentary observations. In: Proc. ACM IMC (2020)
- [37] Leyba, K.G., Daymude, J.J., Young, J.G., Newman, M., Rexford, J., Forrest, S.: Cutting through the noise to infer autonomous system topology. arXiv preprint arXiv:2201.07328 (2022)
- [38] Luckie, M., Huffaker, B., Dhamdhere, A., Giotsas, V., Claffy, K.: As relationships, customer cones, and validation. In: Proc. ACM IMC (2013)
- [39] Marcos, P., Prehn, L., Leal, L., Dainotti, A., Feldmann, A., Barcellos, M.: As-path prepending: there is no rose without a thorn. In: Proc. ACM IMC (2020)
- [40] NCC, R.: RIPE Atlas probes map. <https://atlas.ripe.net/results/maps/network-coverage/> (2023)
- [41] Oliveira, R.V., Pei, D., Willinger, W., Zhang, B., Zhang, L.: In search of the elusive ground truth: the internet’s as-level connectivity structure. ACM SIGMETRICS Performance Evaluation Review **36**(1) (2008)
- [42] Orsini, C., King, A., Giordano, D., Giotsas, V., Dainotti, A.: Bgpstream: a software framework for live and historical bgp data analysis. In: Proc. ACM IMC (2016), <https://bgpstream.caida.org/>
- [43] Our World in Data: Human Height. <https://ourworldindata.org/human-height> (2019)
- [44] PeeringDB: The Interconnection Database. Available at <https://www.peeringdb.com/> (2023)
- [45] Prehn, L., Feldmann, A.: How biased is our validation (data) for as relationships? In: Proc. ACM IMC (2021)
- [46] Prehn, L., Lichtblau, F., Dietzel, C., Feldmann, A.: Peering only? analyzing the reachability benefits of joining large ixps today. In: PAM conference. Springer (2022)
- [47] Reuter, A., Bush, R., Cunha, I., Katz-Bassett, E., Schmidt, T.C., Wahlsch, M.: Towards a rigorous methodology for measuring adoption of rpkis route validation and filtering. ACM SIGCOMM CCR **48**(1) (2018)
- [48] Roughan, M., Tuke, S.J., Maennel, O.: Bigfoot, sasquatch, the yeti and other missing links: what we don’t know about the as graph. In: Proc. ACM IMC (2008)
- [49] Roughan, M., Willinger, W., Maennel, O., Perouli, D., Bush, R.: 10 lessons from 10 years of measuring and modeling the internet’s autonomous systems. IEEE Journal on Selected Areas in Communications **29**(9), 1810–1821 (2011)
- [50] Sacharidis, D., Mouratidis, K., Klefogiannis, D.: A common approach for consumer and provider fairness in recommendations. In: Proc. ACM RecSys (Late-breaking Results,) (2019)
- [51] Sermpezis, P., Kotronis, V.: Inferring catchment in internet routing. Proc. ACM SIGMETRICS (2019)
- [52] Sermpezis, P., Kotronis, V., Arakadakis, K., Vakali, A.: Estimating the impact of bgp prefix hijacking. In: IEEE/IFIP Networking (2021)
- [53] Sermpezis, P., Kotronis, V., Gigis, P., Dimitropoulos, X., Cicalese, D., King, A., Dainotti, A.: Artemis: Neutralizing bgp hijacking within a minute. IEEE/ACM Trans. on Networking (TON) **26**(6) (2018)

⁷The main insights of this paper are not expected to deviate significantly, since we have not identified any counterintuitive findings in our analysis.

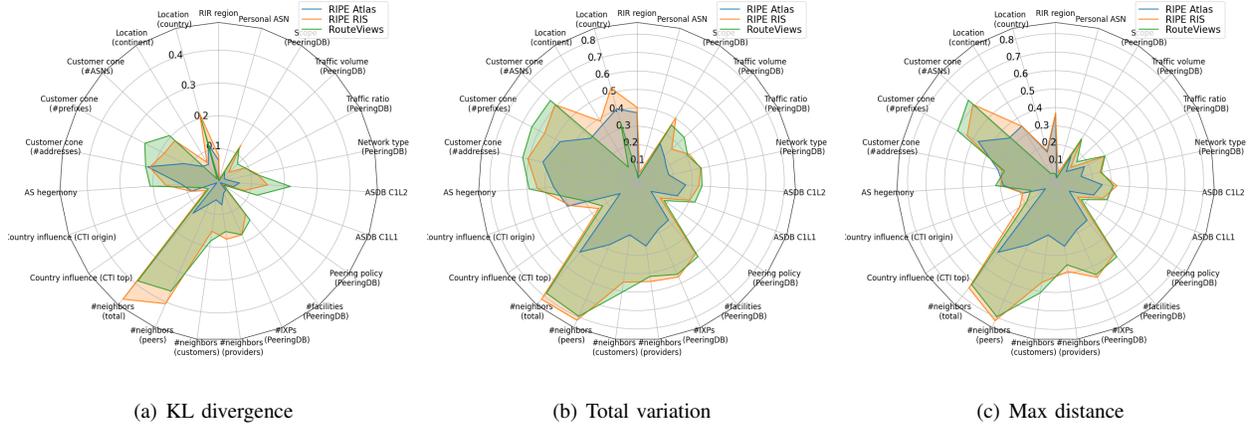


Figure 7: Radar plot depicting the bias score for different bias metrics.

[54] Shah, A., Fontugne, R., Aben, E., Pelsler, C., Bush, R.: Disco: Fast, good, and cheap outage detection. In: 2017 Network Traffic Measurement and Analysis Conference (TMA). pp. 1–9. IEEE (2017)

[55] Steck, H.: Calibrated recommendations. In: Proc. ACM RecSys (2018)

[56] Streibelt, F., Lichtblau, F., Beverly, R., Feldmann, A., Pelsler, C., Smaragdakis, G., Bush, R.: Bgp communities: Even more worms in the routing can. In: Proc. ACM IMC (2018)

[57] Testart, C., Richter, P., King, A., Dainotti, A., Clark, D.: Profiling bgp serial hijackers: capturing persistent misbehavior in the global routing table. In: Proc. ACM IMC (2019)

[58] Verizon: Seeing the World with RIPE Atlas. https://labs.ripe.net/Members/verizon_digital/seeing-the-world-with-ripe-atlas (2017)

[59] Willinger, W., Alderson, D., Doyle, J.C.: Mathematics and the internet: A source of enormous confusion and great potential. Notices of the American Mathematical Society **56**(5), 586–599 (2009)

[60] Winter, P., Padmanabhan, R., King, A., Dainotti, A.: Geo-locating bgp prefixes. In: Proc. IEEE TMA (2019)

[61] Ziv, M., Izhikevich, L., Ruth, K., Izhikevich, K., Durumeric, Z.: Asdb: a system for classifying owners of autonomous systems. In: Proc. ACM IMC (2021)

APPENDIX A

SURVEY ON INTERNET MEASUREMENTS AND BIAS

We conducted an anonymous survey on Internet measurements and bias in IMPs. In this paper, we only provide a pointer to a single high-level finding of the survey in §II-B, and we do not rely any of the content of the paper on it. For completeness, we provide a short description of the survey:

Questions. We have asked participants to indicate measurement use cases, and the insights they aim to get from measurement data. We ask them what measurement types they use (control and/or data plane), what information they collect from them (e.g., latencies, BGP paths), what IMPs they use, and what is their scope (e.g., if they target small or large geographic areas, or network types such as ISPs, CDNs, etc.).

We ask them the question (to whose answers we refer in this paper) "Is there any kind of bias in the measurement data collected for this use case?", giving them three possible answers to select from "Yes / Probably yes", "No / Probably not", "I don't know".

We also ask them if they believe that there are location / network-type biases, and how useful they would find if we

could provide them with analyzes/tools that would show and mitigate bias.

Responses (relevant to this paper). We have received responses from 50 participants, both network engineers/operators (~75%) and researchers (~25%). More than ~80% of the participants said that they are experienced users of IMPs. 70% uses RIPE Atlas in their measurements, and around 50% use RIPE RIS and/or RouteViews.

26 participants (52%) replied "Yes / Probably yes" in the question about existence of bias, 14 (28%) replied "No / Probably not", and 10 (20%) replied "I don't know".

APPENDIX B

COMPARISON OF DIFFERENCE BIAS METRICS

Bias metrics. There are several metrics to quantify the difference between two distributions (i.e., the "bias" in our context).

- Kullback–Leibler (KL) divergence; see §III-A
- Total Variation (TV) distance: $B_{TV} = \sum_{i=1}^K |p_i - q_i|$
- Max distance: $B_{max} = \max_{i=1}^K |p_i - q_i|$

The main difference between KL-divergence and TV distance metrics, is that the former is more sensitive to changes in characteristics of lower probabilities p_i [55]. For example, let $P = [0.6, 0.2, 0.2]$ and two distributions $Q^A = [0.7, 0.1, 0.2]$ and $Q^B = [0.6, 0.1, 0.3]$ that differ by ± 0.1 compared to P . While for the total variation it holds that $B_{TV}(P, Q^A) = B_{TV}(P, Q^B)$, for the KL-divergence it holds $B_{KL}(P, Q^A) < B_{KL}(P, Q^B)$, because the +0.1 was at a characteristic with a lower probability in Q^B .

The main difference between the *Max* distance and the other metrics, is that the former accounts for the "worst case" (i.e., max deviation between two distributions), whereas the latter calculate distances over the entire distribution.

Bias in IMPs for each metric. In Fig. 7 we present the radar plot depicting the bias for the three bias metrics. While the actual values differ for different metrics, the qualitative findings (e.g., which infrastructure set is more biased) remain the same for the majority of dimensions.