
SAARLAND UNIVERSITY

Faculty of Mathematics and Computer Science
Department of Computer Science
Dissertation



An Empirical Evaluation of Messy BGP Data Sources

Dissertation zur Erlangung des Grades des
Doktors der Ingenieurwissenschaften (Dr.-Ing.)

der Fakultät für Mathematik und Informatik
der Universität des Saarlandes

vorgelegt von
Pascal Hennen

Saarbrücken, 2025

Date of the colloquium: 09.02.2026

Dean: Prof. Dr. Roland Speicher

Chairman of the examination board: Professor Martina Maggio, Ph.D.
Reporter: Professor Anja Feldmann, Ph.D.
Dr.-Ing. Tobias Fiebig
Professor Cristel Pelsser, Ph.D.

Scientific Assistant: Dr. rer. nat. Johannes Zirngibl

Notes on style:

As most of the work presented in this dissertation was done in collaboration with other researchers, the scientific plural "we" is used.

Saarland University
Faculty MI - Mathematics and Computer Science
Department of Computer Science
Campus - Building E1.1
66123 Saarbrücken
Germany



Eidesstattliche Versicherung

Hiermit versichere ich an Eides statt, dass ich die vorliegende Arbeit selbstständig und ohne Benutzung anderer als der angegebenen Hilfsmittel angefertigt habe. Die aus anderen Quellen oder indirekt übernommenen Daten und Konzepte sind unter Angabe der Quelle gekennzeichnet. Die Arbeit wurde bisher weder im In- noch im Ausland in gleicher oder ähnlicher Form in einem Verfahren zur Erlangung eines akademischen Grades vorgelegt.

Declaration of original authorship

I hereby declare that this dissertation is my own original work except where otherwise indicated. All data or concepts drawn directly or indirectly from other sources have been correctly acknowledged. This dissertation has not been submitted in its present or similar form to any other academic institution either in Germany or abroad for the award of any other degree.

Saarbrücken, 25.08.2025 gez. Pascal Hennen

Place, Date (Unterschrift / Signature)

Acknowledgments

Look Mama and Papa, I made it!

After many long and sleepless nights, deadlines, and conferences, I would do it all again without a second thought.

Thank you!

Abstract

The Internet is the world’s largest human-build system and as such evolved to be rather complex. Operators use the Border Gateway Protocol (BGP)—the Internet’s de-facto inter-Autonomous System (AS) routing protocol—to enable global connectivity. However, routing on the Internet is evolving. Although the specification of BGP has not changed since decades, its additions and usage patterns have. Thus, BGP has become an important topic to study for researchers. They use BGP data to, e.g., understand routing decisions, map the Internet’s topology, and improve security.

Each AS uses BGP to realize its routing policies based on the business agreements that they have with its neighboring ASes. ASes typically do not share their business agreements publicly. Yet, ASes need to see the effects of a change in their BGP configuration. Route collector projects such as Routeviews and RIPE Routing Information Service (RIS) collect BGP data from as many ASes as possible and make that data publicly available in BGP archives. In addition, data broker services provide interfaces to these BGP archives. Whereas operators use this data to optimize their networks, researchers frequently use this data to study and understand the routing ecosystem.

Until now the consistency and reliability of these data sources was usually assumed to be a given. However, it is not. In this dissertation, we fill this gap by investigating the temporal consistency (are routes recorded *when* they should be) and internal consistency (are routes recorded *correctly*). Furthermore, we evaluate whether a popular BGP route collector data broker (BGPStream’s broker) reliably returns all data files according to supplied search terms.

As a policy-based protocol, BGP is implemented on the border routers of ASes. A border router maintains multiple BGP sessions and selects the best route for a prefix by evaluating all learned routes. This is done via BGP attributes. Adjusting these BGP attributes and / or filtering routes allows an AS to implement its routing policies and manage its relationships with other networks.

It is commonly assumed that ASes use the same BGP policies for all sessions with the same neighbor AS, preferring the same next-hop AS for the same prefix. In this dissertation, we show that this is often not the case—we refer to such ASes as being *heterogeneous*. We propose two inference methods to (i) quantify the number of heterogeneous ASes as observed by the route collectors, and (ii) identify ASes which explicitly diverge from the conventional BGP behavior.

Route collectors yield a public view of the Internet—they do not show privately assigned BGP attributes. Thus, ASes collaborate with each other and operate publicly accessible Looking Glass (LG). LGs are websites that allow other operators to perform queries on a subset of routers within the ASes to gather routing information.

In this dissertation, we collect a LG dataset that focuses on collecting BGP attributes from more than 149 LGs in 154 ASes from 931 routers via scraping LGs. Hereby, the difficulties relate to the non-uniformity of the LGs—most interfaces differ, the fluctuating accessibility of the LGs, as well as the different output formats. To overcome this we combined manual configuration with an automated scraping process followed by careful post-processing and manual checks.

Zusammenfassung

Das Internet ist das weltweit größte von Menschen geschaffene System und hat sich als solches zu einem komplexen Gebilde entwickelt. Betreiber nutzen das BGP—das de-facto verwendete Inter-AS-Routing-Protokoll des Internets—um globale Konnektivität zu ermöglichen. Das Routing im Internet entwickelt sich jedoch weiter. Obwohl sich die Spezifikation von BGP seit Jahrzehnten nicht geändert hat, haben sich seine Ergänzungen und Nutzungsmuster gewandelt. Daher ist BGP zu einem wichtigen Forschungsthema geworden. Forscher nutzen BGP Daten beispielsweise, um Routing-Entscheidungen zu verstehen, die Topologie des Internets zu messen, und Routing sicher zu gestalten.

Jedes AS verwendet BGP, um seine Richtlinien auf der Grundlage der Geschäftsvereinbarungen mit seinen benachbarten ASes umzusetzen. ASes geben ihre Geschäftsvereinbarungen in der Regel nicht öffentlich bekannt. Dennoch müssen ASes die Auswirkungen einer Änderung ihrer BGP-Konfiguration erkennen können. Routensammler Projekte wie Routeviews und RIPE RIS sammeln BGP Daten von so vielen ASes wie möglich und machen diese Daten in BGP Archiven öffentlich zugänglich. Darüber hinaus bieten Datenbrokerdienste Schnittstellen zu diesen BGP Archiven. Während Betreiber diese Daten zur Optimierung ihrer Netzwerke nutzen, verwenden Forscher sie häufig, um das Routing-Ökosystem zu untersuchen und zu verstehen.

Bislang wurde die Konsistenz und Zuverlässigkeit dieser Datenquellen in der Regel als gegeben vorausgesetzt. Dies ist jedoch nicht der Fall. In dieser Dissertation schließen wir diese Lücke, indem wir die zeitliche Konsistenz (werden Routen *zum richtigen Zeitpunkt* aufgezeichnet?) und die interne Konsistenz (werden Routen *korrekt* aufgezeichnet?) untersuchen. Darüber hinaus bewerten wir, ob ein beliebiger BGP Routensammler Datenbrokerdienst (BGPStream's Broker) alle Dateien zuverlässig gemäß den angegebenen Suchbefehlen zurückgibt.

Als Richtlinienbasiertes Protokoll wird BGP auf den Border-Routern von ASes eingesetzt. Ein Border-Router unterhält mehrere BGP Sitzungen und wählt die beste Route für einen IP-Präfix aus, indem er alle erlernten Routen bewertet. Dies geschieht über BGP Attribute. Durch Anpassen dieser BGP Attribute und / oder Filtern von Routen kann ein AS seine Routing Richtlinien implementieren und seine Beziehungen zu anderen Netzwerken verwalten.

Es wird allgemein angenommen, dass ASes für alle Sitzungen mit demselben Nachbarn AS dieselben BGP Richtlinien verwenden und für den selben IP-Präfix das selbe benachbarte AS bevorzugen. In dieser Dissertation zeigen wir, dass dies oft nicht der Fall ist—wir bezeichnen solche ASes als *heterogen*. Wir schlagen zwei Inferenzmethoden vor, um (i) die Anzahl heterogener ASes zu ermitteln, wie sie von den Routensammlern beobachtet werden, und (ii) ASes zu identifizieren, die explizit vom herkömmlichen BGP-Verhalten abweichen.

Routensammler bieten einen öffentlichen Einblick in das Internet—sie zeigen keine privat zugewiesenen BGP Attribute an. Daher arbeiten ASes miteinander zusammen und betreiben öffentlich zugängliche LGs. LGs sind Websites, auf denen andere Betreiber Abfragen zu einer Teilmenge von Routern innerhalb der ASes durchführen können, um Routing-Informationen zu sammeln.

In dieser Dissertation sammeln wir einen LG Datensatz, der sich auf die Erfassung von BGP Attributen aus mehr als 149 LGs in 154 ASes von 931 Routern durch das Scannen von LGs konzentriert. Die Schwierigkeiten liegen dabei in der Uneinheitlichkeit der LGs—die meisten Schnittstellen unterscheiden sich, die Zugänglichkeit der LGs schwankt und die

Ausgabeformate sind unterschiedlich. Um dies zu überwinden, haben wir die manuelle Konfiguration mit einem automatisierten Scraping-Prozess kombiniert, gefolgt von einer sorgfältigen Nachbearbeitung und manuellen Überprüfungen.

List of Publications

Parts of this dissertation are based on pre-published work. These works are co-authored with other researchers as listed below.

International Conference Publications

Poornima Mani, Q Misell, **Pascal Hennen**, Anja Feldmann, and Tobias Fiebig. "*Quis mensurat ipsos menses?*" Evaluating Data Consistency of BGP Route Collectors & Brokers". In Proceedings of the ACM Internet Measurement Conference, ACM, 2026 [Results appear in Chapter 3]

Workshops

Pascal Hennen, Cristian Munteanu, and Anja Feldmann. "BGP AS Paths: Shorter is not Always Better". In Proceedings of the ACM SIGCOMM 2025 Conference NGNO Workshop. ACM, 2025 [Results appear in Chapter 4]

Pascal Hennen, Poornima Mani, and Anja Feldmann. "Peaking Beyond the Best Route: An Extensive Dataset for Looking Glasses". NOMS 2024 IEEE Network Operations and Management Symposium QoDaNet Workshop. IEEE, 2024. [Results appear in Chapter 6]

Doctoral Symposium

Pascal Hennen, Tiago Heinrich, and Johannes Zirngibl. "Measuring Increasing Heterogeneity in BGP". NOMS 2025 IEEE Network Operations and Management Symposium. IEEE, 2025. [Results appear in Chapter 5]

Work under submission

Pascal Hennen and Johannes Zirngibl. "Heterogeneity in BGP: The Missing Hive Mind of Routers". [Results appear in Chapter 5]

Posters

Poornima Mani, **Pascal Hennen**, Tiago Heinrich, and Anja Feldmann. "5G Measurements: Revisited and Extended" 8th Network Traffic Measurement and Analysis Conference (TMA). IEEE, 2024. [Not relevant for this dissertation]

Pascal Hennen, Cristian Munteanu, and Lars Prehn. "Can We Infer Local Preferences From Passively Collected Update Streams?" 7th Network Traffic Measurement and Analysis Conference (TMA). IEEE, 2023. [Ties in with Chapter 4]

Contents

List of Publications	ix
1 Introduction	1
1.1 Scientific Approach	3
1.2 Route Collector Inconsistencies and BGP Data Distribution Software . . .	3
1.3 Measuring the Internet’s Heterogeneity	4
1.4 Extending BGP Data with Looking Glasses	4
1.5 Contributions	5
1.6 Publications and Collaborations	6
1.7 Outline	6
2 Background	8
2.1 Inter-domain Routing	8
2.1.1 Border Gateway Protocol	8
2.1.2 BGP Routing Example	9
2.1.3 Routing Information Base	10
2.1.4 AS Relationships and Policies	10
2.2 Vantage Points and Datasets	11
2.2.1 BGP Route Collectors	11
2.2.2 BGP Route Collector Data Aggregation and Brokering	11
2.2.3 BGP Routing Beacons	12
2.2.4 Looking Glasses	12
2.2.5 PeeringDB	12
2.2.6 CAIDA’s ASRank	13
2.2.7 RPKI and ASPA	13
3 Evaluating Data Consistency of BGP Route Collectors and Brokers	14
3.1 Route Collector Problem Statement	15
3.2 Methodology	16
3.2.1 Route Collector Datasets	16
3.2.2 BGP Stream Dataset	17
3.2.3 Evaluating Route Collector Data Inconsistencies	17

3.2.4	Inconsistency Categorization	18
3.2.5	Possible Inconsistency Sources	19
3.3	BGP Route Collector Inconsistencies	19
3.3.1	RIPE RIS Inconsistencies	19
3.3.2	Routeviews Inconsistencies	19
3.4	Broker Inconsistencies	22
3.5	Discussion	23
3.5.1	Related Work	23
3.5.2	Limitations	23
3.5.3	Service Operators	23
3.6	Summary	24
4	Inferring Routing Heterogeneity from AS Paths	25
4.1	Routing Heterogeneity in Research	25
4.2	Methodology	27
4.3	Datasets	29
4.4	Implementation	29
4.5	Results	30
4.5.1	Inferences Across Time	31
4.5.2	Conflicts Across Time	31
4.5.3	Impact of Beacon / Route Collector Choice	34
4.5.4	Characterization Of ASes With Conflicts	35
4.6	Limitations	36
4.7	Discussion	37
4.8	Summary	39
5	The Missing Hivemind of BGP Routers	40
5.1	Routing Heterogeneity in BGP	40
5.2	Related Work	41
5.3	Detecting Heterogeneity	42
5.4	Adoption of Heterogeneity	43
5.4.1	AS Characteristics	44
5.4.2	ASRank	45
5.4.3	Key Takeaways	46
5.5	Catching Heterogeneity	47
5.5.1	Route Collector Characteristics	47
5.5.2	Route Collector Requirements	49
5.5.3	Key Takeaways	50

5.6	Summary	51
6	Extending Existing BGP Data with Looking Glasses	52
6.1	Building a LG Dataset	53
6.2	Ethical Considerations	54
6.3	Methodology	54
6.3.1	Overview of Data Collection Process	55
6.3.2	Process Pipeline	55
6.3.3	Remarks	57
6.4	Dataset	58
6.4.1	Looking Glasses	58
6.4.2	Performed Queries	58
6.4.3	Used Prefixes	59
6.4.4	Dataset overview	59
6.5	Example Analysis	60
6.5.1	Setup	60
6.5.2	Results	61
6.6	Summary	61
7	Discussion	63
7.1	Summary	63
7.2	Conclusion and Future Directions	64
	List of Abbreviations	68
	List of Figures	70
	List of Tables	72
	Bibliography	74

Chapter 1

Introduction

The Internet is the largest human-made infrastructure [33] that brings many benefits to society, e.g., entertainment, global connectivity. Besides its operation, the Internet impacts society as a whole and is an important subject to study. One important factor for all of these benefits is the underlying infrastructure. Autonomous Systems (ASes)—networks that are operated by one operational entity—are one of the main building blocks of the Internet. As such, operators of ASes are interested to deliver reliable connectivity to their customers. One field that has vast impact on the Internet’s performance is routing via Border Gateway Protocol (BGP)—the Internet’s de-facto inter-AS routing protocol. Commonly assessed research questions about BGP are, e.g., Internet topology discovery [1, 26, 27, 69, 82, 85, 86, 127, 132], prefix hijacking detection / prevention [29, 61, 117, 118, 121], policy inference [15, 39, 68, 82, 129, 130], BGP community classification [71, 72, 123, 124]. All of these studies aim to improve our understanding about routing with BGP and the Internet.

Almost all parties in the Internet (e.g., Internet Service Provider (ISP), Content Delivery Network (CDN), and many more) perform Internet measurements. They do so to improve their own networks, to make them more cost-efficient, and to overall improve their business. Understanding Internet routing is beneficial for all these aspects. However, understanding routing decisions is not trivial. Many researchers are interested in improving our understanding of the Internet’s topology [73]. Improving our knowledge about routing helps researchers to answer many pending questions that improve the Internet as a whole, e.g., measuring Resource Public Key Infrastructure (RPKI) deployments for better security [28, 47, 59, 91], understanding AS relationships for better planning / routing [40, 49, 50, 53, 68, 82], identifying routing anomalies to debug networks [2, 67, 102, 122, 125].

ASes typically do not see the effects of changes in their routing as this information is only visible within other ASes. Yet, for an AS to check their BGP configuration they need to see the effect. The first step to improve the routing with BGP in the Internet is to observe BGP announcements. In the late 1990s / early 2000s [24, 113], two major BGP data collection projects [22, 109] started to collect BGP data. They were first deployed for operators to see how routing decisions (such as traffic engineering) influence established AS paths. The idea of route collectors is to represent snapshots of topology-wise diverse routing states

in the Internet. However, route collectors have many problems, e.g., the collected BGP data is biased [34, 120], there are not enough Vantage Points (VPs) [4], the collected data lacks in quality and must be post-processed [42]. Further, route collectors are operated on physical machines with resource constraints, i.e., bursts of announcements can lead to resource exhaustion. Thus, one underestimated and understudied problem is the consistency of the collected data. Many researchers use the assumption that the process of collecting BGP announcements using route collectors is flawless. To simplify the access of route collectors data, BGP distribution software has been developed. Furthermore, some researchers use BGP data distribution software [10, 96] that can introduce another layer of possible biases and missing data (see Chapter 3).

The available tools to fetch BGP data [10, 96] and models to explain routing in the Internet [1, 82, 92] are not perfect. To simplify (or even enable) their methodologies, researchers often impose assumptions. One such common assumption is that ASes use the same BGP policies for all sessions with the same neighbor AS, preferring the same next-hop AS for the same prefix. This is often not the case. Border routers in an AS run the BGP decision process on their own—they can choose different paths for the same prefix. We find that in almost 10% of observed ASes different paths are used for the same prefix at the same time. We refer to such ASes as being heterogeneous (non-homogenous). Another common assumption is that a shorter path is always better. We identify that more than 12% of observed ASes choose a longer path as their preferred “best” path, calling into question the reliability of these assumptions. To gain a better understanding of the Internet and BGP, it is important to shed some light on the implications of this assumption (Chapter 4 and Chapter 5).

ASes collaborate and operate publicly accessible Looking Glasses (LGs) to further increase the visibility of effects of routing changes. LGs are websites that allow the users to query one or several routers within the AS for routing information. This information may be restricted to BGP routes (routing prefix plus AS path) only or other BGP attributes as well, e.g., Local Preference, Multi-Exit Discriminator (MED), and BGP communities. Such LGs data is required by many BGP topology inference methods either as input or for validation. Hereby, one understudied field of possible data sources are LGs. There are many studies that used [1, 45, 50, 68, 69, 77] LGs, build frameworks around them [48], or built datasets [51]. However, these studies are either old, incomplete, or not extensive enough. We tackle this problem by creating, to the best of our knowledge, the largest LG dataset as of August 2025 (Chapter 6).

The goal of this dissertation is to shed light on the datasets used in understanding Internet topologies in literature, identifying limitations, and improving the datasets available for future empirical work. The main research question in this dissertation hence is:

To what extend are currently used BGP datasets suited for, e.g., Internet topology research: which limitations do they have, how have these limitations impacted prior work, and how can these datasets be improved to form a more robust foundation of empirical work in Internet topology research?

Answering this question is challenging, as artifacts and exact datasets are regularly not made available with published papers [50, 69, 77], research community-wide assumptions may not map to practice in networking [41], and validating new datasets is difficult without ground-truth information. For example, researchers use the available data from route collectors [46, 82, 89], sometimes without assessing its consistency or even quality [29, 123, 127].

To approach this main research question in a structured way, we split it into three individual sub-questions. Starting from the status-quo and existing datasets, our first sub-question is:

(1) How consistent is the currently available BGP data from the major route collector projects and which biases are introduced by BGP data distribution software?

On top of using possibly inconsistent and biased data, researchers often impose the assumption that each AS in the Internet is a homogeneous network [28, 47, 91]. They sometimes do so without acknowledging this assumption as a limitation [59, 114]. This circumstance motivates our second sub-questions:

(2) Can we use the available data to quantify the number of ASes in the Internet that do act in a non-homogeneous (heterogeneous) manner?

And lastly, although there are other sources for BGP data, they are mostly old [90], lack in quantity [97]¹, or are in their early stages of adoption [4]. Our last sub-question that we answer in this dissertation is:

(3) Can we use Looking Glasses to extend the already existing BGP data?

We use the remainder of this chapter to expand upon these three questions in Section 1.2, Section 1.3, and Section 1.4. Afterwards, we introduce the contributions of this dissertation in Section 1.5 and outline its content in Section 1.7.

1.1 Scientific Approach

The methodology that we follow in this dissertation is to first have a look at the commonly used data sources for BGP research, namely the route collectors [22, 109]. We perform an empirical study to assess this data’s quality and pinpoint inherent problems (sub-question 1). Our next step is to use route collector data to investigate the increasing routing heterogeneity in the Internet (sub-question 2). While looking into the first two sub-questions, we find that there are limitations that need to be tackled. We extend the available data sources by investigating LGs (sub-question 3). Hereby, we longitudinally collected LG data to build the currently (August 2025) largest openly accessible dataset². Finally, we tackle the main research question of this dissertation.

1.2 Route Collector Inconsistencies and BGP Data Distribution Software

Although researchers can ingest traffic, e.g., pings or traceroutes, into the Internet [44, 100], some research in BGP relies on available data sources [40, 82, 119]. This implies the importance of the quality behind existing datasets—especially of RIPE’s Routing Information Service (RIS) and Routeviews route collector projects [22, 109]. There are

¹We note here that the Packet Clearing House (PCH) is collecting BGP data from many ASes but only few are full-feeders.

²To the best of our knowledge.

fixed steps that researchers have to do before route collector data can be used: (i) researchers must check if all data is present, (ii) the data must be cleaned, e.g., AS paths with routing loops are removed, and, (iii) project dependent filtering is applied. However, these steps only tackle the immediately visible artifacts in singular snapshots. But there are more problems that can only be uncovered when the dimension of time is taken into consideration. Flavel et al. [42] tackled this problem by designing CleanBGP—a tool that improves the quality and trustworthiness of route collector snapshots. A further level of inconsistencies is given by BGP data distribution software, such as CAIDA’s BGPStream [96] and BGPKIT [10].

In Chapter 3, we continue on a similar but extended path such as CleanBGP. Whereas Flavel et al. introduce an at the time novel way of cleaning the existing BGP data we take a deep-dive into the causes of inconsistencies. We use data from all route collectors and apply the recorded updates to their provided Routing Information Bases (RIBs)—snapshots of a route collector’s routing state. We generate one RIB per route collector using its data. This generated routing state is expected to be equal to the next following RIB snapshot of the same route collector. We use this methodology to evaluate inconsistencies in the route collector’s data over time. Further, we quantify the amount of data that is missing in CAIDA’s BGPStream, or more precisely, their broker service.

1.3 Measuring the Internet’s Heterogeneity

The Internet is an enormous system that offers endless room for customization and optimization. It changes and evolves on a constant basis—this is especially the case for routing with BGP. Whereas routing policies were simple decades ago [94] the Internet of today is more complex [6, 12, 77, 131]. However, models of the Internet still use a simplified view. Researchers often assume that ASes in the Internet act as atomic nodes in a graph. This is not the case in the Internet of today [30, 92]. Nevertheless, this assumption still heavily influences Internet topology studies [1, 69, 127], AS relationship inferences [40, 53, 68], RPKI deployment measurements [28, 47], and many more.

In Chapter 4 and Chapter 5, we introduce the notion of heterogeneity and challenge the assumption that ASes are homogeneous networks. Hereby, we develop two methods that identify ASes that act heterogeneously. This is the first step of improving BGP studies across all fields. The first method focuses on the temporal changes of AS paths for routing beacon prefixes in the Internet. We use the BGP decision process to infer information about certain decision criteria. With this information, we search for conflicting results among the same AS—a clear indication of the usage of diverse routes for the same prefix in this AS. The second method uses RIB snapshots and compares routing states across various locations inside the Internet’s topology. Here, we are not limited to a set of prefixes and can find more heterogeneous ASes.

1.4 Extending BGP Data with Looking Glasses

The main problem behind BGP data sources is the sparsity of available VPs and the "public view" of update messages. BGP, as observed by route collectors, only shows final decisions for routes—private configurations are not shared. Route collector data does not show the private BGP attributes (for example Local Preferences) that operators assign to routes. However, some ASes host LGs which allow for queries to get more information

from BGP routers. LGs are well-known by the research community and used in many studies [1, 50, 77]. Although there are existing frameworks such as Periscope [48], there are many challenges in facilitating LG data access.

In Chapter 6, we introduce the largest known publicly available LG dataset as of August 2025. We use various lists of LGs [9, 64, 64] to compile a list of responsive, query-able, and automatable servers. We further automate the query process and compile a dataset that spans over 6 months of daily queries. In addition, we not only collect the raw output but prepare regular expressions that parse this data into the CSV format. The raw and parsed data is available through our own hosted website [56].

1.5 Contributions

The goal of this dissertation is to investigate decade-old assumptions that are used in Internet measurement studies about BGP. We make the following contributions towards this goal:

Analyzing Route Collector Inconsistencies: First, we raise alarm on the prevalence of temporal and internal inconsistencies in BGP route collector data and their possible impact on research results. Second, we evaluate the reliability of a popular route collector data broker service, finding inconsistencies and missing data. Third, we identify a possible root-cause, i.e., resource exhaustion for data inconsistencies in Routeviews. Last, we communicated with the operators of these services to mitigate such problems in the future.

Investigating Routing Heterogeneity in the Internet (1): We propose a methodology to identify ASes that prefer longer AS paths over shorter AS paths. We argue that such a decision is based on a higher Local Preference for the former route, which can be used to implement heterogeneity. We propose a methodology to identify heterogeneous ASes. In effect, we infer the preferred next-AS-hop as seen by different VPs and search for diversity. We apply our methodology to routes for the RIPE RIS routing beacons [111] as observed at the RIPE RIS [109] and Routeviews [22] BGP route collectors from January 2018 until December 2023. We find that RIPE, followed by ARIN and AFRINIC are the regions with the most diverse ASes.

Investigating Routing Heterogeneity in the Internet (2): We apply a methodology to identify pairs of ASes that show detours in their AS paths for the same prefix at different VPs. We find that 3855 ASes (5% of the assigned Autonomous System Numbers (ASNs) as of April 2025) show heterogeneous behavior, which has increased over the last decade. We analyze the characteristics of heterogeneous ASes, e.g., their region, business sector, rank, and Customer Cone (CC) size [18]. Although it is expected that mostly ISPs and Network Service Providers (NSPs) are heterogeneous, we find that CDNs also show heterogeneity. The increase of heterogeneity is independent of the growth of route collector projects, but is also visible based on stable BGP sessions over the last decade.

Extending Existing BGP Data with Looking Glasses: We overcome the above-mentioned restrictions and gather an extensive LG dataset. Indeed, our dataset is not restricted to a one-time snapshot but rather offers a longitudinal view as we extend our dataset every four hours with another sample. Our dataset covers 169 LGs, 175 ASes, and 604 routers—it contains over 57 M BGP routes. To do this we developed an automated scraper that emulates a browser and queries the LGs. The returned routes are then parsed to extract all available BGP attributes. Lastly, we augment our dataset with meta-information such as the router’s IP address.

1.6 Publications and Collaborations

Parts of this thesis are based on previously published work conducted in collaboration with other researchers. Below, we outline the main contributions of the thesis author to the research incorporated in this document:

Chapter 3: This chapter analyzes the quality of route collector data from various route collector projects, namely, RIPE’s RIS and Routeviews. The work is currently under submission and is co-authored with Poornima Mani, Q Misell, Tobias Fiebig, and Anja Feldmann. The authors main contributions include: *(i)* Data collection and cleaning, *(ii)* joint background research, and *(iii)* data analysis and edge-case explanation.

Chapter 4 and Chapter 5: These chapters investigate the number of heterogeneous ASes in the Internet using two different methodologies. The work from Chapter 4 is published at SIGCOMM’25 in the NGNO workshop with the authors listed in [58]. The work from Chapter 5 is currently under submission. The authors main contributions for both works include: *(i)* Formulating the research questions, *(ii)* data collection and cleaning, *(iii)* background research, and *(iv)* data analysis.

Chapter 6: This chapter introduces a longitudinal dataset for LGs. The work is published in NOMS’25 in the QuDaNet workshop with the authors listed in [57]. The authors main contributions for both works include: *(i)* Data collection and cleaning and *(ii)* data analysis.

No AI tools were used for the analysis or writing of this thesis or publications. All research, data processing, and text were carried out manually by the author and collaborators.

1.7 Outline

We align the structure of this dissertation with the above-mentioned research questions and contributions:

Chapter 2 introduces the required background information for this dissertation. In this chapter, BGP, route collectors, and many other basics are covered.

Chapter 3 takes a closer look at route collectors, their inherent inconsistencies, and flawed BGP distribution software. There, we present a deep-dive into different types of inconsistencies and possible reasons behind them. This work is currently under submission.

Chapter 4 investigates routing heterogeneity in the Internet. We hereby introduce a methodology that uses BGP update messages and RIPE RIS routing beacons to find out which ASes prefer longer over shorter AS paths. This work was published in [58].

Chapter 5 introduces another method of finding heterogeneous ASes in the Internet. This methodology focuses on BGP RIB snapshots and works in a more generalizable manner. This work is currently under submission.

Chapter 6 talks about our own LG dataset that was published in [57] and is available in [56]. We compile, automate, and scrape LG servers across the Internet to create a longitudinal dataset that introduces new BGP data sources.

Chapter 7 concludes this dissertation with a discussion about our findings and impact on the research landscape. This chapter also introduces future work.

Chapter 2

Background

This chapter covers the required background information for this dissertation. Although the following content presents an overview, there are many subtle nuances for each work that are introduced where needed.

We start with a summary of Inter-domain routing in Section 2.1. Afterwards, we introduce the used VPs and datasets in Section 2.2.

2.1 Inter-domain Routing

The Internet is a collection of different networks and other hardware that is governed by multiple entities. This section introduces routing in the Internet, specifics that are needed for the content of this dissertation, and AS relationships / routing policies.

2.1.1 Border Gateway Protocol

The Internet is a network of networks known as ASes. An AS is a collection of IP networks and routers operated by a single administrative organization that presents a unified routing policy to the rest of the Internet. BGP is a policy-based routing protocol, used to exchange reachability information between ASes on a per-network prefix basis. To scale, BGP propagates only the best route for each prefix³.

ASes establish BGP sessions with neighboring ASes to exchange routing and reachability information. Although a single BGP session allows an AS to participate in BGP routing with another AS, in practice, an AS needs at least two BGP sessions with different upstream providers or peers to ensure redundancy and reliability.

In Figure 2.1, AS *X* announces prefix *p*, which is propagated through the network. AS *C* receives two distinct paths to reach prefix *p*: *DX* and *EX*. Due to its routing policy—such as hot-potato routing—AS *C* forwards both paths to AS *A*. AS *A* has multiple BGP sessions between a single border router and several border routers of AS *C*. The

³Although M-BGP [77, 78] can propagate multiple routes, it has not yet been extensively deployed.

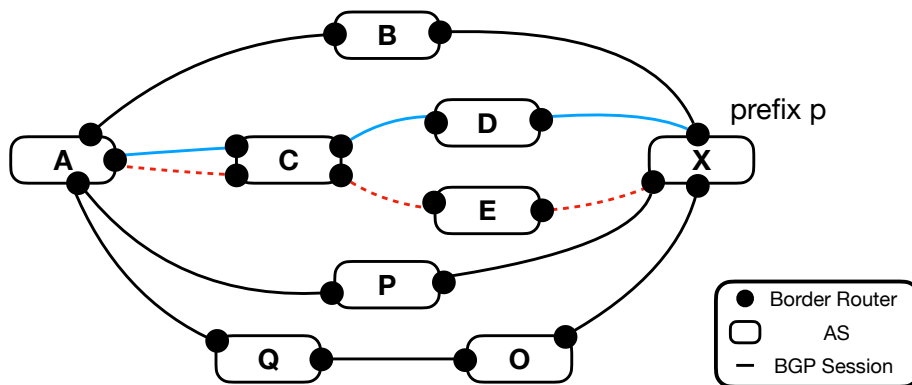


Figure 2.1: A topology that shows multiple paths in between AS *A* and AS *X*. Each line represents a BGP session between two border routers. The blue (solid) and the red (dashed) line between AS *A* and AS *C* show different routes learned by AS *A*'s router for prefix *p*.

border router in AS *A* learns two routes for prefix *p*: the route *CDX* through the session indicated in blue (solid line) and the route *CEX* through the session shown in red (dashed line). AS *A*'s router will then select the best path between the two options.

The BGP decision process, upon receiving newly advertised routes, decides on the most preferred (“best”) path for a given prefix. It consists of a sequence of elimination steps [32, 66, 103] that are based on various BGP (and vendor specific) attributes. The Local Preference, for example, is a numeric value that locally ranks routes. It is one of the first BGP route selection criteria. The next BGP attribute is the AS path of a route, here, shorter AS paths are preferred by default. The following steps include the MED which can be used to rank routes from the same neighboring AS, the Interior Gateway Protocol (IGP) cost which is used to implement hot-potato routing, and a final tiebreaker that decides for the best route by using, for example, the neighboring router’s IP address. This list of filtering steps is not exhaustive.

A router processes BGP route advertisements in three different steps. The first step is to apply import policies—which routes should be filtered and eliminated from consideration. This step also adds / removes / modifies BGP attributes. In the second step, the router applies the BGP decision process. The last step applies export policies to decide which neighboring AS receives which advertisements.

2.1.2 BGP Routing Example

Routing updates usually cause further BGP announcements to be sent, see, e.g., the scenario shown in Figure 2.1. This scenario includes many ASes (AS *A*-AS *E*, AS *O*-AS *Q*, and AS *X*) and many different BGP sessions between these ASes. AS *X* sends a BGP announcement for prefix *p* to AS *B* and AS *O* with the AS path AS *X*. AS *B* sends a follow-up BGP announcement to AS *A* with the AS path *BX* upon receiving the BGP announcement sent by AS *X*. AS *A* forwards this announcement to AS *P*. In this example, AS *O* also receives the update but has to wait until the Min-Route Advertisement Interval⁴ timer expires before forwarding it to AS *Q*. Once AS *A* receives the update with the AS path *QOX*, it realizes that this is a preferred route, updates its routing table, and sends an updated BGP announcement for the prefix *p* to AS *P*. Note

⁴A timer that prevents routers from advertising too often.

that a single announcement that was originated by a single AS—in this case AS X —can cause a sequence of BGP announcements. We refer to such a sequence of announcements as an update sequence. BGP has converged for a prefix if the routing table entries of all involved routers stop changing. This happens within 30 minutes of a prefix being announced for the first time [54, 75].

2.1.3 Routing Information Base

Each BGP router of an AS stores its known routes to all prefixes in its RIB. There are three different layers of the RIB to consider: (i) the RIB-in, (ii) local-RIB, and (iii) the RIB-out. The RIB-in collects routing updates from all established BGP sessions. Ingress filtering is performed on the content inside the RIB-in, e.g., RPKI invalid routes from further consideration, the BGP best path decision process is applied. The local-RIB then stores all these leftover routes in local storage. Finally, egress filters are applied on the content of the local-RIB—the RIB-out contains the results. Egress filters, for example, decide which routes to share over which BGP session. Each BGP router has a BGP configuration and RIB structure of its own.

BGP routers learn routes for prefixes over previously established BGP sessions. These routers then distribute the learned routes via internal Border Gateway Protocol (iBGP) to all other routers in the same AS. Routers inside an AS usually run an IGP to learn the AS's internal network topology. Then, the BGP and IGP information is used to build a Forwarding Information Base (FIB) for each router. The FIB is used to route traffic towards prefixes—the prefixes must not reside inside the AS.

2.1.4 AS Relationships and Policies

BGP, aside from enabling Internet-wide routing, is a policy-centric protocol. This section introduces AS relationships [13] as well as routing policies.

The different relationships in-between ASes are classified in five types [7]: (i) Provider, (ii) Customer, (iii) Route Server (RS), or (iv) RS-Client, and (v) Peer. Research usually uses adjacent classifications, namely: (i) Customer to Provider (c2p), (ii) Peer to Peer (p2p), (iii) Sibling relationship, and (iv) more complex relationships. A customer AS, for example, does not announce prefixes learned from their provider (the customer has to pay for traffic) to their peers (usually no payment for traffic) since that introduces unnecessary cost.

There are public and private BGP attributes. Public attributes are redistributed outside an AS, whereas private attributes are not distributed, that is, they are removed from the route before announcing it externally. Some attributes are chosen by the administrators and assigned via the local BGP router configuration, i.e., the Local Preference value; other attributes are properties of the route, i.e., the AS path. The values of these attributes influence the best path selection of an AS or can be used to filter which routes to accept or announce to which neighbor AS. Thus, ASes can steer how they receive which traffic by choosing appropriate values and filter policies for each prefix—so-called routing policies. Choosing appropriate routing policies is called traffic engineering.

An AS uses routing policies to determine where traffic enters / leaves its network or which traffic should not pass through its network. It implements routing policies by modifying, adding, or removing route attributes and / or filtering routes. Example

routing policies range from hot-potato routing to those necessary to realize AS peering / customer-provider relationships with other ASes.

2.2 Vantage Points and Datasets

The routing environment of the Internet is observed by different systems at topology-wise distributed locations. In this section, we introduce the systems that are used in this dissertation as data sources.

2.2.1 BGP Route Collectors

While researchers and professionals often speak of ‘the’ Global Routing Table (GRT), this is a simplification to easy discussions about ‘the GRT’. In fact, the GRT is the set of all unique routing policies across all BGP routers of all ASes [62]. Each individual AS – and despite differing guidance [55] – each BGP router in an AS might have a different view on the GRT, i.e., sees one perspective of all possible perspectives.

Route collectors have been introduced to allow for building a global aggregation / perspective of different routing tables [105]. To this end, they maintain sessions with a large number of BGP speakers from different ASes [21, 108], and receive (usually) full table feeds from these BGP speakers. ASes voluntarily peer with route collectors and exchange BGP information. They in turn provide interfaces—e.g., RIPE’s BGPPlay [110]—or / and regular snapshots of the collected data in Multi-threaded Routing Toolkit (MRT) file format [11]. As the export of full snapshots is resource intensive, most route collectors currently export their full routing table only in regular, e.g., hourly, intervals, and provide a list of the updates they see in the meantime in, e.g., 5 minute (RIPE RIS), or 15 minute (Routeviews) snapshots [80].

While, by now, several communities and groups operate route collectors, e.g., NLNOG and BGP.tools, the most well-known and frequently used, at least by research projects, are Routeviews [22] and RIPE RIS [109]. Routeviews is a project from the University of Oregon which started in the early 2000s [24]. RIPE RIS is being operated by the RIPE NCC for the RIPE community and the Internet at large since 1999⁵.

There are other BGP data collection services in the Internet as well. One such system is Gill [4]—a fairly new but promising addition. Nevertheless, the Internet-wide adoption of this system is not advanced enough yet to benefit our research. Other systems are the Packet Clearing House (PCH) [97] which feeders are usually non-full-feeders and the now deprecated Isolario project [90].

2.2.2 BGP Route Collector Data Aggregation and Brokering

A raw trace from an route collector can range upwards of 100 MB per hour (1 RIB + updates), i.e., 2-3 GB per day and more. With 71 route collectors currently active, this means terrabyte of data per month. For some research questions, not all of this data is necessary. Instead, only specific time-frames, or specific parts of the data are needed, e.g., information on which ASes announced which prefixes.

⁵<https://data.ris.ripe.net/rrc00/>

To help researchers navigate the huge datasets, several tools and frameworks have been developed to ease access by, e.g., identifying the links to the appropriate RIB and update files (broker services) based on a simple query, or by providing access to pre-processed / pre-aggregated data views. One popular aggregation is the CAIDA AS2Prefix dataset [23], which aggregates the Routeviews dataset to a list of prefixes with all ASes that announce them. A popular service for brokering is CAIDA's BGPStream broker [17], which allows access of the route collector data while applying search filters. Such filters include, e.g., timeframe, route collector projects (e.g., RIS or Routeviews only), specific route collectors (e.g., only those at Internet Exchanges). This broker does not serve BGP data in form of MRT files, but rather points to snapshot files that can be downloaded⁶.

2.2.3 BGP Routing Beacons

BGP routing beacons are a set of IPv4 as well as IPv6 prefixes that have certain criteria. These prefixes are announced from various locations within the Internet ecosystem. In addition, we know exactly when and where they are announced. This dissertation focuses on the RIPE RIS routing beacons [111]. They are a well-maintained RIPE RIS service and constantly include 16 IPv4 and 19 IPv6 prefixes during the period our study⁷ in Chapter 4 takes place. All beacon prefixes are announced by AS 12654 from various geographical locations and to various amounts of peering ASes. Thus, each prefix has a different route. Moreover, BGP beacons have been used extensively in previous work, e.g., [43, 46, 84, 95].

All BGP beacons follow a fixed announcement and withdrawal schedule. They are announced at 00:00 UTC and withdrawn two hours later at 02:00 UTC. This cycle is repeated every four hours. Since we rely on a sequence of BGP announcements we take advantage of this fixed announcement schedule. We can thus use numerous new announcement sequences every 4 hours to perform our analysis. Since previous work has shown that BGP can take up to 30 minutes to converge [54, 75], we consider a sample of updates of 30 minutes for each newly announced prefix.

2.2.4 Looking Glasses

AS operators perform traffic engineering to optimize their networks. However, how certain changes affect the resulting route is not known before BGP has converged for this prefix. Operators thus host so-called LGs as one possible solution. LG are web interfaces that allow outsiders to fill information and perform queries on border routers. One type of query, for example, allows the user to receive RIB entries for a given prefix. Another type of query allows the user to perform ping / traceroute calls from a chosen border router. These queries are usually restricted on a small subset of border routers and in the size of their output.

2.2.5 PeeringDB

PeeringDB is a project that enables ASes to share information about themselves. It is primarily used by ASes to share their peering information, which is used by some ASes to

⁶Brokers may limit the output of a query if the content gets too large.

⁷Their beaconing system was only changed on the 1st of February 2024.

create filter lists. In addition, organizations such as RIPE, Internet Exchange Points (IXPs), and other Internet participants also voluntarily report information about themselves in this database. Although the data is self-entered, its popularity and nature of its use cases allow one to extract information about the business sector that the AS belongs to [81]. PeeringDB distinguishes between 10 different business sectors which include NSPs, ISPs, and CDNs.

2.2.6 CAIDA's ASRank

CAIDA's ASRank [18] is a public ranking of ASes in the Internet. They provide monthly snapshots of, for example, business relationships between ASes, their rankings, and other calculated metrics. Their methodologies use the RIPE RIS and Routeviews route collectors, more precisely, the snapshots of the first 5 days of each month. One benefit of this dataset is that they also classify the Regional Internet Registry (RIR) where ASes are assigned to.

2.2.7 RPKI and ASPA

BGP is a decade old protocol that was designed without security in mind. RPKI [76] is one of the main efforts to make BGP secure. The RPKI consists of Certificate Authorities (CAs) that issue certificates containing a list of IP prefixes to ISPs. These prefixes are allocated to the given ISP. ISPs then use their allocated prefixes to create Route Origin Authorization (ROA) objects that allow certain ASes to announce given prefixes. Aside from the certification, there are ROA objects that allow other ASes to validate BGP routes.

ROA objects can be retrieved by Relying Party Softwares (RPs) and used locally. This is how BGP routers receive data for Route Origin Validation (ROV) purposes. However, this is not the only use-case of ROA objects. There have been many studies that assess the Internet's RPKI coverage [79, 91], measure delay in updates to the RPKI [44], and find vulnerabilities [37]. Although RPKI is a valuable source for BGP studies in general, we could not yet use it to answer our research questions. Previous work [31, 74], for example, already assessed the quality of ROV and ROA objects (sub-question 1). Further, RPKI does only map the AS-level to prefixes which rules out studies about the Internet's heterogeneity (sub-question 2). Lastly, RPKI does, to the best of our knowledge, not have any third party sources of information that could be used to build new datasets (sub-question 3).

Another upcoming and promising security addition to BGP is Autonomous System Provider Authorization (ASPA). With ASPA objects, ASes can publish their relationships towards their peers. Such a mapping with public access benefits all kinds of AS relationship studies, including our own heterogeneity measurements. However, ASPA is at the time of writing (July 2025) still a draft and not deployed yet. For this reason, we cannot use ASPA yet and deem it as out of scope for this dissertation. Nevertheless, ASPA is promising for our future work.

Chapter 3

Evaluating Data Consistency of BGP Route Collectors and Brokers

Route collectors, aside from being important data sources for operators, are useful for studies about BGP. Hereby, researchers access the public data, clean it [42], and use it in their methodologies. However, route collector data is not perfect. There can be many inherent problems due to the public and volunteered nature of feeders. For example, feeders can be misconfigured, route collectors might experience resource issues, or outages can happen while writing snapshot files. To the best of our knowledge, these root causes were not investigated yet.

In this chapter, we tackle the first sub-question of this dissertation: how consistent is the currently available BGP data? We hereby perform an empirical analysis of the data quality of the BGP route collector projects to answer this question. It covers our study about inconsistencies in route collector data, problems in BGP distribution software, and possible root causes. The contributions of this chapter can be summarized as follows:

- Until now the consistency and reliability of route collector data sources was usually assumed to be a given. However, it is not. We fill this gap by investigating the temporal (are routes recorded *when* they should be) and internal (are routes recorded *correctly*) consistency. Furthermore, we evaluate whether a popular BGP route collector data broker (BGPStream's broker) reliably returns all data files according to supplied search terms.
- We find widespread inconsistencies for route collector data, and share our approach to enable researchers to perform the same analysis for spot-checks when working with route collector data. For BGPStream, we find data files being regularly missing since 2017, including all BGP updates for all except one RIPE RIS route collector.

3.1 Route Collector Problem Statement

Recall, the Internet consists of networks – called ASes – that exchange routes using the BGP. The BGP is a highly complex distributed system and has been the subject of much research to fully understand its structures and interactions for several decades, e.g. [43, 79, 82]. There is no central and universal truth – ‘one true’ perspective – of the global routing infrastructure. Each AS has its position in the Internet topology, its *own* perspective of global routing, and thus the GRT - the set of paths to all routes in the Internet.

Route collectors [22, 109] work to create a comprehensive global picture, using the *individual* perspectives of routers within various ASes. The two most well-known route collector projects are RIPE RIS [109] and the Routeviews project [22]. ASes export their Routing Information Base (*RIB*), typically the full routing table, to route collectors. Route collectors do not actively participate either in routing nor in traffic exchange. Their only task is to collect perspectives of the GRT, by recording (i) the stream of BGP announcements from their peers, (ii) snapshots of RIBs from peers, and (iii) state changes of their peering connections. This data is then made available to operators and researchers.

This data has been used to infer the Internet’s topology [82], investigate the Internet’s routing security posture [89], track major outages [63], or to determine biases in the route server system [120]. However, despite crucial role of the route collector infrastructure, there has – so far – been no structural investigation into the reliability and consistency of route collector data and route collector data broker services.

In this dissertation, we investigate RIPE RIS and Routeviews data from 1st of February 2018 to 29th of February 2024 for inconsistencies. We verified the completeness of the archives of both projects, and developed a methodology to identify inconsistencies in this data. We observe not only *temporal inconsistencies* (routes being attributed to incorrect time slices) but also *internal inconsistencies* (subsequent *RIBs* containing or lacking routes not supported by provided update message data). For example, we observe a steady, wide-spread, yet low rate (100s per day), of inconsistencies in Routeviews, likely caused by performance bottlenecks.

Using the raw data can be tedious, thus broker services, e.g. CAIDA’s BGPStream [17, 96], have been developed that aid access to this data. It offers an API which returns the URIs of RC data archive files from search filters. We find that the broker returns incomplete data — it does not return specific update and *RIB* files for Routeviews nor RIPE RIS, including *long time periods*, e.g. single files across years but also the whole of 2022-12 for RIPE RIS RIB snapshots. However, it does not inform about these gaps — users are not necessarily aware of missing data.

When using Routeviews and RIPE RIS we typically assume that they provide a pervasive view of BGP data [60, 119, 128]. Even when providing *all* data, the combination of Routeviews and RIPE RIS only provides coverage of 1% of ASes [4]. Missing, or worse, incorrect data thus further reduces the already limited coverage. An example for impacted fields is the detection of forged-origin hijacks, link-failure localization, and path-manipulation detection [4].

One tool that can be affected is ARTEMIS [119]. It, in the case of RCs that do not provide live feeds, uses update files collected with BGPStream to find prefix hijacks. Simulations using ARTEMIS show that using Routeviews and RIPE RIS data it is possible to detect all high-impact events, however this relies on having full and accurate data, both from

the collectors and BGPStream. The authors themselves note their method is only truly effective when *all* data is available. Similarly to ARTEMIS, the BGP data collection tool BML [60] uses BGPStream to collect and parse update files with the goal of creating a graph of the Internet. However, their methodology of applying BGP announcements to a baseline is vulnerable to temporal inconsistencies, the effects of which are non-trivial, especially since they propose a model that is used by other researchers. A recent RPKI study by Testart et al. [128] uses BGPStream to create an AS level model of the Internet. Missing *RIB* snapshots can severely impact the inferred relationships and traffic simulations, as can temporal inconsistencies impact their method for the generation of 5-minute table snapshots. These projects highlight the importance of measuring the quality of RC data and BGP data distribution software.

We further motivate this work with personal experiences. In late 2023, we conducted a study on best route selection in BGP. For that we, among other things, processed route collector data. We found major disruption in the Internet in December 2022: route selection became unstable. We saw major network fragmentation with routing islands forming; that is, a loss of overall end-to-end reachability, or so we thought.

We—excitedly, given the major effect we observed—contacted an operator and inquired how they experienced these disruptions, and whether they can provide more information on the impact of these events. To our somewhat surprise, they did not share our enthusiasm. Instead, they inquired the exact date range of events, and strongly suggested that we verify that we are actually processing *all* available files from *all* available route collectors. Indeed, using a broker we did not include most RIPE RIS data for December 2023 leading to this incorrect conclusion. While we *should* have been more diligent in checking our data, using external services can hide such errors from researchers.

Positively, all operators of route collector infrastructure (Routeviews and RIPE RIS) as well as data brokers (CAIDA) promptly responded when we notified them of these problems and immediately started to mitigate them.

3.2 Methodology

In this section, we introduce our methodology for analyzing route collector data for inconsistencies. This includes our route collector data retrieval, verification, and preparation processes, as well as the algorithm for analyzing route collector data. Naturally, this includes our approach for classifying types of inconsistencies in route collector data. Similarly, we also describe how we retrieve and analyze the BGPStream Broker data for our subsequent analysis.

3.2.1 Route Collector Datasets

Route Collector Data Archives: The first step is to ensure that that all snapshots and all updates are available. Therefore, we gathered all *RIB* exports and updates (in MRT format) for all route collectors from Routeviews (57) and RIPE RIS (25) route collectors for every February for a 6 year period (from 1st of February 2018 to 29th of February 2024) — a total of 138.28 TB. Due to storage and compute limitations, we limited our analysis time frame.

Based on our experience we strongly recommend that route collector operators provide machine-readable listings of archive contents. Archive contents, for example, can be

presented in JSON format where each entry contains meta-data and the file’s download link. Moreover, to better detect transmission errors, checksums over data files should be included.

Data Cleaning: Given that the data does not contain checksums it is not surprising that we detected multiple corrupt or invalid files. For example for route-view.linx the *RIB* files have a typical size of 145 MB. However, during the period of 16th of October 2023 to 19th of October 2023 all files had a size of 14 bytes. Similar artifacts occur at other time periods for route-view.linx, for other Routeviews route collectors, and for rrc15 in the RIPE RIS dataset between 20th of February 2019 and 22nd of February 2019. After consulting with Routeviews we consider the affected route collectors to be non-functional during these time periods.

In addition, one has to account for downtimes of route collectors. For example, routeviews.telxatl, was down between 20th of February 2020 and 5th of March 2020, leading to empty *RIBs* (down to 300 byte compared to previously 30 MB and afterwards 26 MB). For Routeviews, we manually identified all such events in our dataset, ultimately excluding 67k files.

System Logs/Session State Changes: To be able to faithfully analyze inconsistencies in *RIBs* we must be aware of changes in the BGP state machine for each of the route collectors’ peers. For RIPE RIS, BGP session state updates are included in the update messages files. However, for Routeviews such messages are absent from the update files. To detect these, we had to parse *syslog* messages from the Routeviews route collectors. These logs allow us to also analyze route collectors’ resource utilization.

Route Collector Data Processing: For processing MRT files we use CAIDA’s PyBGPStream [20] — a Python interface to the BGPReader toolchain [16]. Alternative options include, e.g. BGPdump [104], bgpscanner [99], or MRTparse [70]. As noted by Schlump et al. [115, 116], each tool has its strengths and weaknesses. Hereby, BGPReader matches our needs while used as local parser rather than to access CAIDA’s BGPStream broker.

3.2.2 BGP Stream Dataset

CAIDA’s BGPStream [96] includes a broker service that provides URIs of route collector archives matching filters. To evaluate the completeness of the returned file lists, we retrieved the URIs for our measurement period (1st of February 2018 to 29th of February 2024) from BGPStream and compare them to the raw input data.

3.2.3 Evaluating Route Collector Data Inconsistencies

The main challenge for evaluating inconsistencies in route collector’s data is an absence of ground-truth. To assess what route collectors did or did not miss, one need access to the routes exported to the route collector by its AS neighbors at the time the *RIB* exports is created. *This* is exactly the purpose of the route collector datasets and, thus, sadly the dataset we need to validate. Hence, instead, we evaluate the ‘*internal consistency*’ of the gathered data in terms of (i) temporal consistency, and (ii) eventual consistency:

Evaluating temporal consistency is a straight-forward process. Each individual message in an MRT file is time-stamped. Thus, we compare the messages in a file with the time frame covered by the file—as captured by the file name.

For eventual consistency, we rely on the intuition that given two adjacent *RIBs* and the updates in between, the second *RIB* can be derived from the first one by applying all

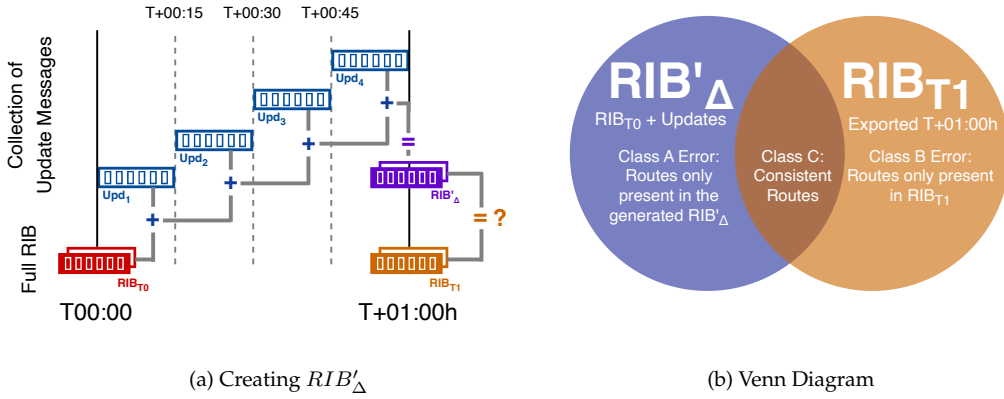


Figure 3.1: The process of combining RIB_{T_0} and update files to create RIB'_{Δ} for comparison with RIB_{T_1} . Then Venn diagram shows the expected overlaps (not to scale): Class A ($r \in RIB'_{\Delta} \setminus RIB_{T_1}$), Class B ($r \in RIB_{T_1} \setminus RIB'_{\Delta}$), or Class C ($r \in RIB'_{\Delta} \cap RIB_{T_1}$).

updates in received between the two snapshots. Hence, for each pair of adjacent $RIBs$, RIB_{T_0} and RIB_{T_1} , we apply all route updates as captured by update files from the period between the two to RIB_{T_0} , to create a derived RIB'_{Δ} . This includes (i) adding newly announced routes, (ii) updating changed routes, (iii) removing withdrawn routes, and (iv) handling (that is, removing) routes that should be cleared due to a session state changes. For RIPE RIS, this information is included in the update files. For the Routeviews dataset, this information is extracted from the log messages. Updates are applied in temporal order in order of the timestamps in the MRT files. Once we generate RIB'_{Δ} we can compare it to RIB_{T_1} . Ideally, the two should be identical, see Figure 3.1a.

3.2.4 Inconsistency Categorization

After creating RIB'_{Δ} we can compare it to RIB_{T_1} . Since $RIBs$ are sets of routes, these are set comparisons, i.e., for each route r , one of $r \in RIB'_{\Delta} \setminus RIB_{T_1}$, $r \in RIB_{T_1} \setminus RIB'_{\Delta}$, or $r \in RIB'_{\Delta} \cap RIB_{T_1}$ holds, see Figure 3.1b. For easier readability, we term these three options 'Class A (Error)', 'Class B (Error)', and 'Class C'. Ideally, all routes r' in RIB'_{Δ} are in RIB_{T_1} , and all routes r_b in RIB_{T_1} are in RIB'_{Δ} , i.e., $\forall r' \in RIB'_{\Delta} \Rightarrow r' \in RIB_{T_1} \wedge \forall r_b \in RIB_{T_1} \Rightarrow r_b \in RIB'_{\Delta}$. That is, there is no inconsistency, $RIB'_{\Delta} = RIB_{T_1}$. Note, for this dissertation we only check for differences in prefixes and AS_{path} . We do not include other attributes, e.g. BGP (large) communities. That is, we key the set by (prefix, peer, AS_{path}).

Class A ($r \in RIB'_{\Delta} \setminus RIB_{T_1}$): This includes routes that are only present in RIB'_{Δ} but not in RIB or those where the AS_{path} differs. This can, e.g., occur if RIB_{T_0} has not processed a withdrawal that occurred before the start of the first subsequent update file, or if a route is announced after the export for RIB_{T_1} started and is, thus, still included in the update file. Alternatively, this occurs when update messages are lost.

Class B ($r \in RIB_{T_1} \setminus RIB'_{\Delta}$): Routes that are in RIB_{T_1} but not in RIB'_{Δ} . Here, we have inverse of Class A Errors, a route is announced after the export for RIB_{T_0} started, but still recorded in the update file prior to RIB_{T_0} or a route still contained in RIB_{T_1} but not recorded in the update file prior to RIB_{T_1} . Just like before another reason may be lost update messages.

Class C ($r \in RIB'_{\Delta} \wedge RIB_{T1}$): Routes present in RIB_{T1} and RIB'_{Δ} .

3.2.5 Possible Inconsistency Sources

There are multiple places where inconsistencies can occur when accessing BGP data: (i) The route collector might not record a route correctly in an update and / or RIB export, (ii) An aggregation service may not accurately process all files from an route collector project, and, (iii) A broker may not provide all available data matching a search query. All can significantly alter the results. However, especially the first two are difficult to detect, as they require processing the raw data.

3.3 BGP Route Collector Inconsistencies

In this section, we analyze inconsistencies in both RIPE RIS and Routeviews.

3.3.1 RIPE RIS Inconsistencies

For RIPE RIS, we found 14.604.493 inconsistencies between 1st of February 2018 and 29th of February 2024, several orders of magnitude less than for Routeviews, see below. All of these were Class A inconsistencies. The only notable exception was `rrc12`, which shows around 15,000 Class A inconsistencies per RIB pair on average between 1st of February 2021 and 2nd of April 2023, when the inconsistencies stopped. All routes with inconsistencies were learned from one neighbor, `80.81.195.243` with `AS48646`, which is listed as announcing a single IPv6 prefix in the RIPE RIS peer list [108]. After contacting the RIPE RIS operators they confirmed that a misconfiguration involving this peer led to incorrect data, and subsequently removed this peer's routes from the database. Due to the comparatively small number of inconsistencies with RIPE RIS we instead focus on the Routeviews dataset.

3.3.2 Routeviews Inconsistencies

Temporal Consistency Errors: Update files should only contain the updates received in the corresponding time period; that is, an update file starting at 08:00 should only include routes up to 08:15 for Routeviews and up to 08:05 for RIPE RIS. However, for Routeviews, it was common to find BGP announcements with timestamps beyond 08:15 or even before 08:00. In fact, we observed instances where routes from earlier updates ended up in the later files and vice versa, see Figure 3.2. Note, we did not observe any misplaced routes for the RIPE RIS route collectors.

We note, that the number of misplaced updates per RIB for individual route collectors is relatively stable. We, thus, conjecture these misplaced updates are an effect of resource starvation given a non-atomic export process. Routeviews confirmed that these issues are known to them and the result of non-atomic exports.

For our subsequent analysis, we account for such temporal inconsistencies by including data from an additional snapshot before and after the RIB snapshots we focus on. We furthermore check for each update that its timestamp falls into the appropriate time window. We first quantify the number of misplaced updates and then measure the impact by reordering them.

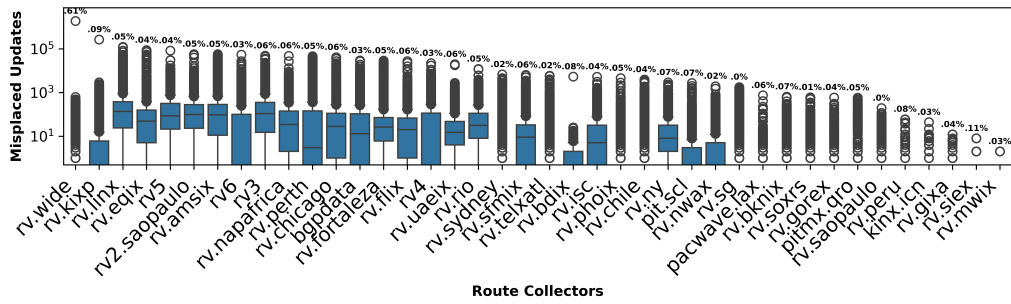
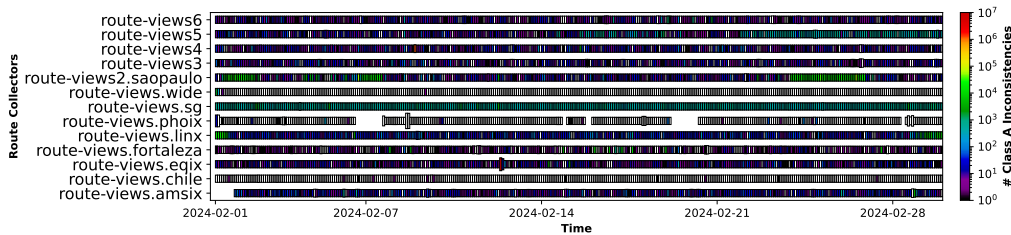
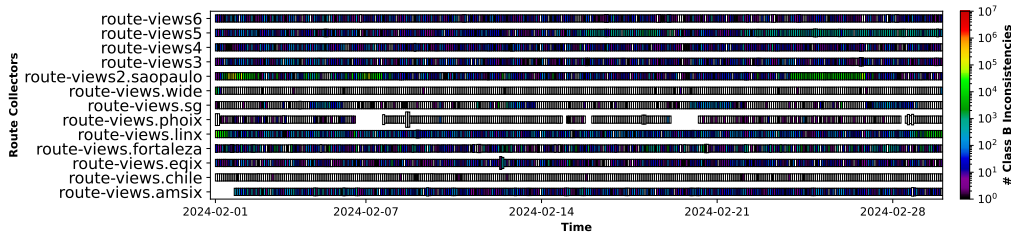


Figure 3.2: Number of misplaced route update messages along with a percentage of all update messages for Routeviews.



(a) Class A



(b) Class B

Figure 3.3: Class A and B inconsistencies for route collectors from Routeviews in February 2024. Box color signifies the number of inconsistencies, whilst box size signifies the change in RIB size from the previous RIB.

Class A Inconsistencies: We show the number of Class A inconsistencies per RIB between 1st of February 2024 and 29th of February 2024 for all route collectors from Routeviews in Figure 3.3a. The largest spike of inconsistencies happens around the 12th for `route-views.eqix`, with a corresponding increase in RIB size, indicating that this route collector became overloaded and was unable to process all updates. A similar, albeit smaller, spike happens on `route-views4` around the 10th. `route-views2.saopaulo` shows several long periods of increased inconsistencies, without a corresponding change in RIB size. Similarly `route-views.linx` shows an elevated level of inconsistencies at the start and end of the month. `route-views.sg` shows a much increased baseline of inconsistencies compared to other route collectors. We however do *not* see spikes of inconsistencies with a temporal correlation across different route collectors. This implies a root-cause of these inconsistencies local to individual route collectors.

Further supporting this, we take a deeper look at `routeviews4`. Here, *RIB* files regularly became smaller whilst update files became larger. Figure 3.4 shows an example

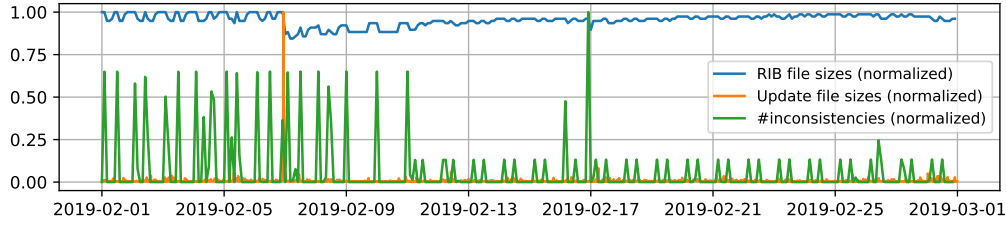
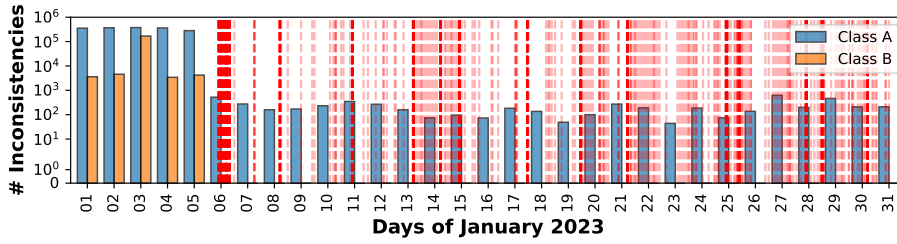
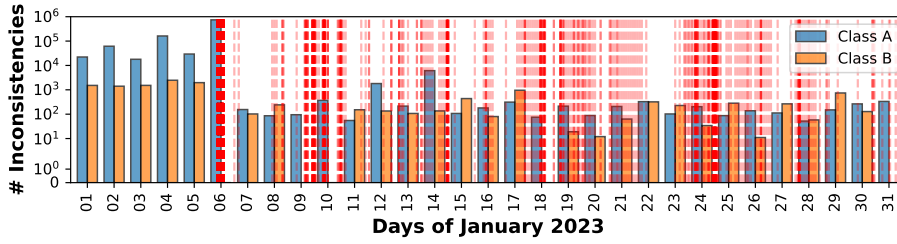


Figure 3.4: An example of how file sizes are an indicator for resource problems thus inconsistencies, using `route-views4`.



(a) `route-views3`



(b) `route-views.amsix`

Figure 3.5: Inconsistencies vs. resource starvation syslog messages for `route-views3` and `route-views.amsix`.

for February 2019. Each dip in *RIB* size corresponds to a spike in update sizes and inconsistencies. Routeviews stated that this is because of resource exhaustion.

There are several route collectors without a baseline level of inconsistencies, nor spikes in inconsistencies with increased *RIB* size, indicating that these route collectors have appropriate resources to deal with incoming traffic volumes.

Class B Inconsistencies: For Class-B inconsistencies (see Figure 3.3b) we find a similar patterns: a steady baseline of about 100 inconsistencies per day per individual route collector with some, yet less pronounced, spikes. The observed increases in Class A and Class B inconsistencies for `route-views2.saopalo`, `route-views5`, and `route-views.linx` overlap. This further indicates that these inconsistencies are caused by resource exhaustion on the Routeviews route collectors. The high baseline in Class A inconsistencies seen on `route-views.sg` is not represented in the Class B inconsistencies, indicating the problems with this route collector are not the result of persistent resource exhaustion.

The Role of Resource Exhaustion: Next, we further investigate the relationship between resource exhaustion and number of inconsistencies. Specifically, we look at the system

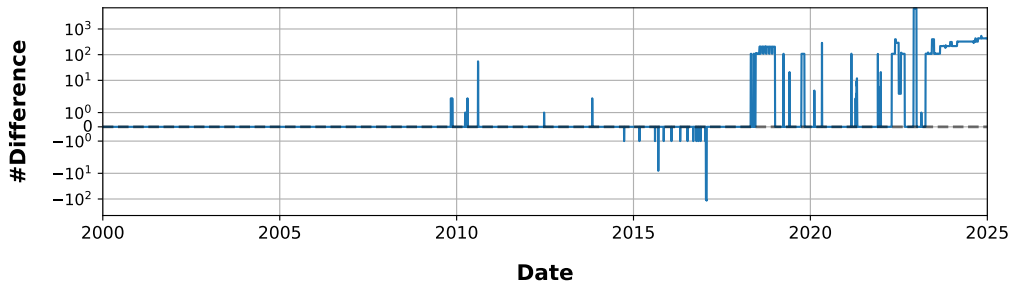


Figure 3.6: Number of URIs advertised by BGPStream vs. those directly retrieved from the archives.

messages from the syslogs that indicate resource exhaustion during periods of increased Class A / Class B inconsistencies. Neither for `route-views3` (c.f. Figure 3.5a) nor `route-views.amsix` (c.f. Figure 3.5b) did we find *any* logs indicating any overload until *after* the exhaustion event. Yet, both route collectors log events indicating thread starvation for FRR – that is, insufficient CPU resources – on 5th of January 2023 (97 on `route-views.amsix` and 48 on `route-views3`) and 6th of January 2023 (245 on `route-views.amsix` and 215 on `route-views3`). This is consistent with a high CPU utilization during the days even when no log messages were generated – or rather, could not be written due to no CPU resources being available – leading to an increase in inconsistencies, as non-atomic exports were delayed, causing accumulated inconsistencies. The synchronized nature in which the effect disappeared again hints at, either, a single peer connected to both route collectors being responsible – e.g. by sending a large number of routes – or a monitoring, management, or backup job / misconfiguration causing these effects.

3.4 Broker Inconsistencies

Next, we focus on inconsistencies in the collected URIs as advertised by BGPStream’s broker vs. those directly retrieved from the archives., see Figure 3.6.

Before 2016, we observe few difference between the broker and the archives. From 2017 onward, discrepancies begin to emerge which aligns with practical realities of academic work, e.g. [41]. Between the start of 2017 and the start of 2018, discrepancies increase as new route collectors are not added to the broker. At the beginning of 2018 these route collectors were added to the broker; however, the old URIs in the broker’s database were apparently updated, as we can still observe the discrepancies to this day.

Between 3rd of December 2022 and 31st of December 2022 only data from Routeviews, and one RIPE RIS route collector (`rrc15`) is available. All other RIPE RIS route collectors are missing. During this time, RIPE RIS changed their web archive format and migrated their web server from Apache to NGINX [93]. The route collector `rrc15` was updated at a later date. The transition likely disrupted the way the Broker collects its information. This took some time to notice and update. During this time the data was not properly processed. Even though this issue was resolved at the end of December 2023, similarly to 2017, BGPStream did not backfill data for December 2023.

At the time of our study, we find four route collectors from Routeviews to be missing (`pacwave.lax`, `kinx.icn`, `pitmx.qro`, and `pit.scl`). These are relatively new route

collectors (setup between 2023 and 2024). We conjecture that this adding route collectors is a manual process which requires maintainer intervention [3].

We contacted CAIDA regarding the observed inconsistencies (missing files and missing route collectors) in the BGPStream’s broker. They promptly responded and are now working on resolving this issue.

3.5 Discussion

3.5.1 Related Work

Related work on not *using* rather than *evaluating* the consistency and reliability of RC data is sparse. Flavel et al. [42] are among the first to examine measurement artifacts in BGP data. However, their focus is on pre-processing and data cleaning. The BGPStream authors themselves [96] discuss some of the kinds of errors we discuss, but do not evaluate their prevalence. Work by Sermpezis et al. [120] evaluates the impact of vantage point selection rather than the consistency of gathered data. Finally, Cittadini et al. [34] identify biases in public BGP datasets and analyze the impact of collector placement. They evaluate the total number of topological position and prefix visibility from collector peers. We could not find a structured evaluation of inconsistencies *within* RC data — it had been missing prior to our work.

3.5.2 Limitations

Our work has several practical limitations: (i) Our evaluation of RC data is limited to one month of each 12 across 6 years due to the major compute requirements of working with such data. Storage was also a limiting factor in this decision, but not as much as the compute requirements. Nevertheless, even this narrow window demonstrates the prevalence and unpredictability of data inconsistencies. (ii) Due to scope constraints, we only evaluate one broker service. Other broker services may have more and/or less inconsistency. (iii) Performing replication studies of related work using cleaned RC or fixed broker data to evaluate the impact of data inconsistencies is a long-term effort beyond the scope of this paper.

These limitations do not impact our core-conclusion — critical inconsistencies exist in BGP RC data and URIs returned by a popular BGP RC data broker. This implies that future work needs careful data cleaning and that we may have to revisit previous work.

3.5.3 Service Operators

We wish to underscore what we feel is the underappreciated work put in by the operators of these RC and broker services in running services so many of us in the research community rely on. However, we highlight what we feel are the important steps to take to improve the quality of these services for future researchers: (i) Provide machine-readable indexes of RC data, enabling automated checking by tools ingesting RC data that they have complete datasets. (ii) Provide data integrity checksums over RC data, enabling checking that the ingested data is free from transmission and/or storage errors. (iii) Ensure sufficient computational resources are provided to nodes, to avoid errors or omissions from resource exhaustion. (iv) Eliminate manual, error-prone, processes, such

as the adding of new RCs to brokers, and implement mechanisms for backfilling when data is missed because of errors. (v) Open source the full stack (such as the BGPStream broker), to allow the community to provide their contributions to improving services and verify implementation details.

3.6 Summary

In this chapter, we present our evaluation of inconsistency in BGP route collector data. We show that inconsistencies and ‘lost’ routes occur frequently. Furthermore, we find notable inconsistencies in data provided by a popular route collector data broker. Given the possible impact of these inconsistencies on future and prior work, it is essential to make them known to the community.

We reached out to the operators about these services, i.e., to Routeviews concerning missing and misplaced route updates leading to inconsistent RIBs, to RIPE RIS concerning the anomaly for rrc12, and to CAIDA concerning up to 89.2% of RIPE RIS updates not being announced for December 2022 along with several inconsistencies since 2016. Positively, all three operators are strongly committed to maintaining reliable data, replied quickly, and started resolving these issues.

Chapter 4

Inferring Routing Heterogeneity from AS Paths

Since we have shown in Chapter 3 that route collector data is inconsistent and that BGP distribution software misses data, we take a step further to answer the second sub-question: which assumptions and limitations impact BGP studies? Whereas data can be filtered / cleaned, these limitations hold back research and introduce artifacts in the results of their methodologies.

In this chapter, we have a deeper look at an assumption that is used in a majority of BGP studies: Routing in the Internet can be modelled as a graph of atomic nodes. Although this assumption simplifies methodologies, it introduces uncertainty to the gained insights. Our work investigates how big this uncertainty is by quantifying the number of ASes that do not follow this assumption.

The contribution of this chapter is that we propose an inference method that uses BGP updates to (i) quantify the number of heterogeneous ASes as observed by the route collectors, and (ii) identify ASes which explicitly diverge from the conventional BGP behavior. We analyze the timeframe from 1st January 2018 to 31st December 2023 and find that the most diverse region is RIPE, with 18% of observed ASes being heterogeneous, followed by ARIN with 10% and AFRINIC with 7%. Among the observed heterogeneous ASes, 21% are NSPs, 11% are ISPs, and 7% are Content Delivery Networks. Our findings suggest that neglecting AS heterogeneity in a study's methodology may result in skewed outcomes or misleading conclusions.

4.1 Routing Heterogeneity in Research

As introduced in Chapter 2, in the Internet, different paths may be selected on a BGP router due to variations in routing policies or as a result of the final tiebreaker in the BGP decision process. An example of such a policy is when a route with a longer AS path is preferred over a shorter one. In this case, we assume that the longer path has been assigned a higher Local Preference, indicating the use of an explicit routing policy

that overrides the usual path selection criteria. If this policy is consistently applied throughout the AS, we would expect the network to behave as a uniform system, with all routers preferring the longer AS path. However, if the AS chooses different routes for the same prefix across its BGP sessions, we consider it *heterogeneous*. In contrast, a system where all routers consistently prefer the same route is described as *homogeneous*.

There is a substantial body of work relating to BGP routing, covering topics such as RPKI measurements [47, 59], policy inference [15, 39, 68, 82, 129, 130], BGP community classification [72, 126], topology inference [85, 86, 127], and route prediction [83, 85, 87]. Nearly half of these studies adopt traditional assumptions such as AS homogeneity and “shorter routes are always better”, while the other half explicitly recognize the existence of heterogeneous ASes and acknowledge this as a limitation.

In the case of RPKI measurements, researchers investigate the current and historical state of ROV deployments in the Internet. Most recent studies in this field [28, 47, 59, 91] build their methodologies on various, non-diverse, AS-level topologies of the Internet. Such methodologies that do not consider AS heterogeneity—together with non-diverse AS-level topologies—lack in precision and build only an abstraction of the Internet. This fact was underlined by Li et al. [79], who hint that ASes might have heterogeneous ROV deployments.

Another field concerns the inference of BGP policies, where Wang et al. [129] created a model of the Internet to understand routing policies. Two decades later, this study was reproduced by Kastanakis [68], who found that the use of Local Preference values has significantly evolved over time. Caesar et al. [15] is first to underline that vanilla “shortest path routing” may no longer be sufficient to handle complex routing policies. These studies acknowledge the diversity of ASes as a limitation of their work and outline its importance.

The study by Krenc et al. [72] classifies BGP communities, a transitive BGP attribute, at the AS-level using a fixed set of VPs. However, by assuming that ASes act homogeneously, the work overlooks the potential to identify variations in BGP communities based on where they are intercepted. This limitation highlights the importance of considering AS heterogeneity in such analyzes.

The existence of heterogeneous ASes has been discussed by Mühlbauer et al. [92], addressing the limitations of existing work on AS topology inference that overlook AS route diversity. Choi et al. [30] examines the next-hop diversity using data from a Tier-1 ISP, shedding light on the effects of AS heterogeneity.

While these studies explore the implications of AS route diversity and suggest that assumptions such as “shorter paths are always better” may no longer hold, to our knowledge, no prior research has quantified the extent of heterogeneous ASes or the number of ASes that prefer longer paths over shorter ones, as observed by route collectors. In this dissertation, we fill this gap by analyzing 6 years of public BGP data from 1st January 2018 to 31st December 2023⁸, and emphasize the importance of route diversity introduced by heterogeneous ASes in the Internet.

In particular, our contributions are:

- We propose a methodology to identify ASes that prefer longer AS paths over shorter AS paths. We argue that such a decision is based on a higher Local Preference for the former route, which can be used to implement heterogeneity. Our results show that 12% of the observed ASes behave like this.

⁸The RIPE RIS beacon setup has changed starting 2024 [112]

	Prefix	AS path	Inference
(i)	p :	ABX	
(ii)	p :	ACDX	$\rightarrow A : C > B$
(iii)	p :	ACEX	$\rightarrow C : E \geq D$
(v)	p :	APX	
(vi)	p :	AQOX	$\rightarrow A : Q > P$

Figure 4.1: BGP update sequence annotated with the inferred preferences. Thus, $B : E > C$ means that AS B has a higher Local Preference for prefix p when received from AS E than from AS C . $B : G \geq E$ means that AS B decided to use AS G as well as AS E as neighbors with a preference towards AS G . The latter preference must not be implemented by the Local Preference value.

- We propose a methodology to identify heterogeneous ASes. In effect, we infer the preferred next-AS-hop as seen by different VPs and search for diversity. Our experiments identify that almost 10% of ASes observed in this study show route diversity.
- We apply our methodology to routes for the RIPE RIS routing beacons [111] as observed at the RIPE RIS [109] and Routeviews [22] BGP route collectors over the last six years. We find that RIPE, followed by ARIN and AFRINIC are the regions with the most diverse ASes.

Structure: This chapter is structured as follows. We first introduce the methodology in Section 4.2, the used datasets in Section 4.3, and the implementation in Section 4.4. Section 4.5 presents the results of our analysis. Lastly, we discuss the limitations of our methodology in Section 4.6 and implications of our work in Section 4.7. We present our conclusion in Section 4.8.

4.2 Methodology

Our inference method uses the sequence of BGP updates resulting from a prefix announcement as observed via one BGP session of one route collector. See Figure 4.1 for an example of such a sequence. Notice that all AS paths included in this sequence were selected as the best path. Thus, if the length of an AS path increases the best path decision was made before the AS path length criterion. One might think that other sources, such as failing links or withdrawn prefixes might cause this behavior. We further explore that this is not the case due to the prefixes we use (see Section 4.3) and further analysis (see Section 4.5). This means that the AS used either the Local Preference or vendor-specific additional attribute (e.g., Cisco weight) for this path decision. Thus, we see the effect of a routing policy—there is a higher Local Preference. We refer to the inference in this scenario as a *strong preference*. To find ASes that use such a routing policy we check the sequence of updates for successive AS paths where the latter AS path is longer. We then remove AS path prepending and locate the first difference in the ASNs, e.g., for update (ii) this is AS A . Algorithm 1 summarizes the corresponding algorithm for inferring strong preferences. For the example in Figure 4.1 we have 2 strong preferences: AS A prefers AS C over AS B in line (ii) and AS A prefers AS Q over AS P in line (vi).

We also see some updates where the AS path length does not change but the used neighbor does, e.g., update (iii). This indicates that AS C decided to use a different neighbor, here AS E instead of AS D . This is often due to an explicit routing policy,

Algorithm 1: The algorithm for the strong inference method.

Data: BGP updates**Result:** Set of strong inferences

```
1 inferences ← set();
2 update_old ← empty update;
3 for update in updates do
4   if  $\text{len}(\text{update\_old.path}) < \text{len}(\text{update.path})$  then
5     remove_prepending(update);
6     remove_prepending(update_old);
7     pref ← method(update, update_old);
8     inferences.add(pref);
9   update_old ← update;
10 return inferences;
```

e.g., MED, lower IGP cost, but can also be due to the final tiebreaker. As such we cannot confirm that the change in the best path is due to Local Preferences. Most of these changes are likely due to routing policies as the tiebreaker should not be the used decision criterion for many route choices [14]. For example, the route change in update (iii) is due to a lower MED value. We refer to such effects of routing policies as *weak preferences*. Notice that ASes that yield weak preferences might yield strong preferences depending on the VP. The algorithm corresponds to the one shown in Algorithm 1 with the one exception that we check if two AS paths are of equal length in line 4. The example in Figure 4.1 in line (iii) contains one weak preference, namely that AS *C* used a different route over AS *E* instead of over AS *D*.

The last case involves successive updates where the latter AS path is shorter. Here, we cannot infer anything about the routing policy of the AS that decided to use a shorter AS path. This is due to shorter AS paths being the base choice within the BGP decision process. Given that shorter AS path lengths are the expected behavior one may not expect to get many inferences for strong and weak preferences. In Section 4.5 however, we show that this does not cover the majority of routing decisions.

Up to this point, we consider the BGP update sequences from a single prefix on a single BGP session. Next, we combine information across prefixes and sessions. Hereby, we find that the inferred preferences either agree or disagree. If the inferred strong preferences disagree this indicates that two routers within an AS⁹ use routing policies to prefer different ASes. This means we see a “*conflict*” in the best path choice for this AS and, thus, refer to this AS as *strongly heterogeneous*. To reiterate, one singular BGP session indicates implemented policies, whereas the combination of multiple sessions can point to diversity (conflicts) in this AS.

If all inferences agree this indicates that we either have too limited information about this AS or that the routing policies of the AS lead to the same best neighbor selection. We refer to this AS as being *strongly homogeneous*. Overall, we note that the estimated number of heterogeneous ASes in this study is a lower bound of the actual number of such ASes. This is the case due to the limited amount of available VPs at hand. This

⁹This can also be the same router with two different BGP sessions as shown in Figure 2.1.

further implies that the estimated number of homogeneous ASes in this study is also an upper bound. We might not be able to catch their heterogeneity due to the same reason. We can apply the same reasoning for the weak preference inferences. Here we refer to the AS as either *weakly heterogeneous* if we find disagreeing inferences or *weakly homogeneous* if they all agree. Here, we do not have strict lower / upper bound properties. Still, if an AS is found to be weakly heterogeneous we have a strong indication that it does use routing policies which result in diverse routes.

4.3 Datasets

In this work, we use a variety of datasets. We use publicly available BGP updates from RIPE RIS [109] and Routeviews [22] to infer preferences and classify ASes. Hereby, we focus on routes for RIPE RIS routing beacon [111] prefixes. Then, we use ASRank [82] from CAIDA to map ASes to RIRs and PeeringDB [98] to classify the occupation of ASes.

BGP updates: We download all used BGP data from the RIPE RIS route collectors [109] and the Routeviews Project [22]. Since BGP is a routing protocol and not a measurement tool, the collected data suffers from misconfiguration, errors, etc. It contains artifacts such as unallocated ASNs and poisoned paths, e.g., AS paths with loops. We clean our dataset by removing such announcements following the instructions in [19].

4.4 Implementation

We implement our methodology from Section 4.2 in Bash and Python. The first step is to gather the various BGP updates from the route collectors. We planned to use CAIDA's PyBGPStream framework [96] to execute our analysis. Unfortunately, using CAIDA's BGPStream broker resulted in unforeseen problems due to missing data (see Chapter 3 for an in-depth analysis). Thus, we collect all dumps ourselves and then use CAIDA's BGPReader tool to extract the relevant updates for the RIPE RIS routing beacon prefixes.

In particular, we process the updates in chunks, whereby each chunk corresponds to one beacon announcement period. In effect, each beacon announcement is one experiment that allows us to identify the preference of ASes. Our data collection starts 5 minutes before the beacon's announcement time and ends 30 minutes after the announcement. Overall, we find that there are only a few updates observed before the announcement as well as after 15 minutes after the announcement. This confirms that the beacons are well synchronized and our assumption that BGP converges within 30 minutes holds. Moreover, we only observe a handful of withdrawal messages during these periods which are likely caused by unrelated outages. We decided to ignore these due to inherent problems caused by BGP [43].

Next, we filter the updates to remove malformed announcements. To optimize compute and storage resources we focus on the AS path only and move from a string representation to a list of integers. Next, we sort the updates by their BGP session. Hereby a session is identified by the route collector, the peer's ASN, and the peer's router IP. Next, we apply the algorithm as listed in Algorithm 1 to identify strong and weak preferences. Finally, we combine the information from all sessions and all beacons to classify ASes as homogeneous or heterogeneous.

Description	Overall		Results per Experiment			ASes per Experiment		
	Tot.	Uniq.	Avg.	Max.	Min.	Avg.	Max.	Min
Updates	554 M	554 M	42.192	1.347.667	36	588	734	7
Strong Inferences	17M	499 K	1.306	4.161	95	215	456	41
Weak Inferences	80 M	1.564 K	6.112	13.073	563	351	484	134
Strong Conflicts	0,36 M	14 K	27	242	0	11	77	0
Weak Conflicts	14.8 M	697 K	1.132	4.700	113	89	281	18

Table 4.1: Overview of the results.

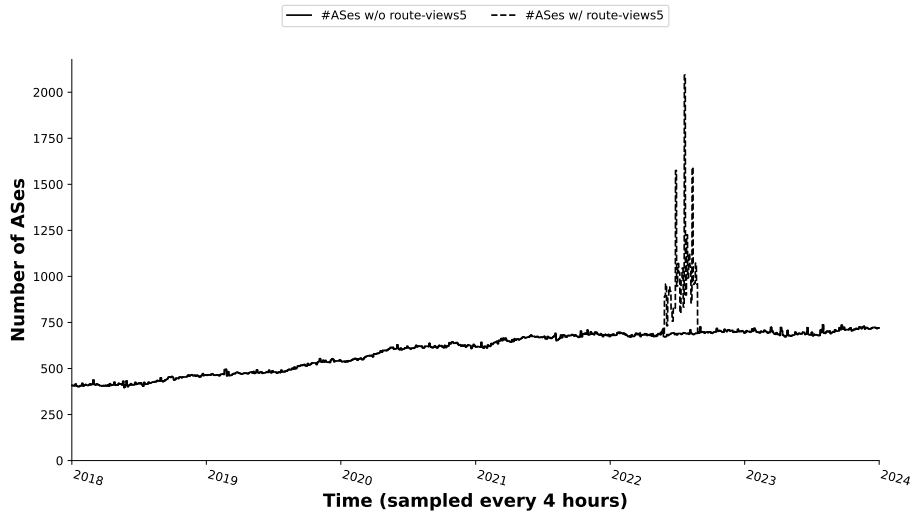


Figure 4.2: Across time: Number of observed ASes for the BGP beacons: (i) for all BGP updates (ii) for BGP updates without route-views5 before September 5th 2022.

4.5 Results

We analyzed BGP data for the 35 BGP beacons from the last 6 years from 1st of January 2018 to 31st of December 2023. In this time we observed a total of 554M updates that, on average, contained 588 ASes in their AS paths. Find an overview of all results in Table 4.1. The total number of ASes in the AS path of these updates is 11872. While this is only roughly 16% of all ASes it still gives us a good basis to test our methodology as it includes all tier-1 ASes and many tier-2 ASes. It does not cover all hypergiants and stub-ASes. Hypergiants are unlikely to be on the AS path of a beacon to a BGP monitor as they are not in the transit business. The same reasoning applies to stub-ASes. Still, the ASes we observe include 53% of the top 1000 ASes according to CAIDA’s ASRank [18]. Table 4.1 shows an overview of the results from our experiments.

Figure 4.2 shows how the number of observed ASes in the AS path varies across time. We find a huge spike which appears on 20th of May 2022 and lasts until 5th of September 2022. Upon closer inspection, we find that Routeviews started to add peers to a fairly new BGP route collector (active since September 2021), namely route-views5, in that timeframe. In particular, it added a peering with AS 40864 which is causing many more

BGP updates than any of the other peers. Moreover, these updates include many ASes that are otherwise not observed. A close inspection of the update sequences showed that AS 40864 is responsible for almost all spikes. These spikes stopped after 5th September 2022 indicating that a misconfiguration of the peering with AS 40864 got fixed. Therefore, for all results below we do not include the BGP updates from route-views5 until 5th of September 2022. After filtering, the average number of observed ASes per experiment is 588 with a minimum of 7 and a maximum of 734. In total, we observe 11872 ASes.

From these updates, we identified preferences for 1478 ASes (1155 strong, 1441 weak). This corresponds to 12% (10% and 12%) of the observed ASes. Using these preferences we find that at least 605 (5%) ASes are strongly heterogeneous and 1039 (9%) ASes are weakly heterogeneous. We note, that we have many more weak inferences than strong inferences. The reason is that the former only focuses on routing policies involving Local Preference while the latter focuses on many more policies. Indeed, the former requires us to observe an update with a shorter AS path followed by one with a longer AS path while the latter just requires updates with equal AS paths.

4.5.1 Inferences Across Time

Figure 4.3 shows how the number of ASes with inferences changes across time both absolute (right y-axis) as well as relative (left y-axis). Notice that the shown values are overlapping and must not add up to 100%. We find that neither the absolute numbers nor the relative numbers vary drastically. Regarding weak / strong inferences we find that on average 37% / 60% of the observed ASes have at least one. As such we can conclude that our methodology is able to highlight that routing policies do play a major role in selecting the best route. Moreover, we find many ASes for which **shorter is not always better** holds.

To move towards understanding routing heterogeneity Figure 4.3 also includes how many ASes we find at least 2 resp. 10 inferences within a single experiment. We find that we quite often find at least 2 inferences (24% strong / 52% weak on average) but that finding 10 inferences is not so common (5% strong / 23% weak on average). Given that we only consider 35 beacon prefixes this is not too surprising. Still, for some ASes, in particular, tier-1 ASes such as Hurricane Electric, we get more than 100 inferences for many of the experiments. Even more ASes have either at least 2 strong or at least 2 weak inferences. This implies that we can study conflicts in the best route choices for a significant number of ASes.

4.5.2 Conflicts Across Time

Thus, Figure 4.4 shows the number / fraction of ASes where we identify a conflict in their best path choice. It is similar in spirit to Figure 4.3.

We note that more than 20% of the involved ASes show route diversity, having a conflict in their best path choice. We also find that not all ASes with multiple inferences use diverse best routes. Indeed, on average across the experiments, we only find 1.5% of observed ASes that are strongly heterogeneous in the sense that they prefer two different longer AS paths. Still, we can identify them and across all experiments identify 605 such ASes. The maximum number of conflicts per AS is 34 for Hurricane Electric. In fact, we found that 60% (35 out of 51) of observed ASes with a CC size over 1024 are strongly heterogeneous. We find that roughly 15% of the observed ASes are weakly heterogeneous. Here, we often observe more than one conflict indicating substantial

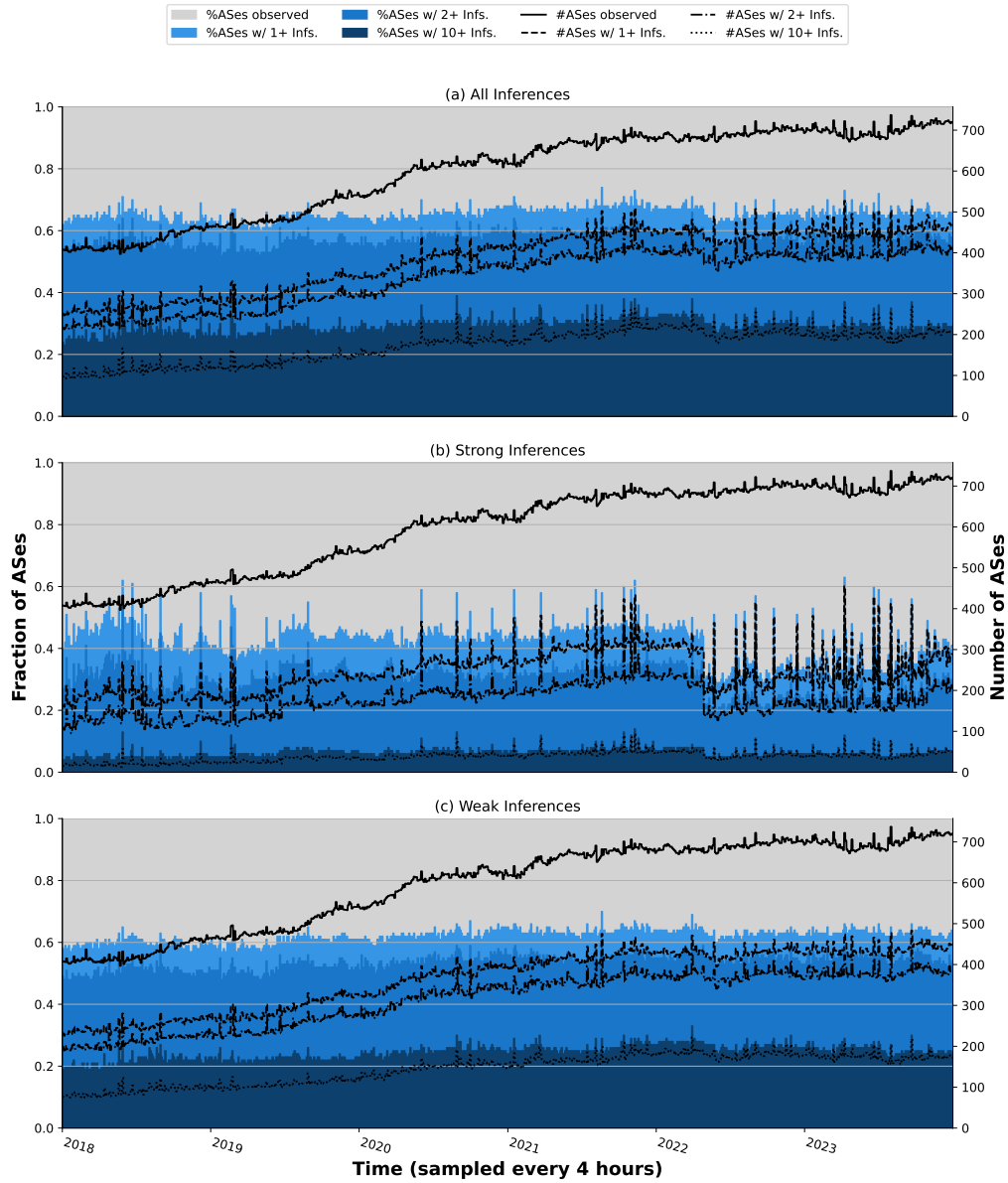


Figure 4.3: Inferences across time: Fraction of ASes (left y-axis) / number of ASes (right y-axis) with inferences of routing preferences—(i) all inferences, (ii) strong inferences, and (iii) weak inferences. To understand how many inferences are made per observed AS we include the total # of ASes the ASes with at least 1, 2, resp. 10 inferences.

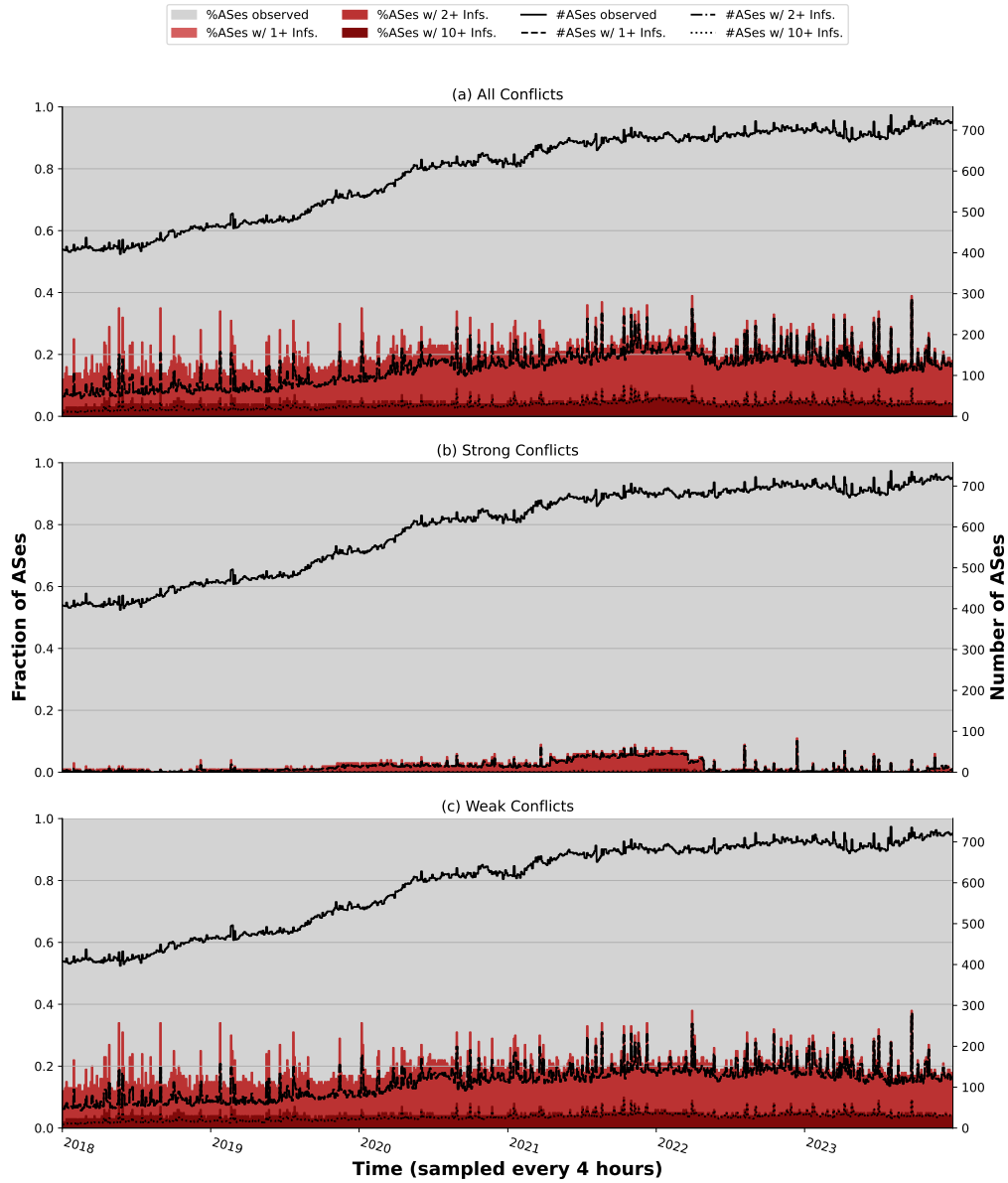


Figure 4.4: Conflicts across time: Fraction of ASes (left y-axis) / number of ASes (right y-axis) with conflicts in their best route choice—(i) all inferences, (ii) strong inferences, and (iii) weak inferences. To understand how many conflicts are observed per AS we again include the total # of ASes the ASes with at least 1, 2, resp. 10 conflicts.

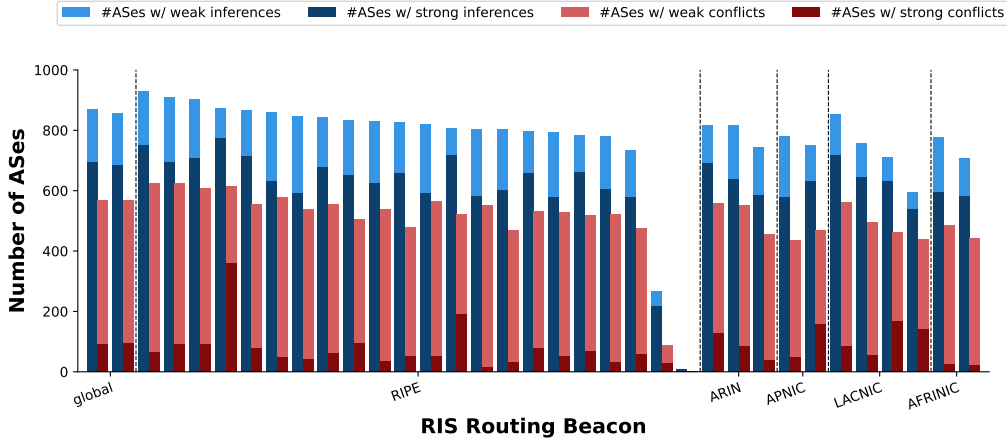


Figure 4.5: Inferences and conflicts across BGP beacons: #ASes with inferences/conflicts grouped by the beacon’s RIR.

diversity in route choices. Indeed, we find a significant fraction (15.4% on average) of ASes with more than 2 conflicts. Thus, we conclude that routing heterogeneity is not unusual and should not be ignored.

4.5.3 Impact of Beacon / Route Collector Choice

Next, we explore if different beacons lead to different results. Figure 4.5 shows the total number of weak and strong inferences and conflicts per beacon. We, hereby, group the beacons by the region where they are announced according to the RIPE RIS’s website. Overall, we find that the total number of ASes with inferences does not vary drastically. Moreover, a Pearson correlation analysis shows a strong correlation between the data collected by each routing beacon for the fraction of weak inferences per beacon ($R=0.95$, $p\text{-value}<0.001$) highlighting that our methodology is not biased by any specific beacon.

However, the number of ASes with strong conflicts does vary. The two globally announced beacons are announced via BGP multihop to more than 100 BGP peers and help us identify many more ASes which are strongly heterogeneous. We observe similar behavior for the BGP beacons announced within the RIPE region via BGP multihop. We also note that the beacon announced within AFRINIC yields fewer ASes with weak inferences. Overall, this highlights that path choices are more limited in this region.

Another source of bias might be the route collector. Thus, Figure 4.6 shows the total number of weak and strong inferences and conflicts per route collector. We again group the route collectors by region where they are located according to the route collector websites. Here, the results show a significant skew. The global route collectors together with the RIPE region allow us to identify many more ASes than those in the other regions. Again, the Pearson correlation test shows a strong correlation between the fraction of ASes with weak inferences and conflicts per route collector ($R=0.96$, $p\text{-value}<0.001$).

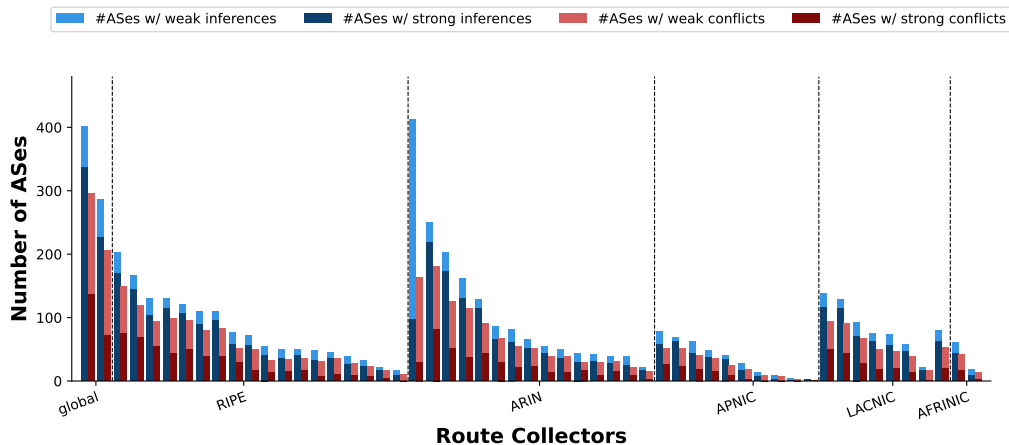


Figure 4.6: Inferences and conflicts across BGP route collectors: #ASes with inferences/conflicts grouped by the route collector’s RIR.

4.5.4 Characterization Of ASes With Conflicts

Finally, we characterize the ASes with inferences resp. conflicts. In the first step, we use PeeringDB to categorize ASes by type. we distinguish between NSPs, ISPs, CDN, Educational (Edu.), Non-Profit (NP), Enterprise (Ent.), Network Services (NS), RS, Governmental (Gov.), and Other. Figure 4.7 shows a barplot of the number of ASes with inferences resp. conflicts by network type. The majority of identified ASes are either ISPs or NSPs. Next, we have CDNs, Educational networks, and the “other” networks. One may be surprised to find that CDNs are prominent, given that beacons are not originated at CDNs and that CDNs are typically not offering transit. However, many networks offer diverse services including peering and content hosting and some of these are classified as CDNs by PeeringDB.

Note, just as before the proportion of conflicts to inferences is similar for all types of AS. As we do not find a specific bias this suggests that AS path diversity may not be related to the main activity of the ASes.

Next, we put our observed ASes in relation to the total number of ASes. Hereby, we split the ASes by RIR. Figure 4.8 shows the percentages per region. We get the number of registered ASes per RIR from RIPE’s weekly Routing Table analysis—we used the report from the 9th of August 2024. To identify the RIR of ASes with inferences and conflicts we use CAIDA’s ASRank data. In the Figure 4.8 100% is marked in grey and corresponds to all ASes in a region. The fraction of observed ASes is shown in dark grey on top of the grey base bar. The fraction of ASes with inferences (both strong and weak) is shown on top of the dark grey bar in blue. The fraction of strong and weakly heterogeneous ASes (ASes conflicts) is shown in red (on top of the blue bars).

As expected, RIPE has the largest number of ASes (11925 total, 3004 observed). Our method shows that 19% of observed ASes in the RIPE region are heterogeneous. Interestingly, we find that ARIN (4958 total, 1631 observed) has fewer ASes than APNIC (9715 total, 2363 observed) and LACNIC (9046 total, 4035 observed) but has a larger fraction of heterogeneous ASes 10% (6% and 2.6% respectively). APNIC’s and LACNIC’s smaller percentages are in part due to the limited coverage by the BGP routing beacons. While it may be expected to see a sizable fraction of heterogeneous ASes in RIPE we were

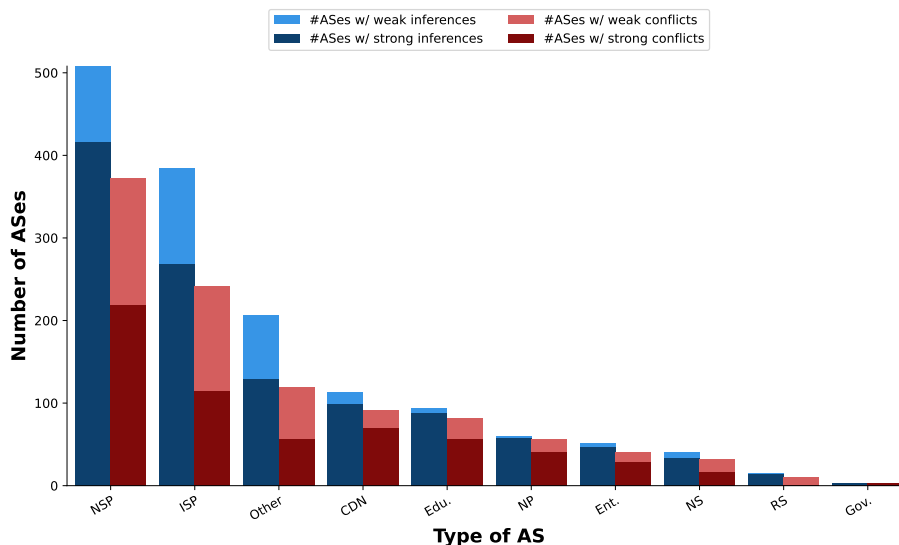


Figure 4.7: Number and fraction of ASes with inferences / conflicts per AS type.

surprised to see AFRINIC (1230 total, 473 observed) with a relatively large percentage of 7%. However, we want to point out that this may be an artifact of the smaller number of ASes in AFRINIC. Thus, even finding a few heterogeneous ASes increases the fraction substantially.

Overall, we find that routing heterogeneity is quite prominent and that our methodology is only biased by the coverage of the routing system provided by the beacons, in the sense that we sample certain regions better with the BGP route collector system.

4.6 Limitations

So far our experiments are restricted to the RIPE beacon prefixes. These are few in number. Still, they let us infer preferences for a surprisingly large number of relevant ASes (according to ASRank). However, the coverage is somewhat limited since the BGP beacons are not routed via hypergiants or stub ASes. Thus, the routing policies of these ASes are harder to infer. To generalize our study one should explore how to extend the methodology to take advantage of BGP updates for other prefixes as well.

Moreover, we depend on the BGP route collector system to accurately reflect the state of the routing system via BGP updates. If this is skewed, e.g., due to some misconfiguration or other reason—recall the large spikes generated by routeviews-5—the results can be misleading.

Another limitation is that we assume that there are no additional events that influence the BGP update sequences of the RIPE RIS routing beacons. In effect we assume that the Internet does not suffer outages, does not add / remove links, and does not change its routing policies—at least not in a way that it impacts the convergence process for the routing beacons. However, we acknowledge that this may not always be true. One such indication is if we infer the same preference multiple times for the same session, i.e., a next-hop AS was less preferred and then again more preferred. This can happen if a

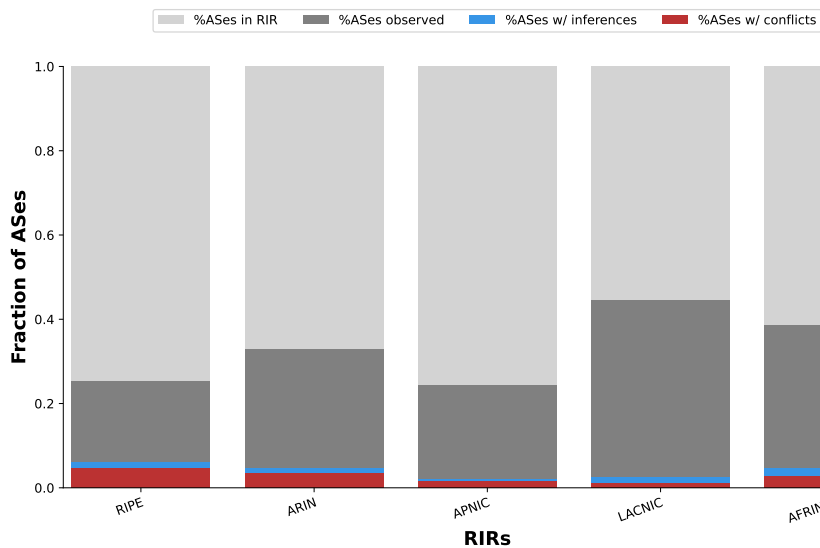


Figure 4.8: ASes by RIR: Fraction of ASes with inferences / conflicts vs. total number of ASes in RIR.

route becomes unavailable for a short time for some reason. Figure 4.9 shows how many of the total number of inferences are repeated inferences (the dashed line). While there are some clear spikes visible in the number of repeated inferences they are almost all short-lived and only affect a small number of our inferences. Thus, we conclude that our inferences are mostly not affected by other events.

4.7 Discussion

Overall, our results indicate that the number of heterogeneous ASes is substantial, which challenges a common assumption made in prior work—namely, that the routing decisions of ASes are homogeneous. This finding is significant because much of the research on Internet routing protocols has relied on this assumption for simplicity. The assumption of homogeneity implies that all routers of an AS behave similarly when making routing decisions, often leading to simplified models and conclusions. However, our data suggest that this may not be the case, particularly as the Internet continues to evolve and undergo a "flattening" process [6, 12]. The increasing deployment of large IXPs and the proliferation of peering arrangements introduce new complexities into how ASes interact, as noted in previous studies [101].

Our findings align with this observation, as we confirm that Tier 1 ASes, as well as most Tier 2 ASes, maintain multiple peering connections and, consequently, exhibit diverse path selection. These ASes are at the core of the Internet's infrastructure, and their heterogeneous routing behaviors can have far-reaching implications for the global network. Traffic engineering, congestion management, and security policies can all be shaped by the diversity in AS routing, which suggests that simplifying assumptions about AS behavior may overlook important nuances that affect the overall performance and resilience of the Internet [5]. As a result, it is essential that researchers *(i)* adjust their methodologies to account for AS heterogeneity, *(ii)* re-run their analyzes to ensure

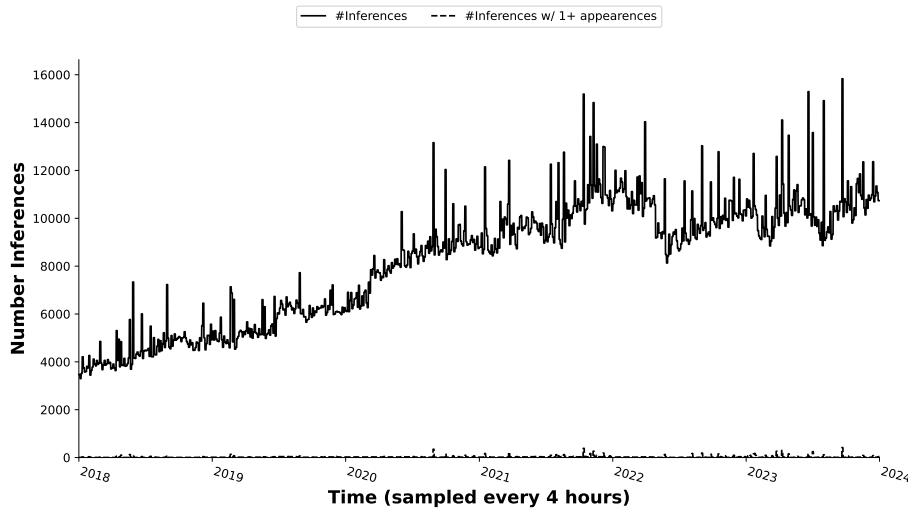


Figure 4.9: Number of inferences across time: total number as well as repeated inferences.

that their findings hold under more realistic conditions, and *(iii)* compare the results to understand the impact of assuming homogeneity versus acknowledging heterogeneity. These steps will help the community understand the full extent to which this assumption has influenced previous work and whether new insights can be gained by adopting more nuanced models. We note that this type of re-examination is beyond the scope of the current paper and require significant effort in future research.

One method to further validate our findings and address the heterogeneity of AS routing policies is through the use of LGs. ASes operate LGs to allow network operators—and the general public—to observe the effects of routing policies in real-time. These LGs are web-based tools that provide insights into BGP routing decisions by querying routers within an AS and displaying the available routes for specific prefixes. By using these tools, researchers can observe multiple routes per prefix, and sometimes, key BGP attributes that affect route selection, such as Local Preference, AS path, and MED. In our initial validation using LGs, we are able to confirm a few inferences from our study, demonstrating the utility of this method. There are several challenges when relying on LGs for a comprehensive analysis.

One key limitation is that LGs typically do not provide historical data, making it difficult to conduct longitudinal studies or validate results over time. Since most LGs only offer a real-time snapshot of an AS’s routing policy, researchers are limited in their ability to track how routing decisions evolve or react to external events such as outages or policy changes. However, a recent study that collects LG data over time [57] suggests that it may be possible to overcome this limitation by building a more extensive historical dataset. Another challenge lies in the diversity of LGs. They vary considerably in the level of detail they provide, as well as in the geographic distribution of the participating ASes. This variability can result in incomplete or biased insights, with certain regions or networks being overrepresented in the data. These factors complicate efforts to generalize findings and make LGs less reliable as a standalone tool for validation. Despite these limitations, LGs remain a valuable resource for real-time observation of routing policies, offering critical insights into the behavior of ASes.

An alternative approach for understanding AS routing policies is analyzing withdrawals in the BGP protocol. While much attention has been focused on announcements in previous studies, withdrawals may offer additional insights into how routing policies are implemented. When a prefix is withdrawn, it signals that a particular route is no longer available, prompting ASes to explore alternative paths. However, our initial exploration indicates that the signals from withdrawals are not as straightforward as those from announcements [43]. The BGP path exploration process that follows a withdrawal is inherently more complex. As such, drawing clear inferences from withdrawal data is more error-prone, especially when considering the heterogeneity of AS routing policies. This complexity adds another layer of difficulty to use withdrawals as a reliable source of information on AS behavior.

4.8 Summary

In this chapter, we present a methodology to identify ASes that for some prefixes prefer routes with a longer AS path over those with shorter AS paths. As such “shorter is not always better”. This indicates that the route choice is the result of a more complex routing policy. Moreover, our methodology allows us to infer ASes with non-homogeneous, i.e., heterogeneous routing policies. This means that we identify ASes where two border routers prefer routes with different AS paths to reach the same prefix.

Our analysis shows that $\approx 5\%$ of ASes that are observable when applying our methodology to the RIPE RIS beacon prefixes are heterogeneous and ≈ 9 are likely to be heterogeneous. We also observe that the ARIN and RIPE regions host more heterogeneous ASes compared to the rest of the RIR regions. While the choice of routing beacons does not skew our results they are skewed in terms of BGP route collectors. Route collectors with many peers are able to infer more inferences. However, the ratio of inferences vs. conflicts roughly stays the same. Lastly, we find that Internet and Network Service providers tend to be more heterogeneous than other ASes.

Chapter 5

The Missing Hivemind of BGP Routers

In this chapter, we continue with the second step to answer the second sub-question: which assumptions and limitations impact BGP studies? Although the methodology from Chapter 4 already provides solid evidence of the existence of heterogeneous ASes, it only touches routing updates and not RIB snapshots.

The contributions of this chapter are as follows. We propose another inference method that uses RIB snapshots instead of routing updates. With this method, we analyze routing data of the past decade to find trends in the adoption of heterogeneity. As such, we observe an increase of 198% in the number of heterogeneous ASes since 2014. The increase is especially the case for ISPs (271%) and NSPs (175%), but also surprisingly for CDNs (323%). Our results show that all ASes with an ASRank of 75 or better are heterogeneous. In the case of the CC size, all ASes with more than 10k ASes are heterogeneous. Lastly, we analyze the route collector setups that enabled us to find heterogeneity. The findings of this study are another strong indication that neglecting AS heterogeneity may result in skewed or misleading conclusions.

5.1 Routing Heterogeneity in BGP

As previously introduced, all routers inside an AS are operated by a single administrative organization that externally looks like a network with a unified routing policy. The BGP exchanges reachability information between ASes on a per-prefix and per-network basis. ASes implement routing policies to determine the path through the network or if traffic should not pass through it. BGP offers adaptable attributes associated with each route to implement policies and perform route filtering. An example of routing policies is hot-potato routing, where traffic should leave the AS's network as fast as possible.

Although BGP routers learn about multiple routes to the same prefix, they only propagate the best route for scalability reasons (M-BGP [77, 78] is not yet widely deployed). BGP attributes are used to iteratively filter and select the most appropriate path for each prefix. However, each router performs this decision on its own. Thus, individual routers can choose different routes for the same prefix, given their configuration. We consider an

AS *heterogeneous* if it chooses and propagates at least two paths for the same prefix. In contrast, a network where all routers propagate the same route for a prefix is called *homogeneous*.

Heterogeneous ASes can not always be considered as one uniform entity by research, e.g., while analyzing RPKI coverage as highlighted by the *following example*: We find that AS 8888, a full-feeder for multiple route collectors, shows heterogeneous RPKI coverage. The AS peers with `route-views.sydney` and `route-views.amsix`. While it only propagates RPKI valid Cloudflare prefixes to the route collector in Sydney, it propagates both valid and invalid Cloudflare prefixes [36] to the route collector in Amsterdam. We observe different BGP paths for the valid prefixes visible at both route collectors. These findings indicate that AS 8888 has a heterogeneous behavior, and RPKI is impacted based on the path (before and / or within the AS).

Given the potential impact of heterogeneous ASes, it is important to get an overview of the general number of heterogeneous ASes, regional developments, and business sector adaptations. This study provides a longitudinal evaluation of heterogeneous ASes, their frequency, and characteristics. *In particular, our contributions are:*

(i) We apply a methodology to identify pairs of ASes that show detours in their AS paths for the same prefix at different VPs. We find that 3855 ASes (5% of the assigned ASNs as of April 2025) show heterogeneous behavior, which has increased over the last decade. The increase of heterogeneity is independent of the growth of route collector projects, but is also visible based on stable BGP sessions over the last decade.

(ii) We analyze the characteristics of heterogeneous ASes, e.g., their region, business sector, rank, and CC size [18]. Although it is expected that mostly ISPs and NSPs are heterogeneous, we found that CDNs also show heterogeneity. We also find that most heterogeneous ASes have lower ASRanks. However, all higher-tier ASes with a rank of 75 or better are heterogeneous. For all ASes with an ASRank of 500 or better, only 40 are not heterogeneous.

(iii) We provide a detailed analysis of our findings for December 2024 to find how well we can catch heterogeneity. We find that one route collector is enough to find 73.4% of heterogeneous ASes, but it takes 35 route collectors to find more than 99%. Similarly, one announced prefix is enough to find 26.8%, whereas more than 1M prefixes are needed to find 99%.

5.2 Related Work

BGP is a diverse research area that covers topology inference [85, 86, 127], BGP community classification [72, 126], route prediction [83, 85, 87], policy inference [15, 39, 68, 82, 129, 130], and RPKI measurements [47, 59]. Early work by Mühlbauer et al. [92] and Choi et al. [30] evaluated the homogeneity limitation of, for example, AS topology measurements. However, most recent studies adopt the assumption that ASes are homogeneous systems, while few recognize and accept this assumption as a limitation of their work. Especially simulation-based RPKI measurements rely on the assumption that ASes are homogeneous nodes in a graph. Recent studies in this field [28, 47, 59, 91] consider a non-diverse AS-level topology of the Internet for their simulations. These methodologies refer to an abstraction of the Internet. Removing AS heterogeneity from a simulation—especially from the used AS topology—reduces precision. Li et al. [79] underline this fact by deploying a real-world measurement system for RPKI. They measure ROV deployments in the Internet with distributed virtual VPs. Whereas they can fully

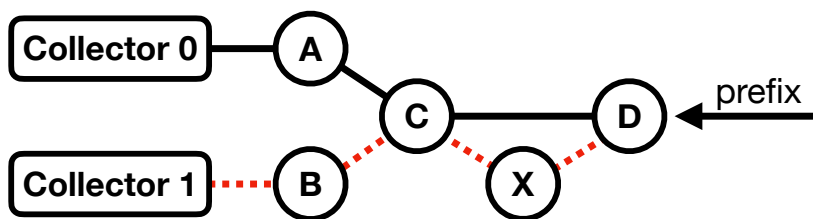


Figure 5.1: An AS topology where AS *C* announces different routes to different peers: Directly to AS *D* and a detour via AS *X*.

classify ROV deployments in 95.1% of the cases, there are 4.9% that show mixed results that heterogeneous ASes could explain.

This work highlights the prevalence of heterogeneous ASes today, especially for top ranked ASes, and provides insights into their characteristics. The Internet was, to the best of our knowledge, not investigated in terms of heterogeneity for over a decade—it has vastly evolved. Our work aims to show trends in adoption of heterogeneity since 2014.

5.3 Detecting Heterogeneity

Our methodology identifies detours—paths that are only visible at specific VPs—between two ASes. Figure 5.1 shows an example of a heterogeneous AS, namely AS *C*. AS *D* announces a prefix to its neighbors, and BGP converges for this prefix. AS *C* sees two routes towards the prefix, directly via AS *D* or AS *X*. AS *C* announces the direct route to AS *A*, which shares *ACD* with Collector0. However, AS *B* receives a different announcement from AS *C* and shares the path *BCXD* with Collector1. AS *C* decided to announce different routes to its neighbors—a clear sign of heterogeneity. In this case, AS *C* used a detour over AS *X* to reach AS *D*. These detours can occur due to different policies within AS *C*, e.g., due to regional differences and different routers.

We identify ASes as heterogeneous if different routes towards a prefix are visible at the same time. This requires VPs that collect routing information from different peers or at different topological locations at the same time. We rely on RIPE RIS [109] and Routeviews [22]. These services deploy geographically distributed BGP routers (route collectors) that are connected to ASes. They receive their peers’ best routes for prefixes and locally record all information in their RIB. The route collectors dump their RIB at fixed intervals providing a view of the Internet via different peers and at different VPs. We use all available route collector RIB snapshots (that are released at the same time) to find detours between ASes. Although the chosen RIB snapshots only present a fraction of the available data, we argue that most of the content in successive RIB snapshots is redundant. Due to resource constraints, we decided to present an overview of heterogeneity in the Internet and a longitudinal view over 10 years with this study and leave a detailed analysis and short-term effects for future work. In this work, we do not use GILL [4] since their peers do not publish RIB snapshots at fixed times and available data does not allow a longitudinal study yet.

As compared to the methodology introduced in Chapter 4, this methodology has many advantages: The methodology in Chapter 4 relies on finding conflicts in update sequences from the routing beacon prefixes, this methodology can use all data for inference,

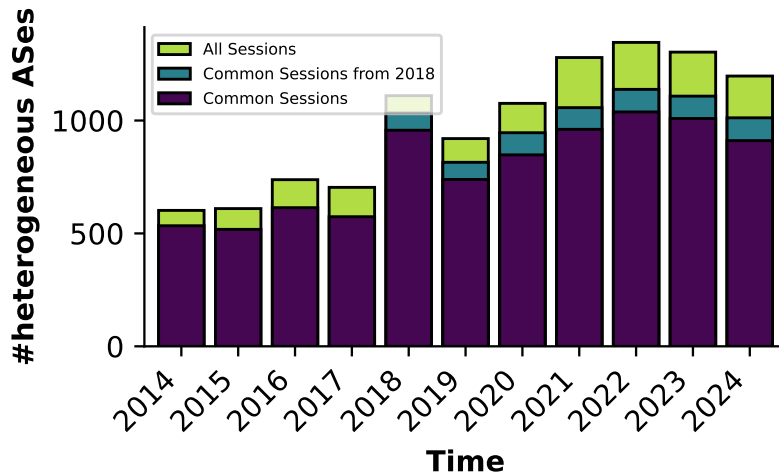


Figure 5.2: Number of heterogeneous ASes based information from Routeviews and RIPE RIS route collectors. The full bar is based on all data. The blue subset is only based on sessions stable throughout the 10 year period.

The methodology in Chapter 4 distinguishes between strong and weak conflicts—only strong conflicts are a clear pointer to heterogeneity. This methodology yields strongly heterogeneous ASes. We argue that both methodologies are important and mutually exclusive in their findings.

5.4 Adoption of Heterogeneity

We used our methodology on historical data from the beginning of 2014 until the end of 2024. We analyze 4 full days of data per month. We analyzed each month’s 1st, 9th, 17th, and 25th. Since the snapshot timings of the route collector projects differ, we took matching periods. This equates to 3 sets of snapshots daily, i.e., at 00:00, 08:00, and 16:00. In total, we applied our methodology to 1440 sets of snapshots.

For this time frame and all available route collectors, we found a total of 3855 heterogeneous ASes. Figure 5.2 shows the development over 10 years. The number of heterogeneous ASes has doubled since 2014. However, the number of route collectors has increased, especially from 2018 onward. Thus, ASes might have been heterogeneous, but we could not detect them. We want to make sure that the increase in heterogeneity is not only the case due to more VPs in later years. Figure 5.2 shows the number of heterogeneous ASes that we found with: (i) all BGP sessions, (ii) BGP sessions that are present in all 10 years of our analysis, and (iii) BGP sessions that are present since 2018. The latter two cases do not rely on an increasing number of VPs to find more heterogeneity. Applying our methodology to all sessions that stayed constant since 2014, we find an increase of 170% in heterogeneous ASes. Thus, while new VPs reveal further heterogeneity, our findings show an increase in heterogeneity in the Internet in general.

In Figure 5.2, we can see a spike in findings starting in 2018. There is no apparent reason for the increased number of heterogeneous ASes in 2018. Our evaluations show that no AS, organization or prefix stands out, but more heterogeneity is visible. A reason might

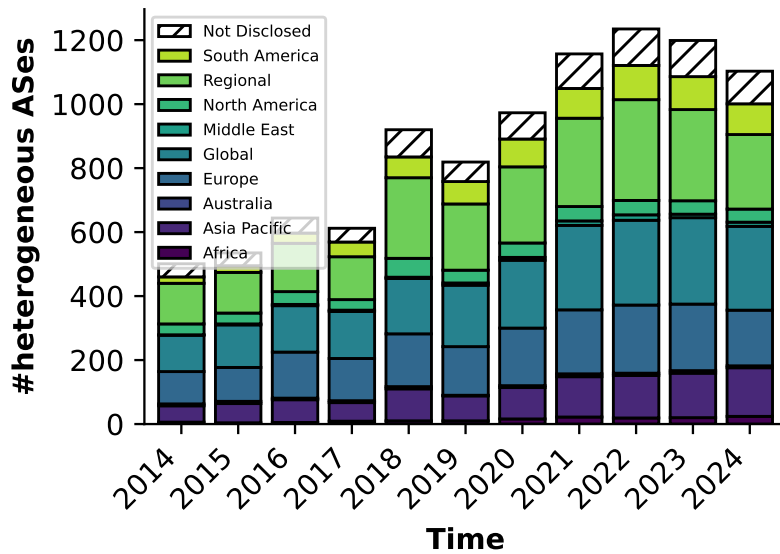


Figure 5.3: Regionality of heterogeneous ASes based on PeeringDB. For 440 ASes no information is available.

be our sampling due to resource constraints. Our findings are lower bounds to the actual number of heterogeneous ASes in the Internet.

Our methodology finds heterogeneity by comparing detours. For December 2024, we observed 1 379 113 detours, of which 966 detours have a length of 10 or more. 79.9% (772) of these detours are originated by CERNET and always have a similar AS path. All these paths contain the sub-path *AS 6939 AS 24239*, which introduces the CERNET paths to the global Internet. Another 5.8% (56) of these detours originate from AS 207744, which is hosted in Russia. Our findings that have longer detours yield 20 heterogeneous ASes which are also covered by other results. These paths do not introduce new heterogeneous ASes to our study and are a research artifact instead.

5.4.1 AS Characteristics

In addition to the prevalence of heterogeneity, understanding heterogeneous ASes and their characteristics is important to understand their impact in the Internet and research. We use PeeringDB [98] to assess regionality and occupation. We use a PeeringDB snapshot from 2024 to make the results comparable. Later snapshots contain more data and, ASes do not frequently change their regionality and / or occupation [81].

Interestingly, diversity is large and not only limited to large *global* ASes. A total of 2.6k ASes mention in PeeringDB that they operate on the global level and enable routing for smaller ASes. Even regional ASes can be heterogeneous based on their routing policy. Most ASes (17k) are comparably small and operate on a continent (or regional) level. Figure 5.3 shows the regionality of heterogeneous ASes for the last decade. All regions became more heterogeneous over time. This is primarily visible in the regions of Asia Pacific (increase of 298%), Europe (172%), and South America (480%). Globally (increase of 229%) and regionally (183%) operating ASes have become more heterogeneous as well.

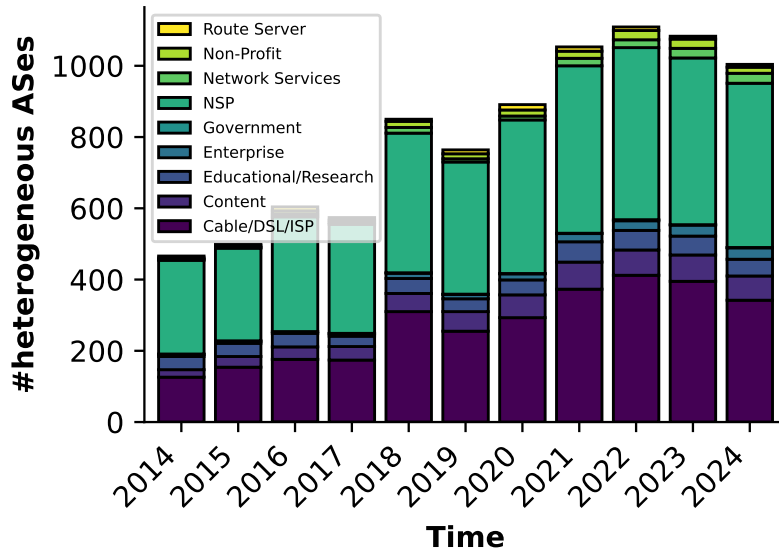


Figure 5.4: Occupation of heterogeneous ASes based on PeeringDB. For 440 ASes no information is available.

We examine the occupation of heterogeneous ASes to see whether certain types, e.g., NSPs are heterogeneous. Figure 5.4 shows the occupation of heterogeneous ASes and their development over time. Occupations that offer network services are more prone to being (and becoming more) heterogeneous. This is the case, for example, for ISPs (increase of 271%) and NSPs (increase of 175%). For them, we see a continuously increasing trend. However, ASes that belong to other occupations are also heterogeneous. Interestingly, we find heterogeneous CDNs (323% with 68 ASes in 2024). In theory, CDNs should not transit data—data transit is a base requirement for our methodology. We argue that the occupation field in PeeringDB abstracts a lot of information away. An example is AS 24429 (Alibaba Cloud), which we find to be heterogeneous. Although this AS is labeled as a CDN in PeeringDB, it offers many different services to its customers, one such service being transit.

5.4.2 ASRank

In addition to self-reported characteristics from PeeringDB, we use CAIDA’s ASRank and CC size [18] to understand the importance of heterogeneous ASes as part of the Internet. We used a snapshot of CAIDA’s ASRank dataset from 2024. An AS is more influential if it has a high ASRank and a large CC. Figure 5.5 shows a Cumulative Distribution Function (CDF) of the fraction of heterogeneous ASes and detours per ASRank. The step-like features are due to ASRank, which groups ASes with the same CC into one rank. Figure 5.5 shows heterogeneous ASes across all ranks and in general more lower-tier heterogeneous ASes. It is important to mention that this is the case due to the larger amount of ASes in lower ranks. Table 5.1 shows the fraction of heterogeneous ASes per rank (and CC size) for the whole time of our analysis. For all ASes with an ASRank of 500 or better, only 40 are not heterogeneous and all ASes with a rank of 75 or higher are

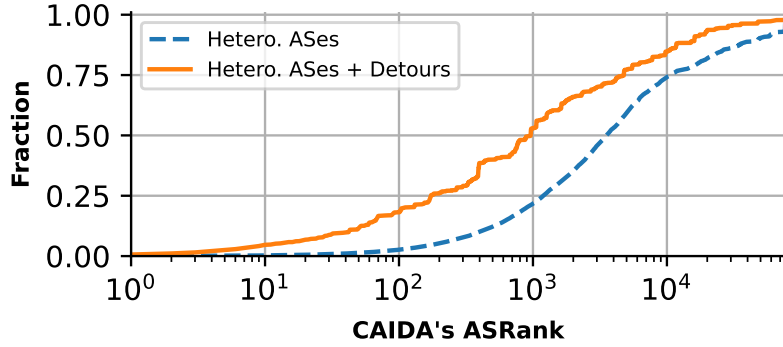


Figure 5.5: CDF for the fraction of observed heterogeneous ASes and detours per CAIDA's ASRank.

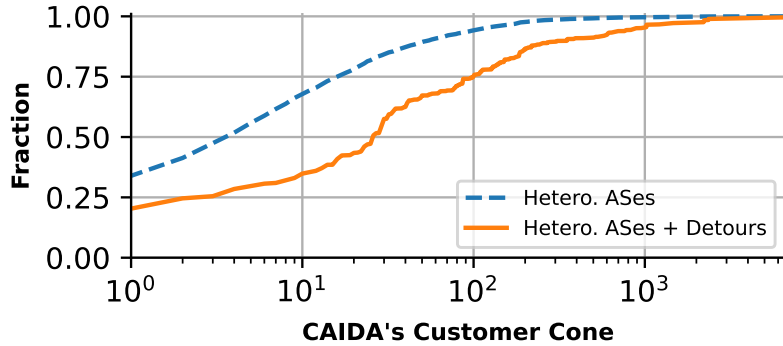


Figure 5.6: CDF for the fraction of observed heterogeneous ASes and detours per CAIDA's CC.

heterogeneous. The fraction decreases for lower-tier ASes, with only 0.9% of ASes being heterogeneous for ranks below 10k.

Metric	0	<10	<100	<1k	<10k	<100k
ASRank	/	100.0%	98.9%	79.7%	21.9%	0.9%
CC Size	0.8%	16.7%	60.4%	91.5%	100.0%	/

Table 5.1: Fraction of ASes per ASRank and CC size that are heterogeneous.

Similarly, larger CCs are expected to indicate heterogeneity. Figure 5.6 shows a CDF of the fraction of heterogeneous ASes and detours per CC size. Surprisingly, 12.5% of our inferred heterogeneous ASes have a CC size of zero. These ASes have multiple peering sessions with the route collector projects, revealing heterogeneity within their network. Table 5.1 shows that this only covers 0.8% of all ASes with this CC size.

5.4.3 Key Takeaways

Heterogeneous ASes are frequent and the numbers increase over time. Furthermore, they are not limited to specific regions, occupations, or sizes. Even ASes without CCs

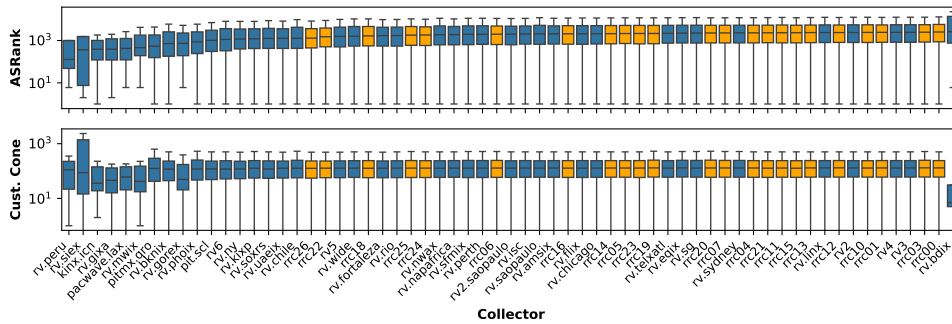


Figure 5.7: Boxplots about the ASRank (and CC size) of heterogeneous ASes found by Routeviews (blue) and RIPE RIS (orange) route collectors.

are internally heterogeneous. Future research should consider heterogeneity and extend existing Internet models to capture it. Due to the characteristics of ASes, their impact can be diverse, depending on the research.

5.5 Catching Heterogeneity

Besides the prevalence and characteristics of heterogeneous ASes, we shed light on the importance of route collectors to catch heterogeneity. The operation of route collectors faces many challenges, e.g., the physical machines or the amount of data [107]. For this reason, many route collectors have started implementing stricter peering policies or stopped operating altogether [90]. One question to ask is whether the choice of route collectors and their peering influence whether we find heterogeneous ASes.

5.5.1 Route Collector Characteristics

Figure 5.7 shows boxplots of the ASRank and CC size of the heterogeneous ASes over 10 years of data per route collector. Route collectors from the Routeviews project are displayed in blue and RIPE RIS in orange. Some outliers that defy the overall pattern can be seen for `rv.bdix` and `rv.siex`. The former finds heterogeneous ASes with a lower ASRank and thus smaller CCs, whereas the latter finds the opposite.

The route collector `rv.bdix` only started operation in 2021 and is not peering with full feeders. The peers do not share their full routing tables; thus the route collector only receives a small number of routes. From this route collector, we mainly found low-tier ASes that are heterogeneous. One possible cause is that we only found heterogeneity from two peers, namely AS 38493 and AS 9825. These ASes are relatively small. In conclusion, this route collector can still observe heterogeneity, but does not yield many results. Similarly, peers at `rv.siex` are not full feeders, and we only found heterogeneity through two peers, namely AS 212271 and AS 50839. Both belong to the same organization, a major infrastructure hosting service based in Italy. Due to this, they might be well connected and allow us to observe top-tier, heterogeneous ASes, such as Hurricane Electric.

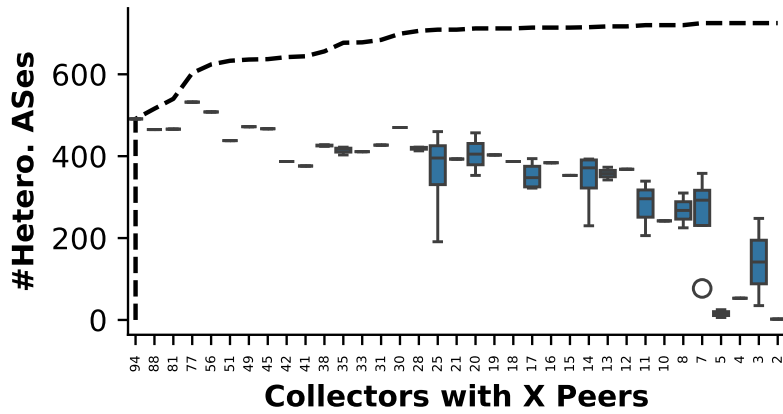


Figure 5.8: Boxplot for #Hetero. ASes found by route collectors that have X peers (ASes) as of December 2024.

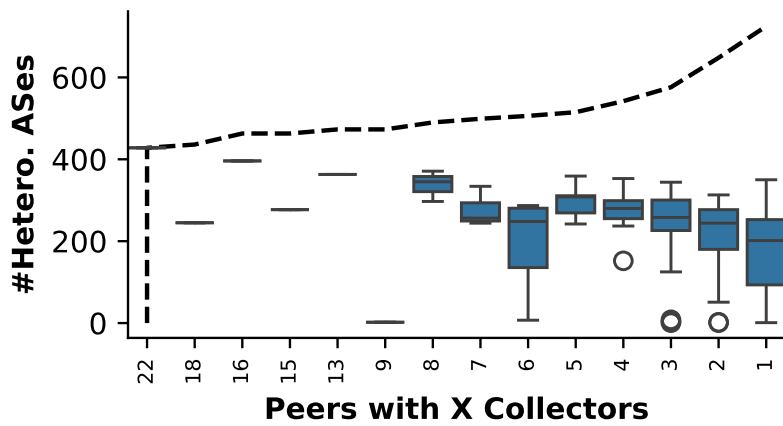


Figure 5.9: Boxplot for #Hetero. ASes found by peers that reveal heterogeneous ASes and peer with X route collectors.

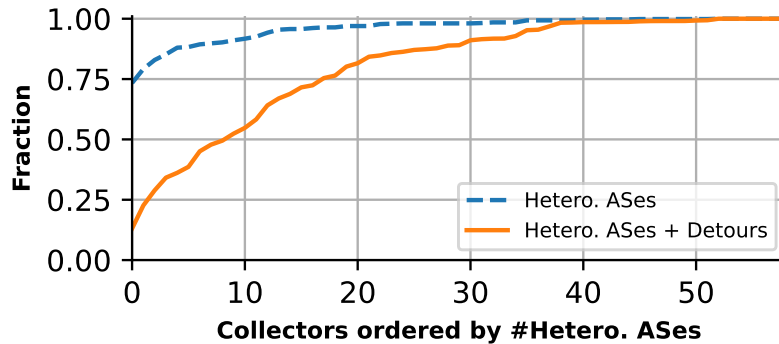


Figure 5.10: CDF for the fraction of observed heterogeneous ASes and detours per route collector.

5.5.2 Route Collector Requirements

Current route collector setups are struggling due to resource utilization and cost cuts [106]. Thus, it is of interest to study the aspects of these setups that allow us to catch heterogeneity. We do so for the entire data for December 2024. In this period, no route collectors or peers are added or removed. Firstly, we evaluate how the number of peers per route collector influences our methodology. Figure 5.8 shows the route collectors ordered by their number of peers on the x-axis. The y-axis shows boxplots and a CDF of the number of heterogeneous ASes for each group. We find that heterogeneity is best observed with route collectors that have many peers. However, all route collectors contribute further insights. For example, the route collector `rrc03` has 94 peers and finds 491 heterogeneous ASes. In contrast, smaller route collectors also provide value, e.g., `route-views.kixp` with only 8 peers still sees 310 heterogeneous ASes.

Secondly, we evaluate the influence of ASes that enter peering with multiple route collectors. Figure 5.9 shows our results for peers reversely ordered by the number of route collectors they are connected to. AS 199524, a CDN, peers with 22 route collectors, enabling us to find 428 heterogeneous ASes. Peers that are connected to a single route collector help to find heterogeneity. This is probably caused by the comparably large number of ASes (383) that peer with only one route collector. These ASes enable us to find 89.6% of heterogeneous ASes. In comparison, we find heterogeneity through 164 further ASes that peer with two or more route collectors.

Figure 5.10 shows a CDF of what fraction of heterogeneous ASes and detours we observe per route collector in December 2024. Since there is no obvious order in which the route collectors should be added to the CDF, we ordered them based on the number of heterogeneous ASes they found. The route collector `rrc25` observes 73.4% (532) of heterogeneous ASes, even though it is only rank 4 based on peers (77 peers, see Figure 5.8). However, the same route collector only sees 13.1% (180 591) of the detours. This makes sense since route collectors are deployed at vastly different locations around the globe. Further, they mostly do not have the same peers and thus have different views of the Internet. Thus, we must include as many route collectors as possible in our methodology to find heterogeneous ASes. We further break it down to BGP sessions (see Figure 5.11). We can see a similar trend to that in Figure 5.10. The same reason applies: route collectors cover diverse locations. Especially since each route collector can have multiple sessions, there are multiple sessions "at the same location". We found that it is enough to consider 742 BGP sessions (42.2%) to already observe more than 95% of

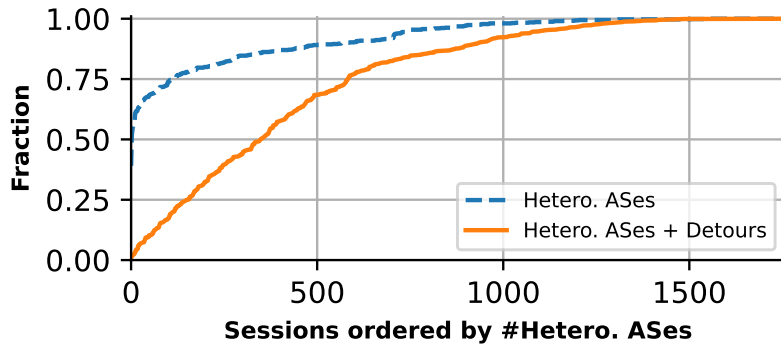


Figure 5.11: CDF for the fraction of observed heterogeneous ASes and detours per BGP session.

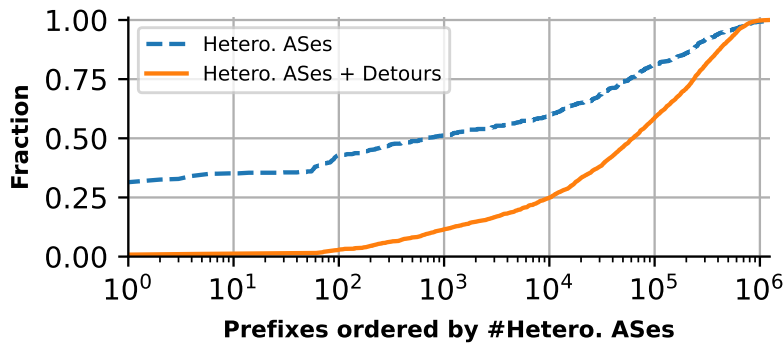


Figure 5.12: CDF for the fraction of observed heterogeneous ASes and detours per prefix.

heterogeneous ASes. The most granular view we can have is to investigate the prefixes the route collectors see. These prefixes are propagated over the Internet, create detours, and thus allow us to find heterogeneous ASes. As seen in Figure 5.12, even one single prefix that is observed at many route collectors is enough to find 26.8% of heterogeneous ASes. Surprisingly, this prefix (2804:2058::/32) belongs to IPv6 and is hosted by an AS in Brazil. Nevertheless, the number of observed detours for one single prefix is minimal.

5.5.3 Key Takeaways

Well-connected route collectors and few prefixes are enough to identify heterogeneous ASes. We argue that (at least to catch heterogeneity) route collector projects should not decrease location diversity but can reduce peering. One singular route collector, namely `rrc25`, already enables us to find a majority of heterogeneous ASes due to 77 peers. Diversity is the key for valuable measurements. Information from additional route collectors, peerings and prefixes allows a better view and differentiate between detours. Preserving this information to support future research is important.

5.6 Summary

The Internet is a complex network that changes its routing, peering, and policies at a fast pace. In this dissertation, we focused on measuring the adoption of heterogeneity—using different routes for the same prefix. We use all route collectors to find detours between pairs of ASes and find that heterogeneity is increasing. Although all regions are becoming more heterogeneous, Asia Pacific (increase of 298%), Europe (172%), and South America (480%) stand out. Especially ASes that describe themselves as operational in a regional scope (183%) also increase in heterogeneity.

As expected, ASes focusing on infrastructure as their business model, such as ISPs (271%) and NSPs (175%), adopt heterogeneity in their networks. Surprisingly, CDNs (323%) also show a growing trend. This is the case since PeeringDB offers fixed occupations for ASes, but CDNs are diverse in their business model. We also found that for all ASes with an ASRank of 500 or better, only 40 are not heterogeneous. These ASes are also likely heterogeneous but can not be detected due to a lack of VPs.

Lastly, we evaluate the impact of current route collector setups and how they enable us to find heterogeneity. We found that even a single route collector (rrc25) allows us to find 532 out of 725 heterogeneous ASes in December 2024. Whereas one route collector with many peers enables us to find heterogeneity, a singular AS peering with many route collectors is not. We argue that route collectors should be diverse in locations.

Chapter 6

Extending Existing BGP Data with Looking Glasses

The previous chapters show that BGP data is inconsistent (Chapter 3) and that more VPs are required to improve BGP studies (Chapter 4 and Chapter 5). As such, the current route collector setup is limited. Further, the BGP data from these route collectors does only show the "public view" of BGP—configurations and customized BGP attributes stay hidden. This is however not the case for LGs. Indeed, LGs extend our "public" view of the Internet and distribute important information that can be used to make methodologies more robust towards the Internet's heterogeneity.

So far, we checked how reliable BGP data sources are and how many ASes defy the popular assumption that the Internet is a homogeneous network. In this chapter, we try to answer the last part of our overarching research question: can we use LGs to extend BGP data sources? We do so by building, to the best of our knowledge, the largest public LG dataset. The analysis in this chapter was performed at the time of submission of the corresponding paper, nevertheless, we continued collecting data for several more months. We adapted the dataset statistics to reflect its current state but left the analysis untouched. Our contributions for this chapter can be summarized as follows:

- We collect BGP attributes from more than 169 LGs in 175 ASes from 931 routers via scraping the LGs. Hereby, the difficulties relate to the non-uniformity of the LGs—most interfaces differ, the fluctuating accessibility of the LGs, as well as the different output formats. To overcome this we combined manual configuration with an automated scraping process followed by careful post-processing and manual checks.
- Our current dataset covers over six months of continuous data collection every 4 hours and contains over 57 M BGP routes. We describe both our collection pipeline as well as initial analysis results which focus on route diversity for ASes with multiple LGs. We find that (in February 2024) up to 43% of these ASes use diverse routes to at least one of their peers. Routes can differ in Local Preference 40%, AS

paths 37%, or BGP communities 41%. Note that routes can differ in more than one aspect.

6.1 Building a LG Dataset

The challenges in putting together a LG dataset include, but are not limited to, the fact that (i) most LGs use different forms with varying mandatory and optional elements (ii) some LGs offer access to multiple routers in different locations (iii) some LGs offer only subsets of the queries (iv) LGs rate limit queries and restrict for which prefixes information can be accessed (v) the output format differs and may only include subsets of the route attributes (vi) LG availability and stability is limited.

Dataset: In this dissertation, we overcome the above-mentioned restrictions and gather an extensive LG dataset. Indeed, our dataset is not restricted to a one-time snapshot but rather offers a longitudinal view as we extend our dataset every four hours with another sample. Our dataset covers 169 LGs, 175 ASes, and 604 routers. To do this we developed an automated scraper that emulates a browser and queries the LG. The returned routes are then parsed to extract all available BGP attributes. Lastly, we augment our dataset with meta-information such as the router’s IP address.

LG Dataset Usage: ASes use BGP to reflect their business relationships with other ASes, i.e., customer-to-provider (c2p), provider-to-customer (p2c), peer-to-peer (p2p), or more complex ones, such as sibling relationships. Since ASes often neither publicly announce their neighbors nor their business relationships, these may have to be inferred. The first problem is called the Internet AS-level topology inference the latter is referred to as AS relationship inference.

The state-of-the-art inference method by Luckie et al. [82] uses an algorithm that relies on observed AS paths and a set of assumptions about AS relationships in the wild. Khan et al. [69] and Ahmed et al. [1] use LGs to augment their data and are able to discover thousands of additional p2p links. Both rely on BGP route queries to collect information about the routing policies inside the AS operating the LG. Giotsas et al. [50] use BGP communities to infer p2p links. They use LGs to validate their results.

Chang et al. [25] use the traceroute feature of LGs to perform traceroute queries from many ASes to detect missing BGP connectivity. The deployment of multipath BGP (M-BGP) was investigated by Li et al. [77, 78]. They use BGP summary queries to gather a list of peer ASes and BGP route queries to monitor M-BGP deployments. One common theme of these studies is that they are one-shot studies.

Other research focuses on discovering LGs. For example, Zhuang et al. [133] develop a classification-based method to find obscure LGs. They discovered close to 1K LGs servers in the wild that were previously not known. Their work facilitates many other studies, including ours. Giotsas et al. [48] focus on developing a LG query tool that covers many LGs. They combine automated classification with manual inspection to overcome the vast differences in LG web interfaces and outputs. Their tool allows, in principle, online access to 297 ASes and 1691 routers. However, over time the number of available LGs decreased. By now, Jan 2024, only 112 LGs seem to allow BGP queries according to the Periscope web site. Moreover, due to rate limiting, only a small subset—typically less than 25, see, e.g., [119]—of the total LG servers can be used. The declining numbers of LGs, as well as the lack of a longitudinal view, motivated us to collect the dataset presented in this dissertation. Therefore, our focus lies on collecting BGP data for a persistent set of prefixes from a large number of LGs across time.

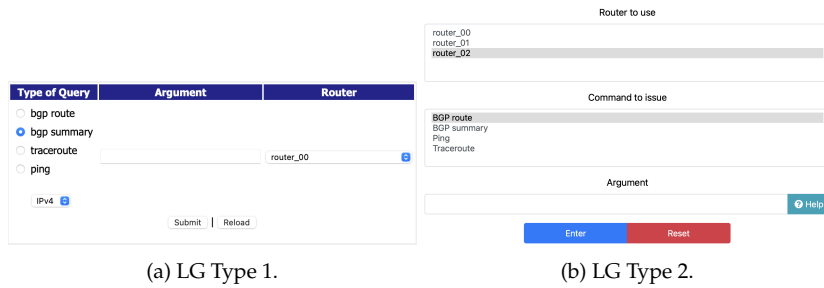


Figure 6.1: LG web interface examples.

Structure: This chapter is structured as follows. In Section 6.2 we discuss the ethical concerns that our work may raise. Next, we explain how we gather our dataset in Section 6.3. Then, in Section 6.4 we give an overview of our dataset. In Section 6.5 shows some initial results based on our dataset. Finally, in Section 6.6 we conclude and outline future work.

6.2 Ethical Considerations

Firstly, using a scraper automatically raises ethical concerns (*i*) regarding the load it introduces and (*ii*) whether operators even allow scraping. To address these concerns we honour the limits imposed by the website via the `robots.txt` file, we pace our requests, and limit them to a small set of prefixes. Moreover, we offer LG server operators to opt-out from our scans via (*i*) a blacklist and (*ii*) the ability to request the removal of all collected information from the dataset. Together, this allows us to collect valuable data for researchers while minimizing the impact on the LGs.

Secondly, some may consider the collected data as private data as it captures the effects of routing configurations. However, we note that the LGs make this data accessible without restrictions. Still, adding a longitudinal as well as latitudinal view may increase privacy as well as security concerns. For example, the longitudinal view enables to study changes in, e.g., business relationships. Thus, we decided to make the dataset as well as the scripts and patterns available upon request on a per-project basis using a license that gives us the ability to remove LG data upon an operator’s request. Thus, the usage license limits the ability to further distribute the dataset.

6.3 Methodology

LGs in the Internet use many different implementations. Many of the LGs (42%) have a visual appearance similar to the anonymized screenshot of the LGs shown in Figure 6.1a. These use the software `LG` [38]. Others (22%), share the look and feel of the LGs shown in Figure 6.1b. These use the software `looking_glass` [52]. Other ASes implement their own LG which leads to a large variety of LG web interfaces with highly varying HTML code.

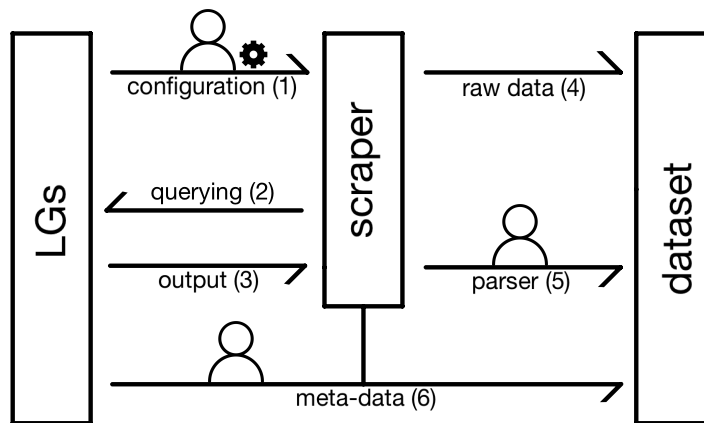


Figure 6.2: Overview of data collection pipeline.

ASN:

```
URI: http://lg.xyz.net
queries:
  show_route: //input[@name="route"]
  show_summary: //input[@name="summary"]
  ping: //input[@name="ping"]
  traceroute: //input[@name="traceroute"]
parameter: //input[@name="parameter"]
routers: //select[@name="routers"]
submit: //button[@id="submit"]
output: //pre
```

Figure 6.3: Example LG configuration in YAML format.

6.3.1 Overview of Data Collection Process

To address the large variability, we use a configuration-based scraper to interact with the LG web interface, as well as a configuration-based parser to collect the data. More precisely, we first create the configurations (Step 1) for (i) the scraper to locate and interact with the LG web interface (ii) for parsing the output, as well as for collecting the meta-data. Then the scraper periodically queries the LG by submitting a filled form and collecting the content of the output (Steps 2 and 3). We then augment our dataset with both the raw output (Step 4), as well as the parsed output (Step 5), which consists of the BGP route including all available BGP attributes. Finally, we augment our data with meta-data (Step 6), which is partially collected beforehand and partially collected from the LG at query time. Figure 6.2 gives an overview of the data collection pipeline.

6.3.2 Process Pipeline

Step 1: We generate the configurations manually, as we noticed too many differences, many of which are often small but subtle. The configuration for the scraper is captured in YAML. See Figure 6.3 for an example. The configuration for the parser consists of 14 patterns to be able to extract the BGP attributes. The extraction of the meta-data is discussed later.

Step 2: We base our scraper on Playwright [88] which is a debugging tool that can emulate a browser. It uses the XML Path Language (XPath) [35] for interactions with the website. The first step for each query is to open the base LG URI in Playwright. Then, the scraper uses its configuration to locate the appropriate HTML elements via their XPaths. Next, the text inputs are filled, the selection fields chosen, and / or the input buttons clicked to gather the desired information.

While the XPath elements differ the logic is consistent and includes the following steps:

1. Choose query type: input, button, or select.
2. Choose IP version (if present): input, button, or select.
3. Enter prefix: text input.
4. Select a router (if present): input, button, or select.
5. Click submit button: input or button.
6. Record the output: code, span, div, ...

Unfortunately, the time required by the LG to answer the queries differs. While most respond within 30 seconds, some may take unpredictably longer. Thus, we use a timeout of 30 seconds and record the query as disrupted. If the load of the base page fails the query is recorded as failed.

Step 3 and 4: The output of the query is either shown in a separate output window or on refresh on the website itself. We use Playwright to record the raw output for each query using a separate file. This enables debugging and, if necessary, an augmentation of the parsing process. Each file contains information about the LG, the router, and the query.

Step 5: Next, we parse the output to extract the BGP attributes using a small Python script which heavily relies on the above-mentioned patterns. Using Python, we first check using specific markers, e.g., unique strings such as `A\s+V\s+Destination\s+P Prf\s+Metric 1\s+Metric 2\s+Next hop\s+AS path` which of the 14 base patterns applies. Then we use regular expressions assigned to the pattern to extract BGP attributes for each route. These include BGP prefix, AS path, Local Preference, MED, and BGP communities. We also extract route meta-information such as the age of the route or the date of the last route update. We then add the route with all its attributes to our dataset, which consists of one CSV file per snapshot. Note, the above works even if multiple routes are included in one output. To account for changes in the output we check if (*i*) we are able to extract the expected data or if it is, e.g., empty, and if the extracted data adheres to the expected output type, e.g., if an ASN is indeed a number in the expected value range. Moreover, we check for consistency, e.g., if on day X we can extract attribute A. We expect that it is also available on day X+1. Any errors are noted and allow us to double-check and update our patterns. These can then be applied to the raw output which is also stored.

Step 6: Next, we augment each BGP route entry with additional meta-data. The output file creation date, for example, is equal to the time a query ended which is added to the dataset. All routers in an LG's web interface have names that are assigned by the operators. We use the XPath locator for the router selection fields to collect them and add them to the dataset as well. While most LGs belong to a specific AS, some LGs offer access to routers in other ASes. This can be inferred via the router name, some text field entry, or via the selection optgroups of the router selection field. We manually created

a list of such routers and include this information in the dataset. Next, we add the IP address of the router that answered the query, if possible, since this may be helpful for IP-based geolocation methods. We extracted the IP addresses by using the ping/traceroute query interface of the LG and extracting the information via TCPdump [65] on one of our servers. Finally, we add the identifier of the parsing pattern that was used for data extraction to enable easier debugging.

Process: To gather one snapshot our scraper software uses one process per LG. To avoid spamming any of the LGs we add a 0.5-second spacing between queries. To avoid overusing the resources of the dedicated server running the scraper we limit the number of parallel scraping processes to 80. This ensures that we do not encounter memory or bandwidth limits. Still, this ensures that a single snapshot for, e.g., 35 BGP prefixes, across all LGs can be taken within a few minutes. Currently, there is an imbalance in runtime which is due to the different number of routers per LG which ranges from 2 seconds to a maximum of 1 hour. In the future, we will rebalance this by assigning routers per scraper instance. The drawback of this is that one LG may receive multiple queries within a relatively short time period.

6.3.3 Remarks

Creating the LG configurations and collecting the meta-data has been a multi-week process. One requires the configurations to develop the scraper but needs to run the scraper to test the configurations. Furthermore, we pace our calls, which increases the debugging time.

Another complication is that each LG's performance differs and that it even varies for routers in the same LG. The query time ranges from milliseconds to minutes. This implies that we had to eliminate some LGs due to their bad/variable performance.

Another aspect of LGs is that they are not high-priority services of ASes. LGs are usually not well maintained. Thus, there are many inconsistencies, even visible at a single LG. Routers, for example, might be added or removed to the LG without updating the web interface. This is visible in many outputs that our scraper yields. Overall, 69% of routers result in outputs, whereby the rest fails to answer the queries. Around 14% of those failing queries return some kind of error indication, 17% of them return general error messages, 9% return a closed file handle, and 26% show that the query in use is deprecated. The remaining 34% of the erroneous outputs show miscellaneous messages. Nevertheless, we collected 2.1M outputs with content so far and parsed a total of 10.2M BGP routes from them.

6.4 Dataset

#LGs left	Filter	Explanation
3800	/	/
2317	timeout	All URIs that result in a timeout.
1903	not working	All URIs that cannot be resolved or do not respond with an HTTP code of 200.
1273	no BGP	All URIs that do not contain the buzzword 'BGP'.
1080	no LG	All URIs that do not point to a Looking Glass, for example, 'bgp.he.net/...'
812	make unique	Remove URIs that point to the same LG.
608	no automation	All LGs that, either verbatim or via Captcha, do not allow automation.

Table 6.1: Filtering steps for the list of used LGs.

In this section, we give an overview of our dataset.

6.4.1 Looking Glasses

To create a list of possible LGs to be used for our dataset we used the multiple websites <https://www.bgp4.as/> [9], <https://www.bgplookingglass.com> [8], and <https://whois.ipinsight.io/looking-glass/> [64]. In addition, we added the collected LGs by Zhuang et al. [133]. Overall, our initial list consisted of 3800 LGs. However, not all of these are active, offer LG services, etc. As such we post-processed the list to identify possible LGs for our dataset. The corresponding filters are shown in Table 6.1 and resulted in a list of 608 possible LGs.

As of February 2024 we have prepared configurations for 205 of these 608. For 32 we encountered problems when using XPath. Accordingly, we excluded them. Moreover, 17 LGs no longer respond to our requests. These were excluded as well. This leaves us for now with a set of 156 for our dataset. Note, that including a LG in the dataset does not imply that we get answers from them for every possible snapshot.

6.4.2 Performed Queries

The focus of this dataset is on BGP. Accordingly, we focus on the BGP route and BGP summary queries. The BGP summary option does not require any argument and returns BGP route statistics, often including the number of active and inactive peers, the number of exchanged BGP routes, and traffic meta-data. Instead of further parsing this information, we only gather the raw output. BGP route queries require a prefix as argument and return BGP route information which includes the various BGP attributes. We store this information in raw format as well as the parsed output in a CSV format.

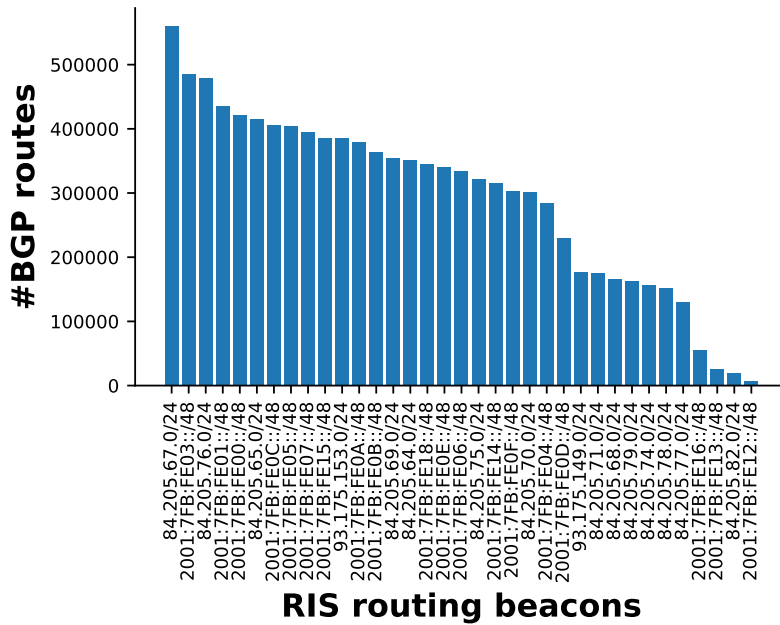


Figure 6.4: Number of collected BGP routes per beacon over the whole duration of data collection.

6.4.3 Used Prefixes

Regarding our choice of prefixes, we are looking for v4 as well as v6 prefixes that originate from different locations within the Internet and that can capture some dynamics of the Internet ecosystem. Given these requirements, we choose the RIPE RIS routing beacons [111].

The goal of this project is to take a snapshot of the stable routing system. We aligned our measurement period with the beacon schedule and take a snapshot every 4 hours starting 30 minutes after the beacons are announced or put differently 90 minutes before they are withdrawn. This leaves us with 90 minutes for the snapshot collection process which is sufficient to perform the roughly 33K queries for all 35 prefixes from all LGs selected for the dataset.

6.4.4 Dataset overview

Our data collection for the following analysis started on the 8th of December 2023 and lasted until the 26th of January 2024. Using our 155 LGs we have, in principle, access to 931 routers. Unfortunately, 377 never responded to our queries. Still, we did not remove them from the data collection process as they may be revived at some future time. 88 routers always responded with an error message, either due to login failures (5), command problems (25), or generic errors (58). Again we did not remove them for the same reason. 148 of the LGs were at some point unreachable by our scraper as well as 479 of the LG routers. 527 routers at some point responded with an error message, either due to login failures (45), command problems (72), or generic errors (410). Moreover, 22 LGs provide access to routers in multiple ASes.

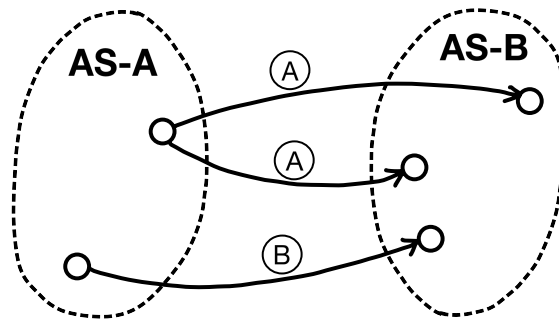


Figure 6.5: Example demonstrating route diversity from AS A.

Overall, we issued 9.3M queries. However, errors are common as discussed above and we do not expect this many data points. Still, a single data point can consist of multiple routes and Figure 6.4 shows the number of collected routes per beacon in our dataset. Notice the large differences. For four of the beacons (3 IPv6 and 1 IPv4 prefix) we were able to only collect a significantly small number of routes. Note, that two of these prefixes are from Africa. It appears that some ASes on the path to the LGs are filtering these prefixes. Interestingly, there appears to be no major difference between the other IPv4 and IPv6 prefixes. Overall, routes are rather stable. Our analysis shows that per LG router we on average see 5 different routes.

6.5 Example Analysis

To highlight what kind of analysis is enabled by the dataset presented in this paper we study route diversity. As depicted in Figure 6.5, we have route diversity if two different routers within an AS within a snapshot observe differences in the BGP attributes of the route for a specific prefix AS peer pair. Hereby, we distinguish between differences in (i) AS paths, (ii) Local Preferences, and (iii) BGP communities.

6.5.1 Setup

The first step of our analysis is to identify those ASes for which we have multiple LG routers using the meta-data. Then we identify all possible AS neighbors using the AS path BGP attribute from the dataset. This gives us a list of AS pairs with multiple LG routers. Next, we iterate over all prefixes and snapshots. If for a prefix we do not have information for that prefix from at least two LG routers and AS neighbors within the snapshot we ignore this sample.

Then we check for each prefix, AS pair, dataset snapshot, and considered BGP attribute if the values of the BGP attribute are equal or differ. Hereby, a difference in AS paths exists if the number of routes per prefix differs, if the AS paths differ, or if the router chooses a different best route. A difference in Local Preference exists if we get Local Preference information and there is a difference in the set of (AS path, Local Preference) pairs. A difference in BGP communities exists if we observe a difference in the set of BGP communities for the AS pair. We label an AS as heterogeneous if for any of the AS pairs we see diversity in the corresponding BGP attribute.

6.5.2 Results

Overall, we have multiple LG routers in 219 ASes and can, thus, include them, in principle, in our analysis. However, we only get reliable feedback from 148 of these ASes. For the other ASes, at least one of the LG routers did not respond reliably. For these 148 ASes, we consider 4118 AS pairs.

Overall, we find differences in the AS path for 55 of the 148 ASes (37%) we study. The reason for this diversity relates to the fact that ASes typically consist of multiple routers in diverse locations. As such BGP selection is biased by Local Preference, location, routing policy, etc.

We get Local Preference information for 142 of the LG ASes. Of these 51 (36%) show AS path diversity and 56 (40%) have differences in Local Preferences. This indicates that Local Preferences might be one of the elements used in the routing policies of these ASes.

BGP communities are tags that can be added to BGP routes and allow ASes to tag routes with information. There is a huge variety of possible BGP community values. As such we expect there to be significant differences. However, for the 88 LG ASes where the BGP community attribute is included, we only see differences for 36 ASes (41%).

6.6 Summary

Multiple ASes in the Internet host LGs, most of which provide internal information about their BGP routes. We develop a web scraper that, using periodically announced prefixes, queries a list of LGs, and records their outputs. The recorded output is parsed into CSV files, each file corresponds to one snapshot. We will make our dataset available upon request. Overall, we address the following three questions in this chapter:

How can such a dataset be created? Our longitudinal dataset relies on prefixes that are announced periodically: the RIPE RIS routing beacons. The scraper is aligned with the announcements of the beacons and uses manually tuned (per LG) configurations to gather the data. This data is then parsed using regular expressions. The output is checked for possible errors. Reoccurring errors lead to readjustments of either the configurations or the regular expressions.

What does the data, raw and parsed, look like? The raw data is a collection of files which contain the output of the scraper. This data is parsed by a Python script using various regular expressions. Its output contains (sometimes multiple) BGP routes, each of which is added to our dataset. Overall, a new snapshot is added to the dataset every 4 hours. We augment each BGP route with meta-information that was either collected beforehand or extracted via scraping. This includes, for example, router names and router IP addresses, but also the time a query was performed.

What is the use case of this dataset? The dataset enables a view of an AS's internal BGP routes (more than just the distributed route) and their attributes. This enables, for example, validation of BGP attribute inference methods. Another use-case is the inference of various hidden links in the Internet, the AS topology inference. In addition, the time-based information contained in our dataset enables the analysis of changes in business relationships. In this chapter, we use our dataset to look into differences in BGP routes inside one AS. We find that 37% of ASes differ in the distribution of their AS paths, 40% differ in their Local Preferences, and 41% differ in BGP communities. Those first two results may be expected, however the latter is surprising to us. We argue that the

performed analysis in this chapter is rather a use-case example of this dataset and not a standalone methodology to find heterogeneous ASes (as compared to Chapter 4 and Chapter 5).

Chapter 7

Discussion

In this dissertation, we question the quality of well-known data sources, quantify number of heterogeneous ASes, and extend the currently available data sources by scraping LGs. In Chapter 3, we investigate the data quality of a major source for BGP studies: the various route collectors. Chapter 4 and Chapter 5 propose the notion of heterogeneity—ASes that use different routes for the same prefix at the same time. These chapters also quantify the number of heterogeneous ASes in the Internet. Finally, Chapter 6 introduces a LG dataset that extends our view of BGP.

7.1 Summary

Route Collector Inconsistencies. Route collectors are widely used in studies about BGP. Nevertheless, most studies about BGP do not question the quality and integrity of route collector data. Our findings from Chapter 3, however, show some underlying inconsistencies. In the course of this study, we analyzed 7 years of sampled route collector data. This gave us some insights into the existing inconsistencies and possible reasons for them. However, we expect there to be many more temporal effects to be observable when investigating all data. Although the BGP data cannot be re-collected, it is important to future-proof the current data collection services. We contacted the operators of the two major route collector services and informed them. In May 2025, RIPE RIS launched OS upgrades for all of their route collectors. Our analysis also reveals that BGPStream misses to announce URIs especially after 2016.

Heterogeneity in the Internet. In Chapter 4 and Chapter 5, we make the first steps to quantify the number of ASes that act heterogeneously. We propose and implement two different methods to find such heterogeneous ASes and apply them to historical data. Our results show that all regions became more heterogeneous in the last decade. Asia Pacific (increase of 298%), Europe (172%), and South America (480%) stand out. Especially ASes that operate in a regional scope (183%) increase in heterogeneity. We further find that the most diverse RIR is RIPE, with 18% of observed ASes being heterogeneous, followed by ARIN with 10% and AFRINIC with 7%. Lastly, we identify that more than 12% of

observed ASes choose a longer path as their preferred ('best') path, calling into question the reliability of these assumptions.

We find that more ASes in the Internet are becoming heterogeneous on a yearly basis. As such, we observe an increase of 198% in the number of heterogeneous ASes since 2014. Among the observed heterogeneous ASes, 21% are NSPs, 11% are Internet Service Providers, and 7% are Content Delivery Networks. The increase is especially the case for ISPs (271%) and NSPs (175%), but also surprisingly for CDNs (323%). Our results show that all ASes with an ASRank of 75 or better are heterogeneous. In the case of the Customer Cone Size, all ASes with more than 10k ASes are heterogeneous.

Looking Glass Dataset. In Chapter 6, we collect LG data to extend existing data sources. Hereby, we develop a web scraper that periodically queries a list of LGs, and records their outputs for a set of prefixes. We further parse the output and create one snapshot every 4 hours. Overall, we describe both our collection pipeline as well as initial analysis results which focus on route diversity for ASes with multiple LGs. We find that up to 43% of these ASes use diverse routes to at least one of their peers. Using this dataset, we find that 37% of ASes differ in their AS paths for the same prefix, 40% have different Local Preferences, and 41% differ in BGP communities. While the former is expected the latter is surprising.

7.2 Conclusion and Future Directions

Most studies about routing in the Internet rely on assumptions to create models and make methodologies feasible. Assumptions, however, build an abstraction between the world that we study and the results that we observe. It is in our—the researchers'—responsibility to stay as close to the system we measure as possible.

In this dissertation, our aim is to shed light on the datasets used in understanding routing in the Internet, to identify limitations, and to improve the available datasets for future empirical work. To reiterate, the main research question in this dissertation is:

To what extent are currently used BGP datasets suited for, e.g., Internet topology research: which limitations do they have, how have these limitations impacted prior work, and how can these datasets be improved to form a more robust foundation of empirical work in Internet topology research?

To answer this overarching research question, we address three sub-questions that describe challenges for BGP studies.

(1) How consistent is the currently available BGP data from the major route collector projects and which biases are introduced by BGP data distribution software?

Data collection, distribution, filtering, and processing influence research results. Whereas distribution, filtering and processing lies in the hands of researchers, data collection might not. In this dissertation, we investigated BGP route collectors and their data. We find that this data is inherently flawed—it is not consistent across snapshots.

Although our study in Chapter 3 suggests that there is a minor number of inconsistencies on average, research results can be biased. Route collectors can face resource depletion that causes non-trivial spikes in inconsistencies. Studies that use data from such times are especially affected. We argue that aside from cleaning BGP data, it is also important to find and fix the root causes for these inconsistencies.

Future work should conduct a more in-depth evaluation of inconsistencies in BGP route collector data, also assessing if and how these inconsistencies may have influenced results of prior studies. Furthermore, besides BGPStream, other BGP data broker implementations—e.g., BGPKIT [10]—exist. We strongly recommend checking the reliability of these implementations as well. In any case, authors need to remain vigilant when relying on third party data sources and should always, at least, evaluate the internal validity using spot checks.

Furthermore, fetching the data can be either done manually or automatically using BGP data distribution software. We investigate the tool BGPStream—a tool that is widely used by researchers. Hereby, we find that BGPStream’s broker misses to advertise data. These results suggest that researchers must make sure that they have all data to not introduce biases to their results.

The final takeaway for this project is that questioning and revalidating research results is important—especially if previously hidden flaws are uncovered. For the future, we propose that the quality of BGP data should be questioned and investigated before usage. This is an additional step that comes before cleaning / filtering the announcements. We also propose that tools such as CAIDA’s BGPStream and its broker should be monitored more closely. Their broker service should be updated to include all snapshots, for the past and the future.

(2) Can we use the available data to quantify the number of ASes in the Internet that do act in a non-homogeneous (heterogeneous) manner?

Whereas data can be cleaned and fixed, this does not apply to the aggregated results of methodologies. A majority of BGP studies that focus on, e.g., Internet topology discovery, assume that each AS is a homogeneous network. However, we find that this assumption is not necessarily true. At least 14.3%—a non-trivial fraction—of ASes that announce at least one prefix are heterogeneous. Especially ASes with a high ASRank and a large CC tend to be heterogeneous. Although our studies are limited by sampling and VPs, we argue that the number of heterogeneous ASes that we find is a lower bound.

The results of our studies in Chapter 4 and Chapter 5 show the importance of adapting the existing state-of-the-art BGP studies to include the notion of AS heterogeneity. We expect that many studies that simulate / emulate BGP would show vastly different results. This also includes studies that infer RPKI deployments in the wild—especially if they are based on simulations / emulations. For the future, researchers must keep the evolution of routing in the Internet in mind. The next step is to adapt existing methodologies and reproduce their findings. The current state-of-the-art methods, e.g., CAIDA’s ASRank [18], Internet topology inference [82], are still based on this assumption. It is not clear how much bias this assumption introduces to research results that are based on these methods. For the future, it is important to further question and research the impact of the assumption that ASes in the Internet are homogeneous networks.

The key takeaway for this sub-question is that it is important to build a scalable model of a heterogeneous Internet. Such a model helps in our understanding of routing via BGP in the Internet. Aside from the gained knowledge, operators could better predict routing and optimize their networks. While the Internet is becoming a more complex system, our methodologies need to adapt as well.

(3) Can we use Looking Glasses to extend the existing BGP data?

The amount of available VPs can be extended by using LGs. Previous studies have already used LGs since decades. However, a minority of such studies have published their collected data. This makes the usage of LGs inconsistent across different studies. We combat this by collecting a longitudinal LG dataset as present in Chapter 6. Overall, this dataset simplifies the work of researchers studying BGP by enhancing their view of the Internet.

LGs offer queries to show BGP routes together with their assigned attributes. There are many BGP studies, e.g., BGP policy inference, BGP path prediction, M-BGP and RPKI deployment analysis, that greatly benefit from LG data. A unified LG data source, thus, is the first step towards comparable results. It is on researchers to adapt and improve their methodologies (if applicable) to make use of private attributes.

In principle, our configurations for the LGs can be extended to include other query options such as ping and traceroute commands. This allows us to use the LGs as an active measurement tool and can add VPs. Another future use case of our scraper scripts concerns RPKI. By using RPKI beacons it may be possible to evaluate the delay of validation state changes.

It is important for the future to continue incorporating more diverse BGP data sources. A diverse view of the Internet allows us researchers to improve upon existing models and methodologies. We thus aim to extend our dataset in the future by automating more LGs. Aside from web scraping, we are planning to contact operators in the Internet to gain immediate access to BGP attributes. An immediate cooperation benefits both sides: operators have unified access to routing information and researchers have more data.

List of Abbreviations

AS Autonomous System	v, 1
ASN Autonomous System Number	5
ASPA Autonomous System Provider Authorization	13
BGP Border Gateway Protocol	v, 1
CA Certificate Authority	13
CC Customer Cone	5
CDF Cumulative Distribution Function	44
CDN Content Delivery Network	1
FIB Forwarding Information Base	10
GRT Global Routing Table	11
iBGP internal Border Gateway Protocol	10
IGP Interior Gateway Protocol	9
ISP Internet Service Provider	1
IXP Internet Exchange Point	13
LG Looking Glass	v, 2
MED Multi-Exit Discriminator	2
MRT Multi-threaded Routing Toolkit	11
NS Network Services	34
NSP Network Service Provider	5
RIB Routing Information Base	4
RIR Regional Internet Registry	13
RIS Routing Information Service	v, 3
ROA Route Origin Authorization	13
ROV Route Origin Validation	13

RP Relying Party Software.....	13
RPKI Resource Public Key Infrastructure	1
RS Route Server.....	10
VP Vantage Point.....	2

List of Figures

2.1	Example BGP topology.	9
3.1	Example of combining RIB_{T0} and updates to get RIB'_{Δ}	18
3.2	Number of misplaced route update messages.	20
3.3	Class A and B inconsistencies for Routeviews route collectors.	20
3.4	Example for filesize / resource problem correlation.	21
3.5	Inconsistencies vs. route-views3 and route-views.amsix resource starvation.	21
3.6	Number of BGPStream vs. archive URIs.	22
4.1	Example BGP update sequence.	27
4.2	Number of observed heterogeneous ASes over time.	30
4.3	Number and fraction of ASes with at least 1, 2, 10 inferences.	32
4.4	Number and fraction of ASes with at least 1, 2, 10 conflicts.	33
4.5	Inferences and conflicts across BGP beacons.	34
4.6	Inferences and conflicts across BGP route collectors.	35
4.7	Number and fraction of ASes with inferences / conflicts per AS type.	36
4.8	Number and fraction of ASes with inferences / conflicts per RIR.	37
4.9	Number of total / repeated inferences over time.	38
5.1	Example AS topology that shows routing heterogeneity.	42
5.2	Number of heterogeneous ASes based on constant BGP sessions.	43
5.3	Number of heterogeneous ASes per RIR.	44
5.4	Number of heterogeneous ASes per occupation.	45
5.5	CDF of heterogeneous ASes and detours (ASRank).	46
5.6	CDF of heterogeneous ASes and detours (CC).	46
5.7	Boxplots of heterogeneous ASes and detours (ASRank and CC).	47
5.8	Boxplot for heterogeneous ASes with X peers.	48
5.9	Boxplot for heterogeneous peers that feed X ASes.	48
5.10	CDF of heterogeneous ASes and detours (route collector).	49
5.11	CDF of heterogeneous ASes and detours (BGP session).	50
5.12	CDF of heterogeneous ASes and detours (prefix).	50
6.1	LG web interface examples.	54
6.2	LG scraper pipeline.	55
6.3	Example LG configuration.	55
6.4	Collected routes per beacon prefix.	59

6.5 Example to demonstrate AS heterogeneity. 60

List of Tables

4.1	Overview of inferences / conflicts.	30
5.1	Fraction of heterogeneous ASes (ASRank and CC).	46
6.1	Filtering steps for the list of used LGs.	58

Bibliography

- [1] R. Ahmed, A. Khan, and M. Rizwan. Collection of Autonomous System Level Topology Using Looking Glass Servers. *Pakistan Journal of Science (PJS)*, 2019.
- [2] B. Al-Musawi, P. Branch, and G. Armitage. Detecting BGP Instability Using Recurrence Quantification Analysis (RQA). In *IEEE International Performance Computing and Communications Conference (IPCCC)*, 2015.
- [3] albertogarciam. No updates retrieved for certain collectors/periods, 2019. Available at: <https://github.com/CAIDA/bgpstream/issues/78> Last accessed: 2025-03-17.
- [4] T. Alfroy, T. Holterbach, T. Krenc, K. Claffy, and C. Pelsser. The Next Generation of BGP Data Collection Platforms. In *Proc. ACM SIGCOMM*, 2024.
- [5] H. An, Y. Na, H. Lee, and A. Perrig. Resilience Evaluation of Multi-Path Routing against Network Attacks and Failures. *Electronics*, 2021.
- [6] T. Arnold, J. He, W. Jiang, M. Calder, I. Cunha, V. Giotsas, and E. Katz-Bassett. Cloud Provider Connectivity in the Flat Internet. In *Proc. ACM Internet Measurement Conference (IMC)*, 2020.
- [7] A. Azimov, E. Bogomazov, R. Bush, K. Patel, and K. Sriram. Route Leak Prevention and Detection Using Roles in UPDATE and OPEN Messages. RFC 9234, IETF, May 2022.
- [8] BGP Looking Glass. BGP Looking Glass. Available at: <https://www.bgplookingglass.com> Last accessed: 2025-06-06.
- [9] BGP4. BGP Looking Glasses for IPv4/IPv6, Traceroute & BGP Route Servers. Available at: <https://www.bgp4.as/looking-glasses/> Last accessed: 2025-06-06.
- [10] BGPKIT LLC. BGPKit. Available at: <https://bgpkit.com/> Last accessed: 2025-03-07.
- [11] L. Blunk, M. Karir, and C. Labovitz. Multi-Threaded Routing Toolkit (MRT) Routing Information Export Format. RFC 6396, IETF, Oct 2011.
- [12] T. Böttger, G. Antichi, E. L. Fernandes, R. di Lallo, M. Bruyere, S. Uhlig, and I. Castro. The Elusive Internet Flattening: 10 Years Of IXP Growth. *arXiv e-prints*, 2018.
- [13] R. Bush, K. Patel, P. Smith, and M. Tinka. Policy Based on the Resource Public Key Infrastructure (RPKI) without Route Refresh. RFC 9324, IETF, Dec 2022.
- [14] R. Bush, C. Pelsser, K. Patel, and O. Maennel. BGP Decision Statistics: A First Experiment. Available at: <https://iepg.org/2010-11-ietf79/101107.iepg-statistics.pdf> Last accessed: 2025-08-25.

- [15] M. Caesar and J. Rexford. BGP routing policies in ISP networks. *IEEE Network Magazine (IEEE Netw.)*, 2005.
- [16] CAIDA. BGPReader. Available at: <https://bgpstream.caida.org/docs/tools/bgpreader> Last accessed: 2025-05-11.
- [17] CAIDA. BGPStream Broker API. Available at: <https://bgpstream.caida.org/docs/api/broker> Last accessed: 2025-05-12.
- [18] CAIDA. Caida AS Rank. Available at: <http://as-rank.caida.org/> Last accessed: 2025-04-14.
- [19] CAIDA. Path Sanitization and Realtime Processing. Available at: <https://github.com/CAIDA/bgpstream-tma-phdschool/blob/master/exercise-5-path-sanitization/README.md?> Last accessed: 2025-08-25.
- [20] CAIDA. PyBGPStream. Available at: <https://github.com/caida/pybgpstream> Last accessed: 2025-05-06.
- [21] CAIDA. University of Oregon Route Views Project Collector Peers. Available at: <https://archive.routeviews.org/peers//peering-status.html> Last accessed: 2025-05-06.
- [22] CAIDA. University of Oregon Route Views Project Collectors. Available at: <https://www.routeviews.org/routeviews> Last accessed: 2025-05-06.
- [23] CAIDA. University of Oregon Route Views Project Prefix to AS mapping. Available at: <https://www.caida.org/catalog/datasets/routeviews-prefix2as/> Last accessed: 2025-05-06.
- [24] CAIDA. University of Oregon Route Views Project, 1997. Available at: https://gmi3s.caida.org/reports/routeviews_architecture_new_prototype.pdf Last accessed: 2025-06-05.
- [25] D. F. Chang, R. Govindan, and J. Heidemann. Locating BGP missing routes using multiple perspectives. In *Proc. ACM SIGCOMM Workshop on Network Troubleshooting (NetT)*, 2004.
- [26] H. Chang, R. Govindan, S. Jamin, S. J. Shenker, and W. Willinger. Towards Capturing Representative AS-level Internet Topologies. *Computer Networks (Comput. Netw.)*, 2004.
- [27] H. Chang, S. Jamin, and W. Willinger. Inferring AS-level Internet Topology from Router-Level Path Traces. In *Scalability and Traffic Control in IP Networks (STC-IP)*, 2001.
- [28] W. Chen, Z. Wang, D. Han, C. Duan, X. Yin, J. Yang, and X. Shi. ROV-MI: Large-Scale, Accurate and Efficient Measurement of ROV Deployment. In *Proc. Network and Distributed System Security Symposium (NDSS)*, 2022.
- [29] S. Cho, R. Fontugne, K. Cho, A. Dainotti, and P. Gill. BGP Hijacking Classification. In *Proc. Network Traffic Measurement and Analysis Conference (TMA)*, 2019.
- [30] J. Choi, J. H. Park, P.-c. Cheng, D. Kim, and L. Zhang. Understanding BGP Next-hop Diversity. In *Proc. IEEE Int. Conference on Computer Communications (INFOCOM Workshops)*, 2011.

- [31] T. Chung, E. Aben, T. Bruijnzeels, B. Chandrasekaran, D. Choffnes, D. Levin, B. M. Maggs, A. Mislove, R. v. Rijswijk-Deij, J. Rula, et al. RPKI Is Coming Of Age: A Longitudinal Study Of RPKI Deployment And Invalid Route Origins. In *Proc. ACM Internet Measurement Conference (IMC)*, 2019.
- [32] Cisco. Cisco BGP Best Path Decision Process. Available at: <https://www.cisco.com/c/en/us/support/docs/ip/border-gateway-protocol-bgp/13753-25.html> Last accessed: 2025-04-23.
- [33] Cisco Press. Explore the Network, 2024.
- [34] L. Cittadini, S. Vissicchio, and B. Donnet. On The Quality Of BGP Route Collectors For iBGP Policy Inference. In *IEEE/IFIP Networking Conference (NETWORKING)*, 2014.
- [35] J. Clark and S. DeRose. XML Path Language (XPath). Available at: <https://www.w3.org/TR/1999/REC-xpath-19991116/> Last accessed: 2025-06-06.
- [36] Cloudflare. Is BGP safe yet?, 2025. Available at: <https://isbgpsafeyet.com/> Last accessed: 2025-04-23.
- [37] D. Cooper, E. Heilman, K. Brogle, L. Reyzin, and S. Goldberg. On The Risk Of Misbehaving RPKI Authorities. In *Proc. of the Twelfth ACM Workshop on Hot Topics in Networks (HotNets)*, 2013.
- [38] Cougar. LG Software by Cougar. Available at: <https://github.com/Cougar/lg> Last accessed: 2025-06-06.
- [39] W. Deng, W. Mühlbauer, Y. Yang, P. Zhu, X. Lu, and B. Plattner. Shedding Light on the Use of AS Relationships for Path Inference. *Journal of Communications and Networks (JCN)*, 2012.
- [40] X. Dimitropoulos, D. Krioukov, M. Fomenkov, B. Huffaker, Y. Hyun, K. Claffy, and G. Riley. AS Relationships: Inference And Validation. *ACM SIGCOMM Computer Communication Review (CCR)*, 2007.
- [41] T. Fiebig. Crisis, Ethics, Reliability & a Measurement. Network: Reflections on Active Network Measurements in Academia. In *Proc. of the Applied Networking Research Workshop (ANRW)*, 2023.
- [42] A. Flavel, O. Maennely, B. Chiera, M. Roughan, and N. Bean. CleanBGP: verifying the consistency of BGP data. In *IEEE Internet Network Management Workshop (INM)*, 2008.
- [43] R. Fontugne, E. Bautista, C. Petrie, Y. Nomura, P. Abry, P. Gonçalves, K. Fukuda, and E. Aben. BGP Zombies: An Analysis of Beacons Stuck Routes. In *Proc. Passive and Active Measurement (PAM)*, 2019.
- [44] R. Fontugne, A. Phokeer, C. Pelsser, K. Vermeulen, and R. Bush. RPKI Time-Of-Flight: Tracking Delays In The Management, Control, And Data Planes. In *Proc. Passive and Active Measurement (PAM)*, 2023.
- [45] S. Garcia-Jimenez, E. Magaña, D. Morató, and M. Izal. On the Performance and Improvement of Alias Resolution Methods for Internet Core Networks. *Annals of Telecommunications (Ann. Telecommun.)*, 2011.

- [46] A. García-Martínez and M. Bagnulo. Measuring BGP Route Propagation Times. *IEEE Communications Letters (CL)*, 2019.
- [47] Y. Gilad, A. Cohen, A. Herzberg, M. Schapira, and H. Shulman. Are We There Yet? On RPKI’s Deployment and Security. In *Proc. Network and Distributed System Security Symposium (NDSS)*, 2017.
- [48] V. Giotsas, A. Dhamdhere, and K. C. Claffy. Periscope: Unifying Looking Glass Querying. In *Proc. Passive and Active Measurement (PAM)*, 2016.
- [49] V. Giotsas, M. Luckie, B. Huffaker, and K. Claffy. Inferring Complex AS Relationships. In *Proc. ACM Internet Measurement Conference (IMC)*, 2014.
- [50] V. Giotsas and S. Zhou. Improving the Discovery of IXP Peering Links Through Passive BGP Measurements. In *Proc. IEEE Int. Conference on Computer Communications (INFOCOM) Workshops*, 2013.
- [51] GlobalNOC. GlobalNOC Routerproxy. Available at: <https://routerproxy.grnoc.iu.edu> Last accessed: 2025-06-06.
- [52] gmazoyer. LG Software by gmazoyer. Available at: <https://github.com/gmazoyer/looking-glass> Last accessed: 2025-06-06.
- [53] E. Gregori, A. Improta, L. Lenzi, L. Rossi, and L. Sani. BGP And inter-AS Economic Relationships. In *International Conference on Research in Networking (ICRN)*, 2011.
- [54] T. G. Griffin and G. Wilfong. An Analysis of BGP Convergence Properties. *ACM SIGCOMM Computer Communication Review (CCR)*, 1999.
- [55] J. Hawkinson and T. Bates. Guidelines for creation, selection, and registration of an Autonomous System (AS). RFC 1930, IETF, Mar 1996.
- [56] P. Hennen, P. Mani, and A. Feldmann. Our own Looking Glass Dataset Website, 2024. Available at: <https://inet-lg-dataset.mpi-inf.mpg.de> Last accessed: 2025-07-04.
- [57] P. Hennen, P. Mani, and A. Feldmann. Peeking Beyond the Best Route: An Extensive Dataset for Looking Glasses. In *Proc. IEEE Network Operations and Management Symposium (NOMS)*, 2024.
- [58] P. Hennen, C. Munteanu, and A. Feldmann. BGP AS Paths: Shorter Is Not Always Better. In *Proc. ACM SIGCOMM Workshop on Next-Generation Network Observability (NGNO)*, 2025.
- [59] T. Hlavacek, H. Shulman, N. Vogel, and M. Waidner. Keep Your Friends Close, but Your Routers Closer: Insights into RPKI Validation in the Internet. In *Proc. USENIX Security Symposium (Secur. Symp.)*, 2023.
- [60] K. Hoarau, P. U. Tournoux, and T. Razafindralambo. BML: an efficient and versatile tool for BGP dataset collection. In *Proc. IEEE ICC*, 2021.
- [61] T. Holterbach, T. Alfroy, A. Phokeer, A. Dainotti, and C. Pelsser. A System to Detect Forged-Origin BGP Hijacks. In *Proc. USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2024.

- [62] G. Huston. Commentary on Inter-Domain Routing in the Internet. RFC 3221, IETF, Dec 2001.
- [63] IJ Lab. Internet Health Report, 2025. Available at: <https://www.ihr.live/en>
Last accessed: 2025-06-05.
- [64] IPinsight. WHOIS - IPinsight Looking Glasses. Available at: <https://whois.ipinsight.io/looking-glass/> Last accessed: 2025-06-06.
- [65] V. Jacobson, S. Floyd, V. Paxson, and S. McCanne. TCPdump. Available at: <https://www.tcpdump.org> Last accessed: 2025-06-06.
- [66] Juniper. Juniper BGP Best Path Decision Process. Available at: <https://www.juniper.net/documentation> Last accessed: 2025-04-23.
- [67] J. Karlin, S. Forrest, and J. Rexford. Pretty Good BGP: Improving BGP By Cautiously Adopting Routes. In *Proc. IEEE International Conference on Network Protocols (ICNP)*, 2006.
- [68] S. Kastanakis, V. Giotsas, I. Livadariu, and N. Suri. Replication: 20 Years of Inferring Interdomain Routing Policies. In *Proc. ACM Internet Measurement Conference (IMC)*, 2023.
- [69] A. Khan, T. Kwon, H.-c. Kim, and Y. Choi. AS-Level Topology Collection Through Looking Glass Servers. In *Proc. ACM Internet Measurement Conference (IMC)*, 2013.
- [70] T. Kiso. mrtparse. Available at: <https://github.com/t2mune/mrtparse>
Last accessed: 2025-05-06.
- [71] T. Krenc, R. Beverly, and G. Smaragdakis. AS-level BGP Community Usage Classification. In *Proc. ACM Internet Measurement Conference (IMC)*, 2021.
- [72] T. Krenc, M. Luckie, A. Marder, and K. Claffy. Coarse-grained Inference of BGP Community Intent. In *Proc. ACM Internet Measurement Conference (IMC)*, 2023.
- [73] D. Krioukov, K. Claffy, M. Fomenkov, F. Chung, A. Vespignani, and W. Willinger. The Workshop On Internet Topology (WIT) Report. *ACM SIGCOMM Computer Communication Review (CCR)*, 2007.
- [74] J. Kristoff, R. Bush, C. Kanich, G. Michaelson, A. Phokeer, T. C. Schmidt, and M. Wählisch. On Measuring RPKI Relying Parties. In *Proc. ACM Internet Measurement Conference (IMC)*, 2020.
- [75] C. Labovitz, A. Ahuja, A. Bose, and F. Jahanian. Delayed Internet Routing Convergence. *IEEE/ACM Transactions on Networking (TON)*, 2001.
- [76] M. Lepinski and S. Kent. An Infrastructure to Support Secure Internet Routing. RFC 6480, IETF, Feb 2012.
- [77] J. Li, V. Giotsas, Y. Wang, and S. Zhou. BGP-Multipath Routing in the Internet. *IEEE Transactions on Network and Service Management (TNSM)*, 2022.
- [78] J. Li, S. Zhou, and V. Giotsas. Performance Analysis of Multipath BGP. In *Proc. IEEE Int. Conference on Computer Communications (INFOCOM) Workshops*, 2021.
- [79] W. Li, Z. Lin, M. I. Ashiq, E. Aben, R. Fontugne, A. Phokeer, and T. Chung. RoVista: Measuring and Analyzing the Route Origin Validation (ROV) in RPKI. In *Proc. ACM Internet Measurement Conference (IMC)*, 2023.

- [80] F. Lichtblau, F. Streibelt, T. Krüger, P. Richter, and A. Feldmann. Detection, Classification, and Analysis of Inter-Domain Traffic with Spoofed Source IP Addresses. In *Proc. ACM Internet Measurement Conference (IMC)*, 2017.
- [81] A. Lodhi, N. Larson, A. Dhamdhere, C. Dovrolis, and K. Claffy. Using PeeringDB to Understand the Internet Peering Ecosystems. *ACM SIGCOMM Computer Communication Review (CCR)*, 2014.
- [82] M. Luckie, B. Huffaker, A. Dhamdhere, V. Giotsas, and K. Claffy. AS Relationships, Customer Cones, and Validation. In *Proc. ACM Internet Measurement Conference (IMC)*, 2013.
- [83] H. V. Madhyastha, E. Katz-Bassett, T. E. Anderson, A. Krishnamurthy, and A. Venkataramani. iPlane Nano: Path Prediction for Peer-to-Peer Applications. In *Proc. USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2009.
- [84] Z. M. Mao, R. Bush, T. G. Griffin, and M. Roughan. BGP Beacons. In *Proc. ACM Internet Measurement Conference (IMC)*, 2003.
- [85] Z. M. Mao, D. Johnson, J. Rexford, J. Wang, and R. Katz. Scalable and Accurate Identification of AS-Level Forwarding Paths. In *Proc. IEEE Int. Conference on Computer Communications (INFOCOM)*, 2004.
- [86] Z. M. Mao, L. Qiu, J. Wang, and Y. Zhang. On AS-Level Path Inference. In *Proc. ACM SIGMETRICS*, 2005.
- [87] Z. M. Mao, J. Rexford, J. Wang, and R. H. Katz. Towards an Accurate AS-Level Traceroute Tool. In *Proc. ACM SIGCOMM*, 2003.
- [88] Microsoft. Playwright. Available at: <https://playwright.dev> Last accessed: 2025-06-06.
- [89] A. Milolidakis, T. Bühler, K. Wang, M. Chiesa, L. Vanbever, and S. Vissicchio. On the Effectiveness of BGP Hijackers That Evade Public Route Collectors. *IEEE Access*, 2023.
- [90] Minap News. Isolario Project Discontinuation, 2025. Available at: <https://www.minap.it/news/2021.html> Last accessed: 2025-04-17.
- [91] R. Morillo, J. Furuness, C. Morris, J. Breslin, A. Herzberg, and B. Wang. ROV++: Improved Deployable Defense against BGP Hijacking. In *Proc. Network and Distributed System Security Symposium (NDSS)*, 2021.
- [92] W. Mühlbauer, A. Feldmann, O. Maennel, M. Roughan, and S. Uhlig. Building an AS-Topology Model that Captures Route Diversity. *ACM SIGCOMM Computer Communication Review (CCR)*, 2006.
- [93] O. Muravskiy. Format change for index pages of RIS MRT data files and RIPE Atlas daily dataset dumps, 2022. Available at: <https://www.ripe.net/ripe/mail/archives/ris-users/2022-November/000731.html> Last accessed: 2025-05-06.
- [94] NRO. Development of the Regional Internet Registry System, 2010. Available at: <https://www.nro.net/development-of-the-regional-internet-registry-system> Last accessed: 2025-07-18.

- [95] P. Ongkanchana, R. Fontugne, H. Esaki, J. Snijders, and E. Aben. Hunting BGP Zombies in the Wild. In *Proc. of the Applied Networking Research Workshop (ANRW)*, 2021.
- [96] C. Orsini, A. King, D. Giordano, V. Giotsas, and A. Dainotti. BGPStream: A Software Framework for Live and Historical BGP Data Analysis. In *Proc. ACM Internet Measurement Conference (IMC)*, 2016.
- [97] Packet Clearing House. Packet Clearing House (PCH). Available at: <https://www.pch.net/> Last accessed: 2025-07-09.
- [98] PeeringDB. Peeringdb. Available at: <https://www.peeringdb.com> Last accessed: 2025-05-06.
- [99] R. Peruch. bgpscanner. Available at: <https://github.com/rfperuch/bgpscanner> Last accessed: 2025-05-11.
- [100] L. Prehn, P. Foremski, and O. Gasser. Kirin: Hitting the Internet with Distributed BGP Announcements. In *Proc. ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 2024.
- [101] L. Prehn, F. Lichtblau, C. Dietzel, and A. Feldmann. Peering Only? Analyzing the Reachability Benefits of Joining Large IXPs Today. In *Proc. Passive and Active Measurement (PAM)*, 2022.
- [102] J. Qiu, L. Gao, S. Ranjan, and A. Nucci. Detecting Bogus BGP Route Information: Going Beyond Prefix Hijacking. In *EAI International Conference on Security and Privacy in Communications Networks (SecureComm)*, 2007.
- [103] Y. Rekhter, T. Li, and S. Hares. A Border Gateway Protocol 4 (BGP-4). RFC 4271, IETF, Jan 2006.
- [104] RIPE NCC. BGPdump. Available at: <https://github.com/RIPE-NCC/bgpdump> Last accessed: 2025-03-07.
- [105] RIPE NCC. RIPE RIS. Available at: <https://www.ripe.net/analyse/internet-measurements/routing-information-service-ris/> Last accessed: 2025-05-06.
- [106] RIPE NCC. RIPE RIS Budget Cuts. Available at: <https://www.ripe.net/publications/docs/ripe-814/> Last accessed: 2025-05-06.
- [107] RIPE NCC. RIPE RIS Collector Cost. Available at: <https://labs.ripe.net/author/kistel/ripe-ncc-measurement-data-retention-principles/> Last accessed: 2025-05-06.
- [108] RIPE NCC. RIPE RIS Collector Peers. Available at: <https://www.ris.ripe.net/peerlist/all.shtml> Last accessed: 2025-05-06.
- [109] RIPE NCC. RIPE RIS Collectors. Available at: <https://ris.ripe.net/docs/route-collectors/> Last accessed: 2025-05-06.
- [110] RIPE NCC. RIPE STAT BGPplay. Available at: <https://stat.ripe.net/widget/bgplay> Last accessed: 2025-05-12.
- [111] RIPE NCC. RIS Routing Beacons. Available at: <https://ris.ripe.net/docs/routing-beacons/> Last accessed: 2025-05-06.

- [112] RIPE NCC. RIS Routing Beacons Historical Setup. Available at: <https://ris.ripe.net/docs/historical-routing-beacons/> Last accessed: 2025-05-06.
- [113] RIPE RIS. Ripe ris route collector start of rrc00, 1999. Available at: <https://data.ripe.net/rrc00/> Last accessed: 2025-07-09.
- [114] N. Rodday, G. D. Rodosek, A. Pras, and R. van Rijswijk-Deij. Exploring The Benefit Of Path Plausibility Algorithms In BGP. In *Proc. IEEE Network Operations and Management Symposium (NOMS)*, 2024.
- [115] J. Schlamp. Agree to Disagree: On the current state of BGP parsing, 2024. Available at: https://ripe89.ripe.net/wp-content/uploads/presentations/129-LEITWERT_2024-10-31_RIPE89_BGP-Parsing-Agree-to-Disagree.pdf Last accessed: 2025-06-05.
- [116] J. Schlamp. GMI-AIMS-3: Challenges in parsing BGP data, 2024. Available at: https://www.leitwert.net/doc/LEITWERT_2024-06-25_GMI-AIMS-3_FTLBGP.pdf Last accessed: 2025-06-05.
- [117] J. Schlamp, R. Holz, Q. Jacquemart, G. Carle, and E. W. Biersack. HEAP: Reliable Assessment of BGP Hijacking Attacks. *IEEE Journal on Selected Areas in Communications (JSAC)*, 2016.
- [118] P. Sermpezis, V. Kotronis, A. Dainotti, and X. Dimitropoulos. A Survey Among Network Operators on BGP Prefix Hijacking. *ACM SIGCOMM Computer Communication Review (CCR)*, 2018.
- [119] P. Sermpezis, V. Kotronis, P. Gigis, X. Dimitropoulos, D. Cicalese, A. King, and A. Dainotti. Artemis: Neutralizing bgp hijacking within a minute. *IEEE/ACM Transactions on Networking (TON)*, 2018.
- [120] P. Sermpezis, L. Prehn, S. Kostoglou, M. Flores, A. Vakali, and E. Aben. Bias in Internet Measurement Platforms. In *Proc. Network Traffic Measurement and Analysis Conference (TMA)*, 2023.
- [121] T. Shapira and Y. Shavitt. AP2Vec: An Unsupervised Approach for BGP Hijacking Detection. *IEEE Transactions on Network and Service Management (TNSM)*, 2022.
- [122] X. Shi, Y. Xiang, Z. Wang, X. Yin, and J. Wu. Detecting Prefix Hijackings In The Internet With Argus. In *Proc. ACM Internet Measurement Conference (IMC)*, 2012.
- [123] B. A. Silva Jr, P. Mol, O. Fonseca, I. Cunha, R. A. Ferreira, and E. Katz-Bassett. Automatic Inference of BGP Location Communities. *Proc. of the ACM on Measurement and Analysis of Computing Systems (POMACS)*, 2022.
- [124] F. Streibelt, F. Lichtblau, R. Beverly, A. Feldmann, C. Pelsser, G. Smaragdakis, and R. Bush. BGP Communities: Even More Worms In The Routing Can. In *Proc. ACM Internet Measurement Conference (IMC)*, 2018.
- [125] L. Subramanian, V. Roth, I. Stoica, S. Shenker, and R. Katz. Listen And Whisper: Security Mechanisms for BGP. In *Proc. USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2004.
- [126] Y. Tan, W. Huang, Y. You, S. Su, and H. Lu. Recognizing BGP Communities Based on Graph Neural Network. *IEEE Network Magazine (IEEE Netw.)*, 2024.

- [127] N. Tao, X. Chen, and X. Fu. AS Path Inference: From Complex Network Perspective. In *IEEE/IFIP Networking Conference (NETWORKING)*, 2015.
- [128] C. Testart, P. Richter, A. King, A. Dainotti, and D. Clark. To Filter or not to Filter: Measuring the Benefits of Registering in the RPKI Today. In *Proc. Passive and Active Measurement (PAM)*, 2020.
- [129] F. Wang and L. Gao. On inferring and characterizing Internet routing policies. In *Proc. ACM Internet Measurement Conference (IMC)*, 2003.
- [130] F. Wang and L. Gao. On Inferring and Characterizing Internet Routing Policies. *Journal of Communications and Networks (JCN)*, 2007.
- [131] Y. Wang and K. Zhang. Quantifying The Flattening Of Internet Topology. In *Proc. of the International Conference on Future Internet Technologies (CFI)*, 2016.
- [132] B. Zhang, R. Liu, D. Massey, and L. Zhang. Collecting the Internet AS-level Topology. *ACM SIGCOMM Computer Communication Review (CCR)*, 2005.
- [133] S. Zhuang, J. H. Wang, J. Wang, Z. Pan, T. Wu, F. Li, and Z. Zhang. Discovering Obscure Looking Glass Sites on the Web to Facilitate Internet Measurement Research. In *Proc. ACM Int. Conference on emerging Networking EXperiments and Technologies (CoNEXT)*, 2021.