

# MASTER PROJECT

## Measuring real broadband speeds using crowdsourcing data from the Internet Foundation

Author: Ivan Vallejo Vall

Tutor: Iñigo Herguera García

Data Science Program 2016/2017



## **Abstract**

Broadband networks have become a fundamental infrastructure and, as a result, there is a growing need for accurate data on some of their main characteristics. In particular, broadband speed measurements are a key input to several consumer, policy and regulatory decisions. This project reviews the state-of-the-art in broadband measurement platforms, both from a technical and statistical point of view, and applies it to assess whether crowdsourcing Internet data on broadband speed measurements can be used to extract robust inference on real broadband speeds. One hundred million observations from *Bredbandskollen*, the broadband speed measurement platform of the Internet Foundation in Sweden, are processed using big data methods. The analysis provides strong evidence of a selection bias and an unstable sample composition across years, thus indicating that the sample of test users may not be representative of the whole population of Sweden. Nevertheless, the analysis finds that there is a relative equality of speeds across regions and that there does not seem to be an urban/rural broadband speed divide in Sweden.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	1
1.2	Relevance . . . . .	2
1.3	Research questions . . . . .	4
<b>2</b>	<b>State of the art</b>	<b>6</b>
2.1	Measurement methodologies . . . . .	6
2.2	Statistical approach . . . . .	13
<b>3</b>	<b>Data</b>	<b>18</b>
3.1	Processing environment . . . . .	22
<b>4</b>	<b>Results</b>	<b>23</b>
<b>5</b>	<b>Conclusions</b>	<b>37</b>

## List of Figures

1	Big Data for Measuring the Information Society – ITU pilot projects. . . . .	1
2	Internet speed test platform, Sri Lanka. . . . .	5
3	Network diagram – key measurement points. . . . .	9
4	Data sources for official ICT statistics. . . . .	14
5	Number of <i>Bredbandskollen</i> observations, by year and type of service. . . . .	20
6	Data processing diagram. . . . .	22
7	Fixed-broadband download speeds (top) and upload speeds (bottom), Sweden, 2011-2016 – <i>Bredbandskollen</i> test results. . . . .	24
8	Histograms of fixed-broadband download speeds (top) and upload speeds (bottom), Sweden – <i>Bredbandskollen</i> test results. . . . .	25
9	Fixed-broadband download speeds, Sweden, 2012 and 2016 – comparison of <i>Bredbandskollen</i> results against other measurement platforms. . . . .	26
10	Share of <i>Bredbandskollen</i> observations per operator. . . . .	28
11	Share of <i>Bredbandskollen</i> observations per operator (left) and actual market shares (right), 2016. . . . .	29
12	Chi-squared test: share of <i>Bredbandskollen</i> observations per operator and actual market shares. . . . .	29
13	Fixed-broadband download speeds (left) and upload speeds (right), Sweden, 2011-2016 – <i>Bredbandskollen</i> adjusted test results. . . . .	30
14	Fixed-broadband download speeds (top) and upload speeds (bottom), by operator, Sweden 2011 and 2016 – <i>Bredbandskollen</i> test results. . . . .	31

15	Share of <i>Bredbandskollen</i> observations per region, 2011 and 2014. . . . .	32
16	Household density and share of <i>Bredbandskollen</i> observations per region, 2016.	33
17	Chi-squared test: share of <i>Bredbandskollen</i> observations and household density per region. . . . .	35
18	Fixed-broadband download speeds (left) and upload speeds (right), by region, Sweden, 2016 – <i>Bredbandskollen</i> test results. . . . .	36

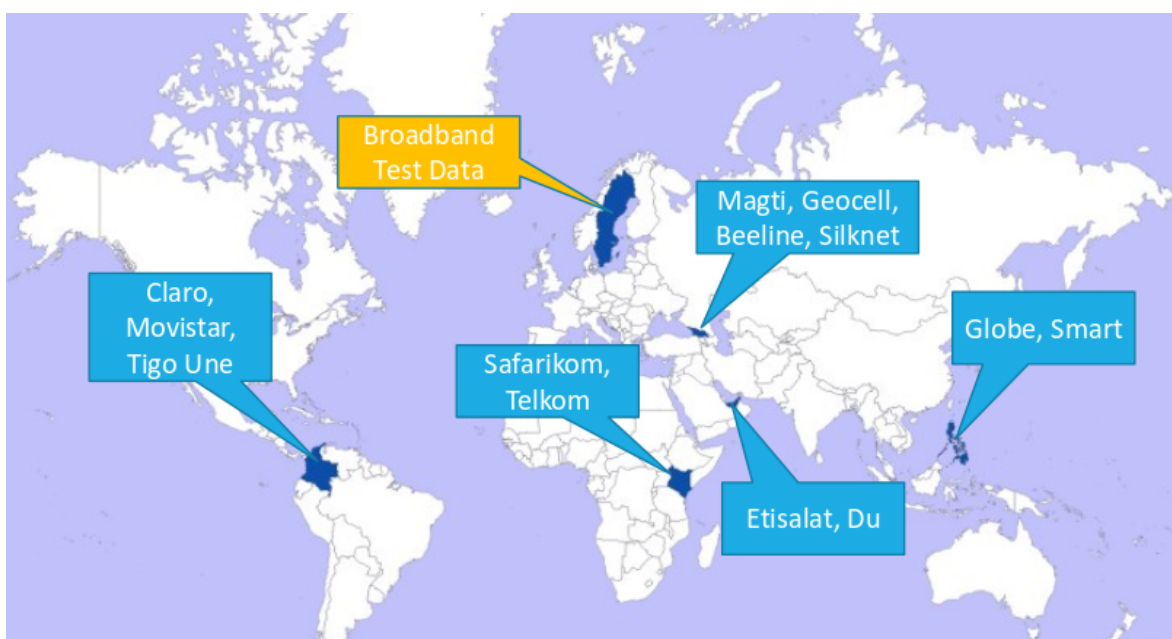
## List of Tables

1	Main differences between selected Internet measurement platforms. . . . .	12
2	<i>Bredbandskollen</i> data record for each observation. . . . .	21

# 1 Introduction

## 1.1 Background

The International Telecommunication Union – ITU, the United Nations specialized agency for information and communication technologies (ICTs)<sup>1</sup> – is carrying out a series of pilot studies under the umbrella of the project Big Data for Measuring the Information Society (Figure 1).<sup>2</sup>



**Figure 1:** Big Data for Measuring the Information Society – ITU pilot projects.

**Source:** Tiru (2016).

The objective of the project is to show how big data from the telecommunication industry can be used to produce new ICT indicators and replace or complement existing ones with a view to measuring the development of the Information Society worldwide.

---

<sup>1</sup>For more information on ITU, see ITU's website: <http://www.itu.int/en/about/Pages/default.aspx>.

<sup>2</sup>For more information on the ITU project Big Data for Measuring the Information Society, see its website: <http://www.itu.int/en/ITU-D/Statistics/Pages/bigdata/default.aspx>.

In particular, these pilot studies aim to be a first step towards filling in the ICT data gaps in the global indicator framework agreed for the monitoring of the 2030 Agenda for Sustainable Development (United Nations Economic and Social Council, 2016), and to inform private and public stakeholders on the current status of the digital divide.

This master thesis analyzes the measurements on broadband speeds carried out by the Internet Foundation in Sweden (IIS)<sup>3</sup> and made available to ITU in the context of the project Big Data for Measuring the Information Society.

## 1.2 Relevance

Broadband speed measurements matter because they are a key input to several consumer, policy and regulatory decisions:

- From a **consumer perspective**, broadband speed is one of the most important factors when choosing an Internet connection. For instance, in the European Union (EU) the download speed is the second most cited factor, after price, when deciding which broadband plan to choose.<sup>4</sup> However, six out of ten EU citizens do not know the maximum download speed of their broadband Internet plan and, among those that know it, a quarter of them believe that their real speed does not correspond to the one specified in their contract. Moreover, four in ten households in the EU admit having experienced difficulties accessing content at home because of speed or capacity issues (TNS Opinion & Social, 2014). It can thus be concluded that Internet speed is both an important and a controversial factor for consumers.

---

<sup>3</sup>The Internet Foundation in Sweden is an independent public-service organization which is responsible for the operation of the top-level domains '.se' and '.nu'. IIS reinvests part of the revenues obtained from the administration of these domains in , *inter alia*, the promotion of research on the Internet. As part of these activities, IIS has developed *Bredbandskollen*, a software-based Internet platform to measure actual broadband speeds. For more information, see IIS website <https://www.iis.se/english>.

<sup>4</sup>On average, 41 per cent of respondents in the EU mentioned download speed as an important factor when subscribing to an Internet connection, compared with 71 per cent of respondents citing price. The figures refer to fieldwork carried out in January 2014.

- From a **regulatory perspective**, broadband speed is one of the parameters often monitored to ensure that telecommunication operators and Internet service providers (ISPs) comply with some minimum quality-of-service (QoS) requirements. For example, Spain regulates the QoS parameters of electronic communication services by means of a service order, which includes specific Internet access parameters (Annex I, Part II in Boletín Oficial del Estado, number 156, Friday 27 June, 2014). In a similar fashion, the Telecom Regulatory Authority of India sets that subscribers should get a minimum of 80 per cent of the speed specified in their contract, as measured from the ISP node to the user (Telecom Regulatory Authority of India, 2006).
- From a **policy perspective**, broadband speeds have wide implications concerning the initiatives undertaken in the telecommunication sector. For instance, the definition of broadband is often tied to a given minimum speed, which is subject to be revised, as was the case in India in 2014 (from 256 to 512 kbit/s) (Telecom Regulatory Authority of India, 2014). In the United States, the Congress asked the Federal Communications Commission (FCC) to evaluate the deployment of *advanced telecommunication capabilities*. FCC considered these capabilities to require 4 Mbit/s download and 1 Mbit/s upload speeds in 2010, but revised the benchmark speeds to 25 Mbit/s download and 3 Mbit/s upload in 2015 (Federal Communications Commission, 2015a). In Europe, Finland was a worldwide pioneer in declaring affordable broadband access a basic right in 2010. Finland's Ministry of Transport and Communications set the threshold for the basic connection to 1 Mbit/s in 2010, revising it to 2 Mbit/s in 2016 (Davies, 2016). All these policy decisions have deep economic implications. Indeed, in most cases they imply the mobilization of universal service funds (USF) or subsidy schemes to meet the targets set in terms of availability of affordable Internet connections at a given speed.

In addition to these consumer, policy and regulatory implications, broadband speeds may also be an important determinant of broadband impact on economic growth (Rohman and Bohlin, 2012). Higher broadband speeds have also been found to be causally linked to increases in the percentage of employees classified as creative class workers (Whitacre et al., 2014).

All these factors motivate the interest in producing accurate data on actual broadband speeds.



### 1.3 Research questions

1. Can crowdsourcing Internet data be used to measure real broadband speeds?
2. Which information can be extracted from online speed measurements to characterize Internet users?

Given the relevance of broadband speeds for consumers, regulators and policy-makers, there is a growing demand for accurate measurements.

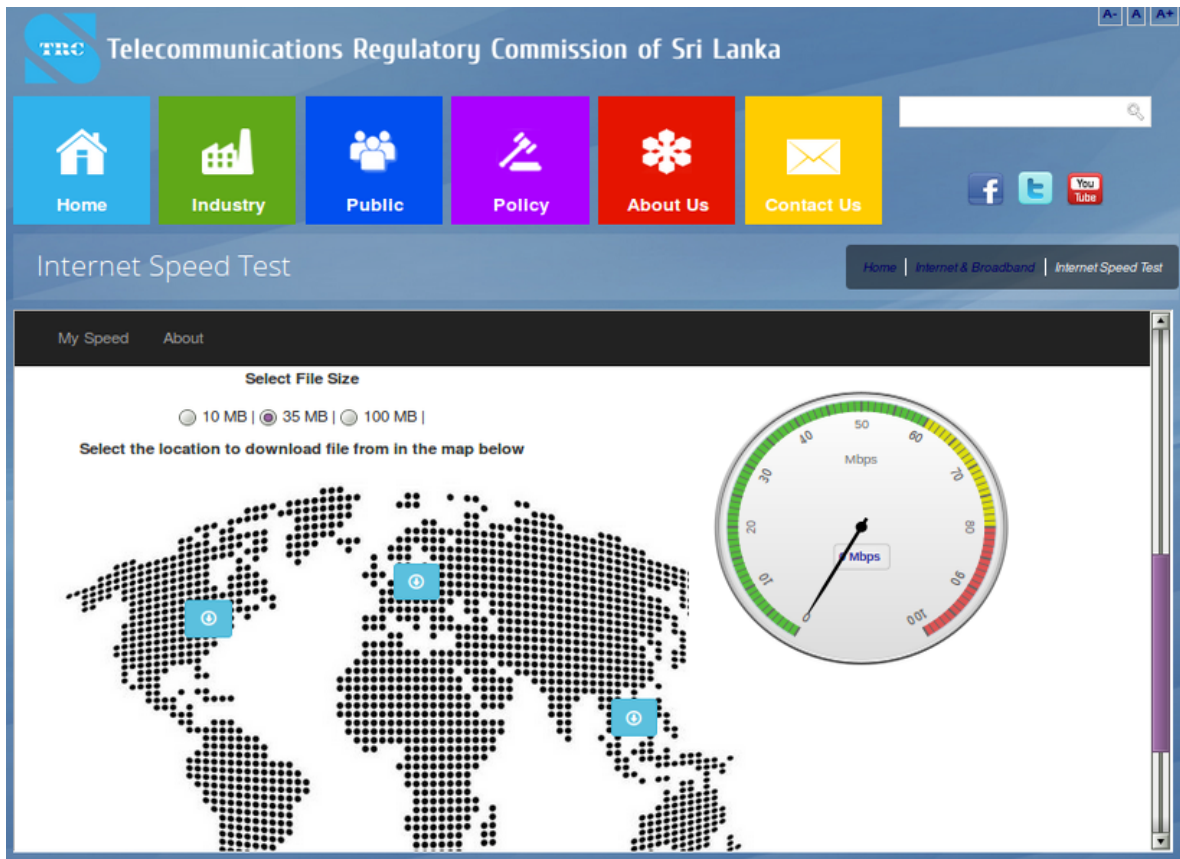
Advertised speeds, as publicized by operators and ISPs, provide only an upper-limit to the actual broadband speeds. On the other hand, precise external hardware-based measurements, such as the ones commissioned to SamKnows by the regulatory agencies in the UK (Ofcom, 2017), the United States (Federal Communications Commission, 2015b) and the European Commission (SamKnows, 2014) are costly. Therefore, they cannot be realistically scaled up to a wider set of countries.

Software-based, crowdsourcing data on Internet speed measurements remains the only possible stable source of real broadband speed information for most countries. Moreover, the low cost of deployment of these measurement platforms makes it possible to envisage its adoption by any interested regulator/policy-making.

Indeed, some regulators in developing countries, such as the Telecommunications Regulatory Commission of Sri Lanka, have already launched their own measurement portal (Figure 2). Moreover, there are private stakeholders, such as Ookla, recording these data at the global level.<sup>5</sup>

---

<sup>5</sup>For more information on Ookla's Speedtest, see Section 2.1 and Speedtest's website: <http://beta.speedtest.net>.



**Figure 2:** Internet speed test platform, Sri Lanka.

**Source:** [Telecommunications Regulatory Commission of Sri Lanka](#).

In this context, the Internet Foundation in Sweden has been a forerunner of public-service broadband measurement platforms with its portal *Bredbandskollen*, which collects data on Internet speeds since 2008 (Davidsson, 2017).

This project takes advantage of the microdata made available by IIS to ITU on the broadband speed measurements from the *Bredbandskollen* platform in the last six years (2011-2016). This large dataset is used as a testbed for determining to which extent this kind of crowdsourcing measurements can be used to provide robust insights into real broadband speeds, as well as information on Internet user behavior.

## **2 State of the art**

Broadband speeds are used by operators and ISPs as a means of characterizing their offers and segmenting their services according to different QoS. Advertised speeds are determined by each service provider according to their own internal methodology, which is not disclosed nor usually inspected by an independent party. However, it is understood that these speeds indicate the maximum or peak data rates that a customer may experience in the link between the customer location and the broadband provider (Bauer et al., 2010).

More generally, advertised speeds are taken by customers to imply something meaningful about their experience when using a given broadband service. Since customers do not know in advance which Internet sites they will want to access, they demand universal connectivity from the ISPs and therefore expect advertised speeds to correspond to their end-to-end experience (see the discussion on the MCI merger in Economides, 2008, pp. 508-511).

For example, a customer contracting a 10 Mbit/s broadband plan with Movistar in Spain will expect to access information from The New York Times website at speeds close to 10 Mbit/s, even though the content of that website may be hosted outside Movistar's network.

This is but one of the subtleties which may have a significant impact on the results from a broadband speed measurement test. Section 2.1 reviews this and other technical factors having an impact on the state-of-the-art methodologies used to measure broadband speeds. Section 2.2 presents how the problem of making statistical inference from crowdsourcing data on broadband measurements has been tackled in similar research efforts.

### **2.1 Measurement methodologies**

From a technical standpoint, the most accurate option for measuring real broadband speeds is to implement in-network measurements, such as real traffic monitoring (using network counters) or test-call routines. The technical issues concerning this type of measurements are well known. Indeed, international standardization bodies, such as the European Telecommunications

Standards Institute, have issued internationally agreed guidelines on how to perform this type of measurements (European Telecommunications Standards Institute, 2005).

In some countries, such as India and Spain, operators are required to self-report in-network measurements and the results are disclosed to the public on a regular basis (Ministerio de Energía, Turismo y Agenda Digital, 2017; Zuhyle and Mirandilla-Santos, 2015).

However, in-network measurements can only be accomplished by those stakeholders having direct access to the network infrastructure and therefore rely on the self-reporting of the concerned network operators. Because of the difficulty of overseeing compliance and the lack of independence of this type of measurements, this approach is usually not considered a solution on its own, but rather a complement to other external measurements.

Therefore, external broadband speed measurements are the most common approach to gaining insights about real broadband speeds. There are a number of private stakeholders performing this kind of measurements and publishing data with a global coverage. These include, *inter alia*, Ookla's Speedtest and Akamai speed reports (Bauer et al., 2010, 2016; Lehr et al., 2013).

In addition, there exist several external broadband speed measurement platforms commissioned or operated by public entities and the academia. Some of them rely on hardware-based approaches, such as the SamKnows' Whitebox (SamKnows, 2013, 2014) or the BISmark router (Chetty et al., 2013; Sundaresan et al., 2014, 2012). However, the majority of external measurements rely on pure software approaches.

Examples of software-based approaches from public entities include the Regulatory Commission of Sri Lanka's Internet speed test platform (Zuhyle and Mirandilla-Santos, 2015), the Internet Foundation in Sweden's *Bredbandskollen* and the Italian Authority for Communications Guarantees' *Misurainternet* platform.<sup>6</sup>

These platforms collect additional broadband performance metrics apart from download speed, such as upload speeds and latency (delay). Depending on the type of online activity performed, some parameters will be more important than others. For instance, latency is very relevant for

---

<sup>6</sup>For more information, see the Misurainternet website: <https://www.misurainternet.it>.

real-time communications, such as VoIP. As a result of the diversity of activities performed online, a complete assessment of user quality of experience cannot be summarized into a single speed metric, but will require several complementary measurements (SamKnows, 2014; Wattegama and Kapugama, 2011; Zuhyle and Mirandilla-Santos, 2015).

Moreover, there exist some application-specific broadband measurements, such as Netflix's Fast.com and Youtube's Speed numbers. These measurement platforms are optimized to reflect what is the actual speed experienced by a user transferring video files with the size and protocols common in Netflix and YouTube, respectively. However, these speed measurements may be a poor approximation of the actual speed experienced when browsing a website or sending a file (Bauer et al., 2010, 2016).

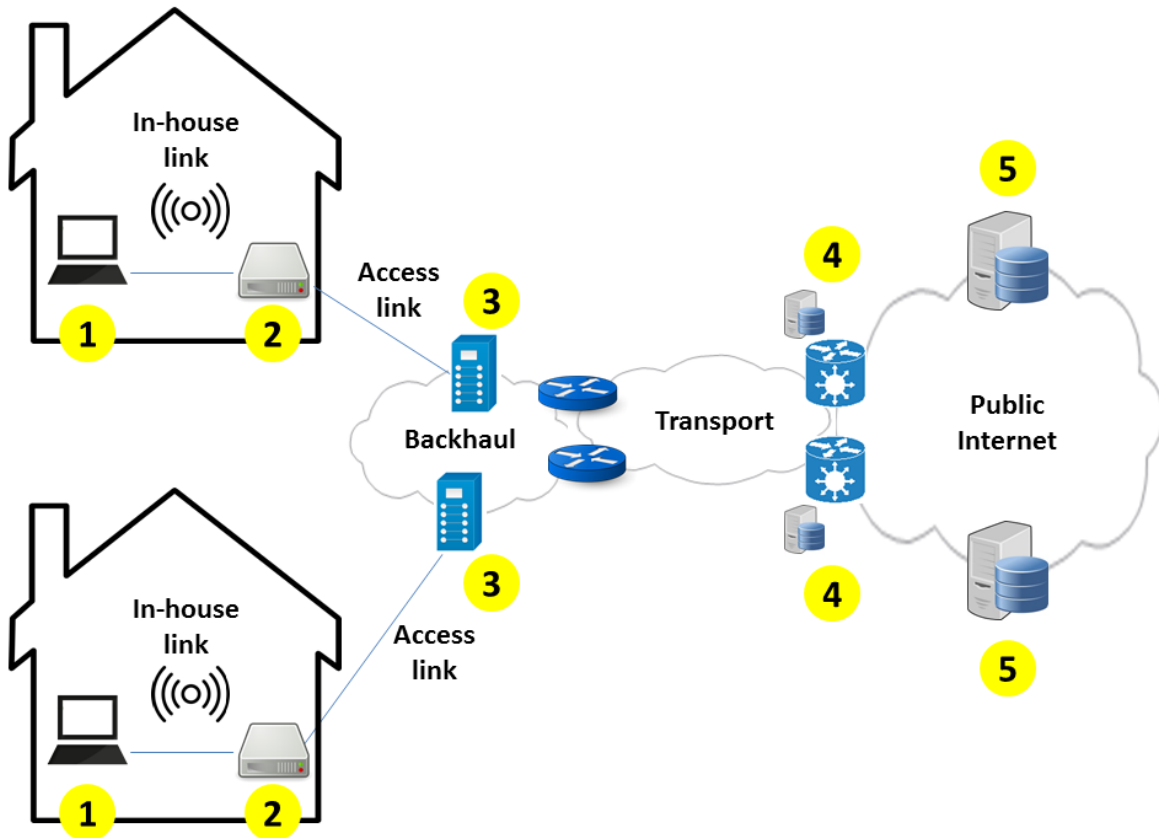
Based on a literature review of the state-of-the-art in broadband speed measurements, it can be concluded that external broadband speed measurement platforms depend on a few main design characteristics that affect their outcomes:

**1 – Measurement path:** from which point to which point of the network the measurement is made (Figure 3). Speeds advertised by operators refer to the maximum achievable over the access link (points 2-3 in Figure 3). In xDSL connections, the access link is not shared, as it is for coaxial cable, but xDSL speeds are more sensitive to the distance to the local exchange (path from 2 to 3).

From point 3 in Figure 3, there is contention for all wired technologies (fibre, coaxial, copper). That is, the transmission capacity is shared by many concurrent users and, depending on the network charge, this may lead to congestion and lower speeds. Even so, from points 2 to 4 the connection remains within the network of the ISP with whom the end user has contracted the service and therefore the speeds depend only on the network dimensioning and management of that ISP.

End users' quality of experience is affected by the end-to-end path (i.e. points 1 to 5). This means that the speed experienced by a user also depends on the quality of the in-house connection (path from 1 to 2). For instance, the WiFi connection may not support high speeds or there might be several users making use of the same connection,

therefore reducing the bandwidth available for the single user performing the test. Home network bottlenecks have been found to be very relevant in those settings in which the access is capable of providing speeds greater than 20 Mbit/s (Sundaresan et al., 2016).



**Figure 3:** Network diagram – key measurement points.

**Source:** Author based on Bauer et al. (2010).

In addition, end-to-end speeds may depend on other ISPs networks if the end user is accessing content hosted on the wider public Internet (off-net) instead of within the home ISP network (on-net). In the off-net case, the connection includes points 4 to 5 which are not under direct control of the home ISP.

However, it is the decision of the home ISP with whom to interconnect and under which conditions (e.g. traffic, peering). Therefore, under the assumption that customers demand universal connectivity from the ISP with whom they have the contract (see

Economides, 2008, pp. 508-511), this ISP could also be hold accountable for the speed delivered over the link from 4 to 5 in Figure 3. Indeed, ISP interconnection may have a substantial impact on consumer Internet performance and business relationships between ISPs are often at the root of this performance degradation, rather than technical issues (M-Lab Research Team and others, 2014).

- 2 – Active versus passive testing:** passive testing measures broadband speeds over end users' normal online activities, whereas active testing relies on standardized tests run independently of each end users' actual online activity. Active testing is usually preferred in benchmarking exercises because it facilitates the comparability between measurements (Zuhyle and Mirandilla-Santos, 2015). In passive testing, if two users are performing very different activities (e.g. heavy video streaming compared with light web browsing) the outcome of the speed measurements may be significantly different even if they have the same type of connection.
- 3 – Voluntary versus automatic testing:** some broadband speed measurements are initiated by the end user (most software-based platforms, including Ookla and Bredband-skollen), whereas a few others are scheduled remotely and take place automatically (e.g. hardware-based platforms, such as SamKnows). Voluntary testing has several drawbacks, including the risk of selection bias. That is, end user's may run the test in a diagnostic fashion when they are experiencing network problems (Bauer et al., 2010). Moreover, actual speeds are sensitive to congestion: they tend to decrease at peak hours or during traffic burst times (Sundaresan et al., 2012; Zuhyle and Mirandilla-Santos, 2015). As a result, broadband speed results of voluntary tests may be biased because of the timing of the tests, which is out of control of the measurement platform and may not be randomly distributed. For instance, there could be more tests run during congestion periods because that is precisely when end users' perceive speed problems.
- 4 – Data protocol configuration:** the data transmission protocol, usually TCP, may actually be the bottleneck if not correctly configured. This is more so in high-capacity networks (e.g. gigabit broadband networks), where a single TCP flow will most likely not be enough to measure the real capacity of the link. Therefore, multiple parallel TCP flows

will be required to produce a reliable measurement in high-capacity networks (Bauer et al., 2010, 2016).

**5 – Statistical aggregation:** results for individual tests are calculated using different methods. These include: (i) total bytes / total time, (ii) total bytes / total time after the ramp-up period, to exclude the initial warm-up period in which transmission is slower, and (iii) payload / bytes, so that only the actual information (and not the header aggregated by the transfer protocol) is considered. In addition, different aggregation methods are used for reporting aggregate statistics. For instance, SamKnows discards the top and bottom percentiles to control for outliers and averages the rest (SamKnows, 2014). Ookla removes the fastest 10 per cent and slowest 30 per cent slices of each measurement and averages the rest (Bauer et al., 2010). *Bredbandskollen* takes the higher between the average speed of the whole measurement period (10s) or the average of the last 8s.<sup>7</sup>

**6 – Other design parameters,** such as test duration time, technology used in the application (e.g. HTML or Flash-based) and the file size used in download/upload measurements (Bauer et al., 2010; Zuhyle and Mirandilla-Santos, 2015).

Table 1 summarizes how this design parameters vary according to the broadband speed measurement platforms chosen in this project to benchmark the results obtained with *Bredbandskollen*.

The main differences would be the following: SamKnows ability to exclude the home network from the measurement and to automatically control the test timing, which will lead to more robust results; Akamai’s passive testing, which may lead to comparing apples and oranges if their results are used to benchmark the outcome of active tests; the different statistical aggregation procedures which may potentially have a significant impact on the results.

Another relevant factor which needs to be considered is the different **sample sizes** of each broadband platform. For instance, concerning data from Sweden (i.e. the object of this project), *Bredbandskollen* collected some 15 million observations per year for download

---

<sup>7</sup>See *Bredbandskollen* methodological [FAQ](#).



speeds, Akamai 20 million and SamKnows 0.52 million. These differences illustrate the trade-off between volumes of data and measurement precision. Indeed, hardware-based measurements are costly and this severely limits the sample size even in developed countries.

	Measurement path	Active passive	Voluntary automatic	Data flows	Statistical aggregation
<b>Akamai</b>	1-4 or 1-5 depending on server load	Passive	Real traffic	Sequential	Unknown
<b>Bredbandskollen</b>	1-4 in most cases. 1-5 for the rest	Active	Voluntary	Parallel	Avg.first 2s or avg 10s
<b>Ookla</b>	1-4 or 1-5 depending on user location	Active	Voluntary	Parallel	Avg. after ramp-up. Excl. top 10% and bottom 30% slices
<b>SamKnows</b>	2-5	Active	Automatic	Parallel	Avg. after ramp-up. Excl. top and bottom percentiles. Separate peak / off-peak

**Table 1:** Main differences between selected Internet measurement platforms.

**Source:** Author based on Bauer et al. (2010, 2016); Canadi et al. (2012); SamKnows (2014) and *Bredbandskollen*.

## 2.2 Statistical approach

Section 2.1 has dealt with what traditional statisticians would consider measurement error. This section looks into the issues related to the statistical treatment of the data. In particular, it looks at how similar studies have dealt with the problem of extracting representative insights about the whole population (or, at least, the in-scope population) based on non-representative Internet data.

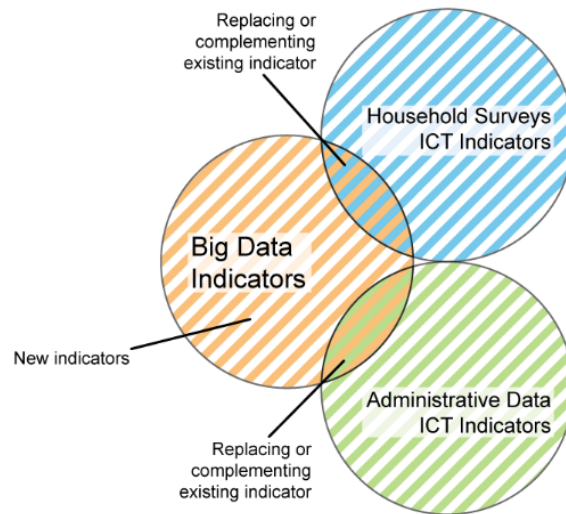
That is, given the data provided by a given broadband measurement platform (*Bredbandskollen* in this project), how can we make statistical inference of the whole population or, at least, about some segments of it?

Official indicators on telecommunications/ICTs are divided into two categories according to the source of the data collection: administrative data, and household surveys or censuses. The methodologies for collecting harmonized and internationally comparable ICT data from these two data sources are well known. See for instance the ITU Handbook (International Telecommunication Union, 2011) and the ITU Manual (International Telecommunication Union, 2014), respectively for administrative and survey ICT data sources.

Big data, such as the millions of records obtained from broadband measurement platforms, are often a middle ground between these two types of data sources (Figure 4).

Despite the ever larger digital datasets which are becoming available, if a given data source does not include all the data for a given type of measurement (e.g. the real broadband speed of all connections according to a given measurement approach), then it cannot be treated statistically as administrative data records and requires some of the techniques usually applied to survey data, such as those concerning sampling (Boyd and Crawford, 2012).

Official ICT statistics from household surveys are grounded on probability-based sample surveys and a frequentist statistical framework (Couper, 2013). However, there is much more statistical uncertainty about the measurement in the data collected from broadband measurement platforms than in the programmed surveys undertaken by national statistical offices.



**Figure 4:** Data sources for official ICT statistics.

**Source:** Tiru (2016).

Indeed, broadband speed measurements are composed of non-representative Internet data and, as such, face similar problems to those found in non-probabilistic samples (Couper, 2013; Zagheni and Weber, 2015):

1. **Quota-based samples:** a panel is selected trying to adjust its composition according to the social, demographic and/or geographic characteristics of the in-scope population. More in general, quota-sampling is based on the selection of the surveyed individuals based on a series of covariates that are believed to explain the target variable. For instance, in its quota-based sampling methodology, SamKnows considers country, ISP and connection technology as the covariates used to define the quotas (SamKnows, 2014). Quota sampling in itself does not solve the issue of selection bias, particularly when participation is voluntary. This is indeed the case in most broadband speed measurement platforms, including SamKnows.
2. **Online voluntary surveys:** on top of the selection bias, unequal participation of the volunteers may add more uncertainty about the representativeness of the results. Indeed, there is evidence from the survey research domain pointing that a relative large numbers

of online surveys are completed by a relative small number of active panelists (Couper, 2013). In the domain of broadband speed measurements, in some platforms where it is possible to identify individual end users, such as in *Bredbandskollen*, there is evidence that some active users may perform thousands of tests per year, whereas some others only a few. In some other broadband speed measurement platforms, however, one-time test users were found to be predominant (Wattegama and Kapugama, 2011).

From all statistical issues, selection bias is the central one when trying to make statistical inferences from the data produced by a broadband speed measurement platform. For instance, Rood et al. (2012) used official statistics on geography, technology, gender and age<sup>8</sup> to test whether self-selected measurements of a broadband speed measurement platform (iPing) were representative of the whole population. The results rejected this hypothesis, and the authors concluded that higher speed connections, males and persons aged 30-60 were overrepresented. Moreover, the paper found that the results of Ookla and Samknows in the Netherlands also risked not being representative of the whole population.

Several approaches have been proposed to deal with selection bias:

1. **Selection panel:** create a proper selection panel from which to draw detailed socio-economic household and individual data in addition to the broadband measurements. This selection panel can then be used to weight the results of the self-selected panel and produce representative results out of it (Rood et al., 2012). This is indeed a reliable but costly option, because it requires the maintenance of a regular panel in addition to the upkeep of the speed measurement application.
2. **Calibration against ground-truth data:** use data from a reliable source (i.e. the so called ground-truth data) to weight the results of the broadband speed test in order to ensure that the results of the two sources correspond for a period in which they overlap. This correction relies on the assumption that there is a (potentially stochastic) structural

---

<sup>8</sup>Official statistics on geography, technology, gender and age in the Netherlands are produced via household surveys by Statistics Netherlands (CBS) and administrative regulatory filings of ISPs to the Dutch regulator OPTA.

relationship between the subgroup of the population included in the broadband speed test sample and those not included. Moreover, this relationship needs to hold in time (Zagheni and Weber, 2015). Couper (2013) warns that there is a risk of in-sample overfitting and type I errors when calibrating large datasets.

3. **Differences-in-differences approach:** when it is not possible to obtain a robust point in time estimation because the previous approaches do not apply, it may still be possible to obtain a sound indication of trends by applying a differences-in-differences approach (Zagheni and Weber, 2015). For such an approach to be feasible, a control group needs to be identified and the selection bias in the control and the measurement group needs to be constant with time or evolve in parallel.<sup>9</sup> In addition, analogous to calibration, a key assumption needed to validate this approach is that relative changes for the subgroups of the population included in the sample should be indicative of trends in the general population.

Based on these points, a key question that needs to be addressed when analyzing broadband speed measurements is how volunteers taking part in the test differ from non-volunteers. This is a common question to most non-probabilistic survey research (Couper, 2013).

In practical terms, the statistical issues highlighted in this section often cannot be addressed *a posteriori*; that is, during the data processing stage where the available information is fixed.

For instance, a common problem when analyzing software-based broadband speed tests is that the information on the number of unique users is not available or not disclosed to the researchers (Canadi et al., 2012). This is also the case in the *Bredbandskollen* microdata that IIS shared with ITU, which for privacy issues does not contain a unique identifier for each user, but rather for each test. The lack of this information makes it difficult to perform any ex-post calibration, because the cross-linking of the observations from the speed measurement platform with external individual population characteristics is not possible.

---

<sup>9</sup>For example, there may be a selection bias in the control and measurement groups towards an over-representation of end users from a given operator. If at some point in the time series the selection bias changes only in the control group (e.g. there is no longer an over-representation of clients from that operator), then the differences-in-differences approach would no longer be valid.

When facing major problems inherent to the structure of the underlying dataset, researchers often follow the approach of describing the limitations of the dataset and publishing the results as such, with little (Chetty et al., 2013; Prasad et al., 2016; Wattegama and Kapugama, 2011) or some admonition about their representativeness for the whole population (Canadi et al., 2012; Riddlesden and Singleton, 2014). In the latter, an approximate correspondence between the global average speed obtained in the study and that obtained in some previous more robust measurement effort is considered as a broad indication of the validity of the results.

More generally, the validity of the statistical inference derived from each broadband speed measurement platform may depend on the purpose of the research. If the objective is to raise awareness that some consumers are not getting the speeds specified in their contracts, showing that this is true for a large enough number of end users' may suffice (Chetty et al., 2013; Wattegama and Kapugama, 2011).

If the purpose of a broadband speed measurement platform is customer self-check, a sound measurement methodology (as discussed in Section 2.1) and consistency in the results for a single user is all that is required. This is, for instance, the primary objective of Ookla's Speedtest. The Italian Authority for Communications Guarantees' *Misurainternet* platform incorporates this logic in the application itself. Indeed, by means of an electronic certificate, legally validated results of the speed test can be communicated to the concerned ISP who then has the obligation to reestablish a minimum speed according to the conditions specified in the contract.

If the objective of the research is to obtain representative information on the real broadband speeds experienced by broadband subscribers in a country, as it is the case in official telecommunication statistics, then all the statistical considerations discussed in this section should be addressed.

Lastly, broadband speed tests can also be a rich source of data for analyzing different Internet conditions and user behaviors across socio-economic strata. An obstacle to mapping broadband speed data to external data sources is the fact that broadband tests usually do not collect demographic information. Indeed, socio-demographic information is key to linking

technical measurements to wider societal issues (for an example, see Chapter 6 in International Telecommunication Union, 2016).

Some broadband speed measurement platforms have bypassed this problem by asking end users to provide extra voluntary information, such as their postal code. Based on this extra information, it is possible to geolocate each test and link it with indicators of socio-spatial structures, such as rurality and levels of material deprivation (Riddlesden and Singleton, 2014). Other platforms ask for demographic data, such as age and gender, on a voluntary basis after or prior to the test. This information allows the linkage of the test results with demographic information from household surveys (Rood et al., 2012).

### 3 Data

This project analyzes the microdata collected by means of the *Bredbandskollen* platform, which was developed and is operated by the Internet Foundation in Sweden (IIS).

*Bredbandskollen* is a free online software tool. Internet users can measure their real broadband speeds for both fixed and mobile connections by running the test either on [Bredbandskollen website](#) or by downloading an app and executing the test from a mobile device. Broadband speed tests are initiated by the end user, whenever they click on the start test button on [Bredbandskollen website](#) or launch the measurement from the mobile app.

*Bredbandskollen* measures the download and upload speed by sending data to and receiving data from the closest Internet exchange point (IXP) to the end user.<sup>10</sup> Currently, *Bredbandskollen* has test servers hosted in five IXPs (in Malmö, Göteborg, Stockholm, Sundsvall and Oslo), against which the tests are run. As a result, measurements often take place in an IXP with which the major ISPs in Sweden will have a direct interconnection.

---

<sup>10</sup> Each download/upload speed measurement is implemented by sending simultaneous HTTP requests from the client to the server via the socket function in Flash. If it fails in Flash, the test is repeated via the browser. The number of concurrent HTTP requests depends on the speed of the connection, the higher the speed, the higher number of HTTP threads that are initiated in order to reach the connection's maximum.

*Bredbandskollen's* download and upload speed tests last ten seconds each. During the test, average speeds are calculated for the first two seconds and for the entire ten seconds. The highest of the two results is retained.

*Bredbandskollen* also measures latency, i.e. the delay or round-trip time to the closest server. Latency is not the object of study of this project and it may require a slightly different approach to that applied to download and upload speeds, given that latency results are very sensitive to the distance from the end user to the test server.

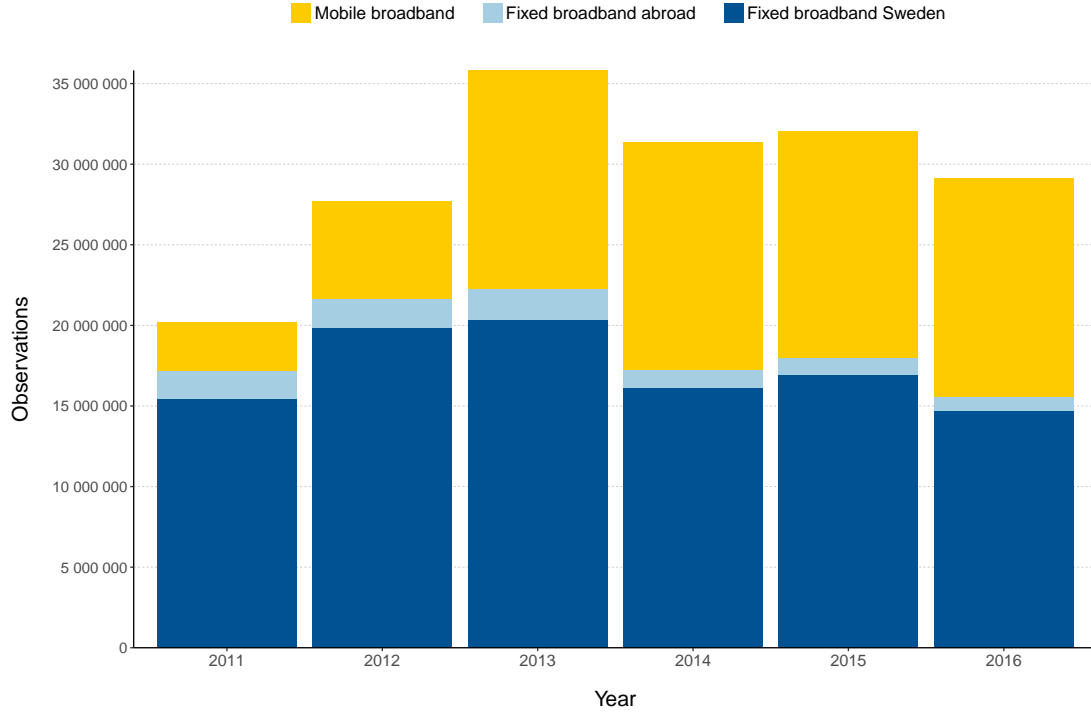
*Bredbandskollen* has been optimized to measure broadband speeds in Sweden and in the neighboring border areas. That is indeed the geographical scope of this project.

The dataset used for the analysis in this project includes around 15 million test observations per year for fixed-broadband connections (Figure 5). Mobile-broadband observations are not within the scope of this project, because they are affected by different technical determinants and would require a different, more challenging processing.

Observations from fixed-broadband tests performed outside the country are also excluded, given that they concern the quality of the access network of another country and therefore require a separate analysis.

Lastly, the fixed-broadband speed tests ran in Sweden are filtered to exclude those in which download speed, upload speed or latency were not greater than zero. These are considered as erroneous measurements and included in the light blue bar in Figure 5.





**Figure 5:** Number of *Bredbandskollen* observations, by year and type of service.

Table 2 shows the information recorded for each observation.

It is important to note that the unique identifier associated with each observation corresponds to a measurement, not to an end user. That is, a user may have tested several times their connection and each test will appear with a different unique identifier. Although *Bredbandskollen* identifies unique users by means of the cookies used on the website, this information has not been shared so far with ITU because of privacy concerns.

For the data analysis carried out in this project, the following fields were used: client, date/time, download speed, upload speed, average response time, country code, region code, network type and operator.

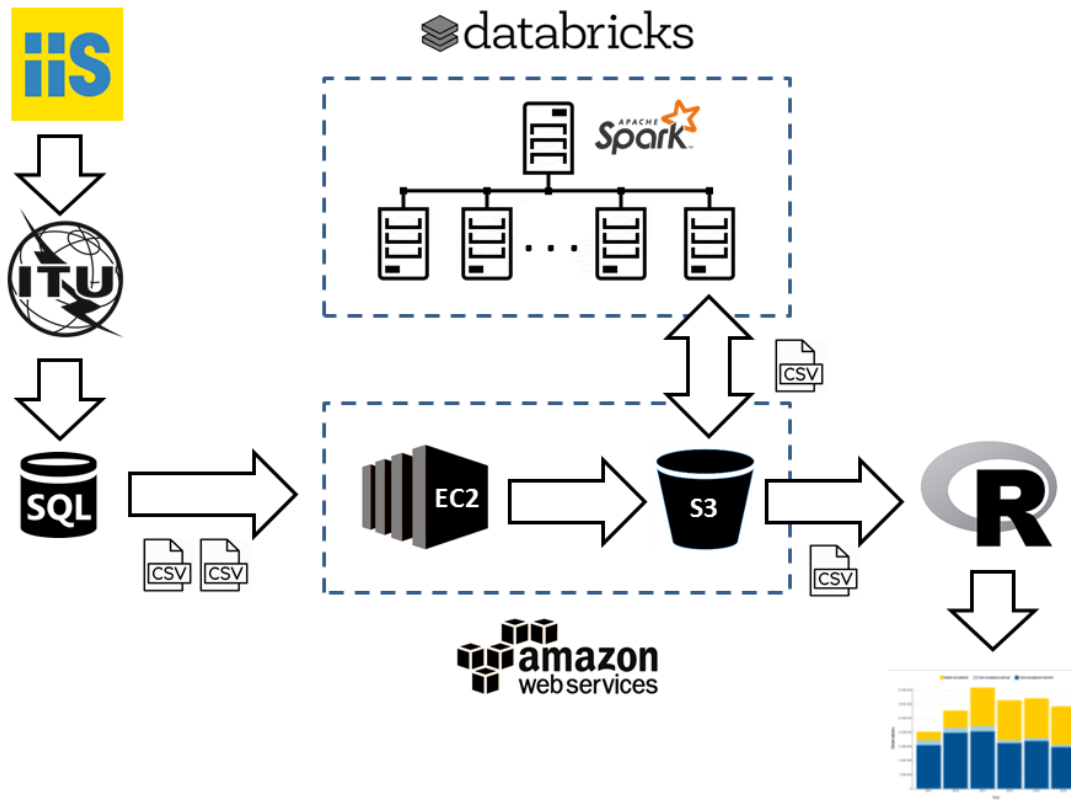
The fields average response time, network type and country code were only used to filter those observations that were valid and had been performed from a fixed-broadband connection in Sweden. The other fields were used to produce the aggregate statistics presented in Section 4.

Field	Description
Unique ID	Unique number identifying each test. A single user running several tests will be assigned several unique IDs, i.e. one for each test.
Client	iOS (mobile app for iOS), Android (mobile app for Android) or Web (measurement from the website).
Date/time	UTC time, i.e. 1 or 2 hours less than <a href="#">Swedish time</a> depending on the period of the year.
Download speed	Result of the test in Mbit/s.
Upload speed	Result of the test in Mbit/s.
Avg. response time	Result of the latency test in milliseconds.
Country code	Country code of the location where the test was conducted. Empty if the location of the measurement is not available.
Region code	<a href="#">Region code</a> of the location of the measurement. Blank if outside Sweden.
Municipality code	<a href="#">Municipality code</a> of the location of the measurement. Blank if outside Sweden.
Municipality	Name of the municipality. May be empty.
Network type	May be empty. If the client type is iOS or Android, the following options are possible: CDMA, GPRS, EDGE, UMTS, HSPA, HSDPA, HSPAP, LTE, Mobile (if the network type is not detected) and WiFi (measurement via WiFi or otherwise not using the mobile network). If the client type is Web, there is some description of the subscription. Often it is based on information provided by the user, so it is not always reliable.
Operator	Operator name, inferred from the IP number. Usually correct, but not always reliable. May be empty.
Phone model	Description of the phone model, if the client type is iOS or Android. If the client type is Web, description of the web browsers.

**Table 2:** *Bredbandskollen* data record for each observation.

### 3.1 Processing environment

Figure 6 shows a graphical summary of the data processing environments used in this project and the corresponding data flows.



**Figure 6:** Data processing diagram.

Given the large size of the dataset (more than 100 million observations corresponding to about 30 GB of data), the data processing could not be performed locally nor with sequential algorithms.

Data were therefore processed in a Spark cluster<sup>11</sup> set in a Databricks environment.<sup>12</sup> The cluster was composed of one driver and up to eight worker nodes.<sup>13</sup> The worker nodes were auto-scaled according to the computing requirements of each job submitted.

The *Bredbandskollen*'s records were obtained in comma-separated values ("csv") format from the SQL database. These files were downloaded directly into an Amazon Web Service (AWS) instance and transferred to an AWS simple cloud storage service (S3) using multipart upload.<sup>14</sup> The S3 data bucket was mounted on the Databricks environment, thus linking the Spark cluster with the raw ".csv" files.

The heavy data processing was carried out in the Spark cluster and the summary results saved back into ".csv" files. These files were transferred to a local computer and used to produce the charts using R.

## 4 Results

Figure 7 shows the average for download and upload speeds of all valid test results performed from a fixed-broadband connection in Sweden. In addition to the average, the median is presented as well as the result of removing the top and bottom percentiles and averaging the rest. This latter approach is proposed by SamKnows (2014) in order screen anomalous or misrepresentative test results.

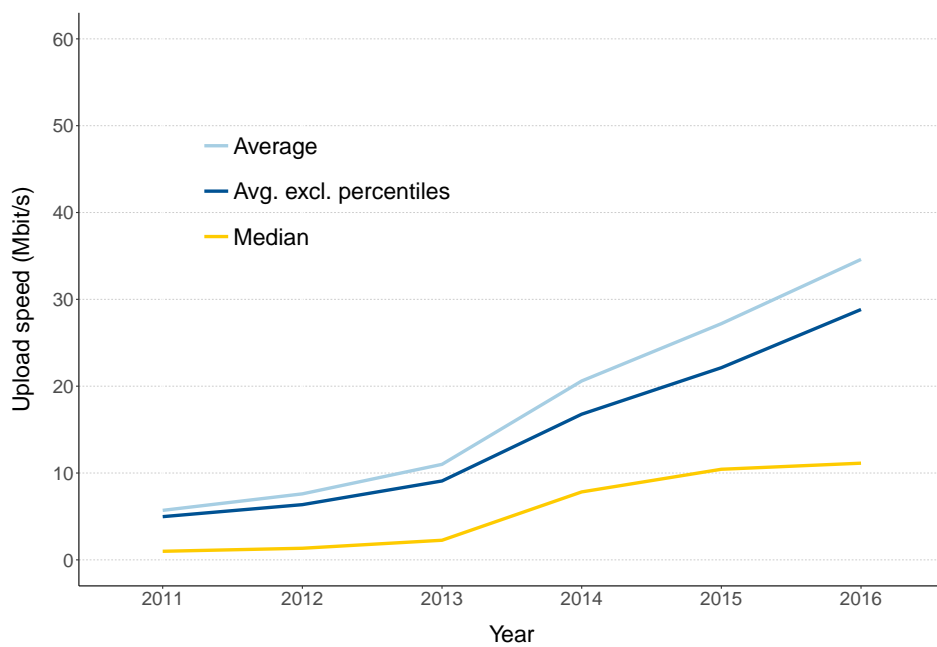
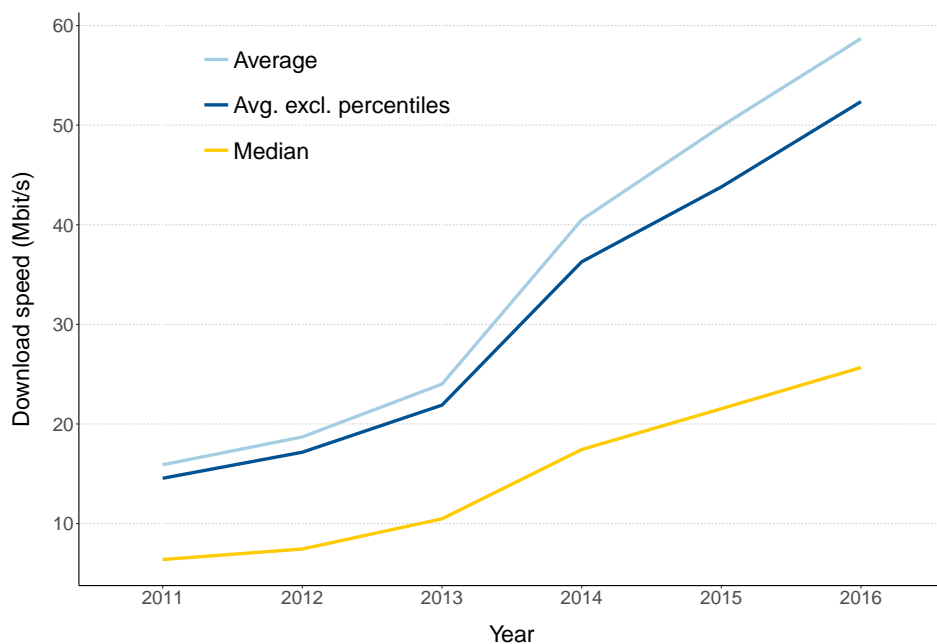
---

<sup>11</sup>Apache Spark is an open-source cluster computing framework that provides an interface for large-scale data processing with implicit parallelism. In comparison to other cluster-computing platforms, such as Hadoop Map Reduce, Spark is significantly faster because, *inter alia*, it exploits in memory computing. Moreover, Spark provides a well maintained Python API (PySpark) which makes it possible to program Spark clusters using Python-like syntax.

<sup>12</sup>Databricks is an online analytics platform that provides environments for cloud-based big data processing using Apache Spark. For more information, see Databricks website: <https://databricks.com>.

<sup>13</sup>Each node was an Amazon Web Service instance of the type r3.xlarge, with 30.5 GB of memory, four cores and one DBU.

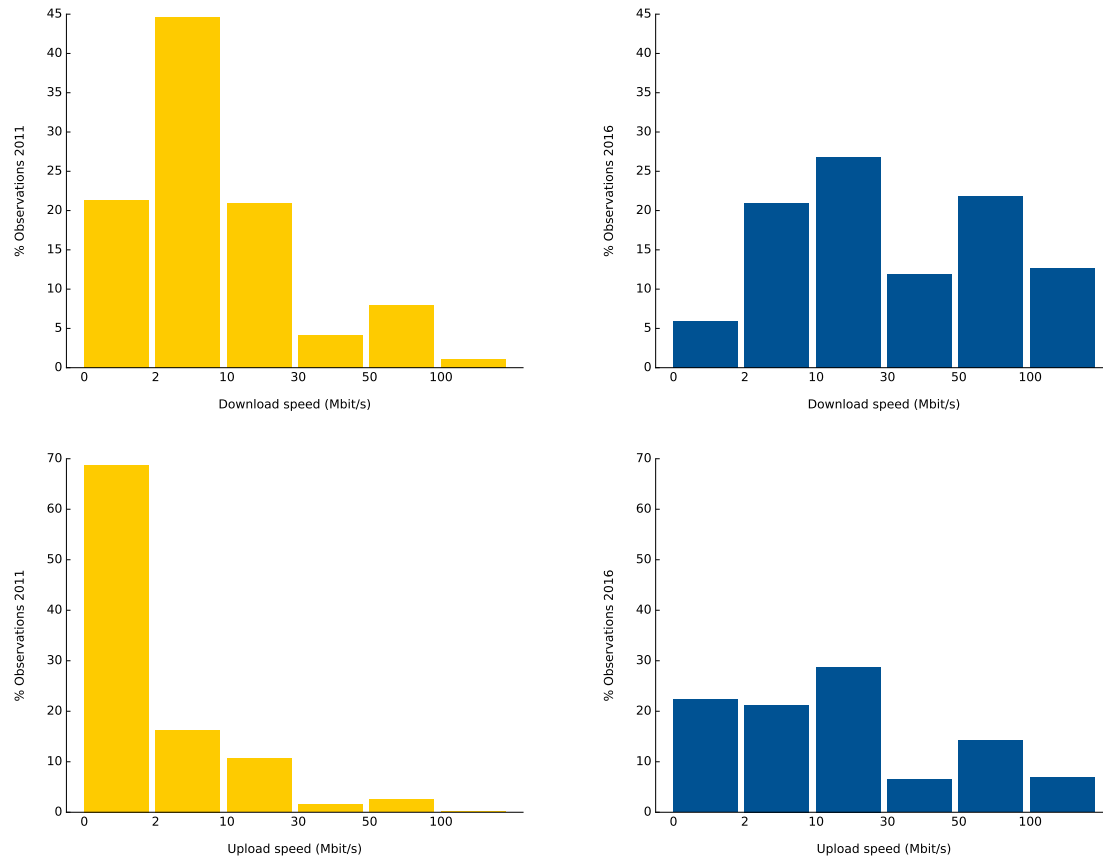
<sup>14</sup>AWS multipart upload enables to upload large files in parts. For more information, see the AWS webpage on this topic: <http://docs.aws.amazon.com/AmazonS3/latest/dev/mpuoverview.html>.



**Figure 7:** Fixed-broadband download speeds (top) and upload speeds (bottom), Sweden, 2011-2016 – *Bredbandskollen* test results.

**Note:** Avg. excl. percentiles refers to the average obtained after removing the top and bottom percentiles.

There is a significant difference between the median and the mean, which is further explored in Figure 8. Indeed, the non-negligible number of observation above 100 Mbit/s in 2016 (blue bars) suggests that the distribution has heavy tails that drag the average up.

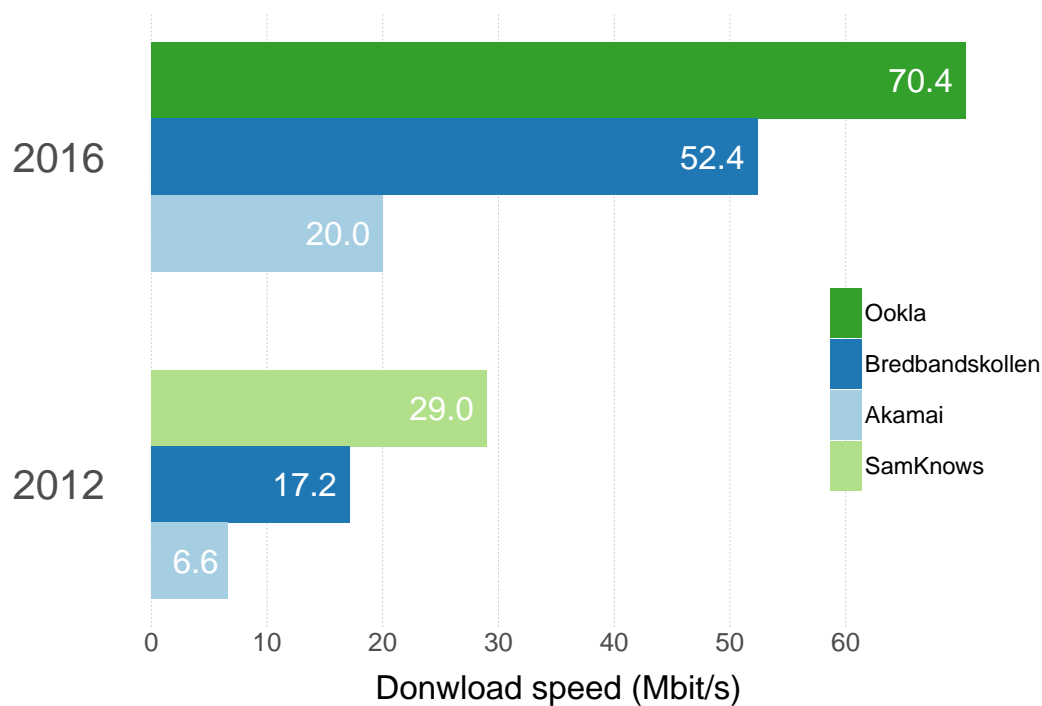


**Figure 8:** Histograms of fixed-broadband download speeds (top) and upload speeds (bottom), Sweden – *Bredbandskollen* test results.

The histograms also show the evolution in the distribution of observed speeds between 2011 and 2016. For instance, download speeds above 30 Mbit/s – the threshold considered by the European Commission for high speeds in the Digital agenda (European Commission, 2010) – represented around 10 per cent of all observations in 2011, whereas in 2016 they surpassed 40 per cent.

As could be expected, upload speeds tend to be lower than download speeds: only 15 per cent of all observations had upload speeds above 30 Mbit/s in 2016, and in 2011 most observations had upload speeds in the interval 0-2 Mbit/s.

Figure 9 compares the results of *Bredbandskollen* with those obtained by other broadband speed measurement platforms. For *Bredbandskollen*, the average excluding the top and bottom percentile is used hereafter in this section. This aggregation method is preferred because it is the same as applied by (SamKnows, 2014). Some filtering is required to screen anomalous results and, indeed, most broadband speed measurement platforms apply some sort of filtering (see Table 1).



**Figure 9:** Fixed-broadband download speeds, Sweden, 2012 and 2016 – comparison of *Bredbandskollen* results against other measurement platforms.

There are large differences between measurement platforms: Ookla’s and Samknow’s averages are, respectively, 34 per cent and 69 per cent higher than *Bredbandskollen*’s. Akamai’s average is about 62 per cent lower than *Bredbandskollen*’s.

For the only two measurement platforms for which we have several years of overlapping results, Akamai and *Bredbandskollen*, the relative difference between the two averages remains rather constant with time. Indeed, Akamai speed averages are in the range of 56 to 64 per cent lower than *Bredbandskollen* in the period 2011-2016. Both platforms have large and comparable sample sizes (above 15 million observations per year), and therefore the difference observed is probably due to structural persisting disparities in the measurement methodology. For instance, Akamai's passive monitoring of speeds as opposed to *Bredbandskollen*'s user-initiated active testing.

Another relevant finding is that the average speed measured by *Bredbandskollen* is significantly lower than that obtained by SamKnows in 2012. SamKnows is often taken as ground-truth data in studies targeting countries with large SamKnows deployments, such as the United States or the United Kingdom, because of the robustness of their hardware-based methodology (Canadi et al., 2012; Riddlesden and Singleton, 2014).

In the case of Sweden, it is not warranted that SamKnows can be used as ground-truth data because of its relative small sample size in Sweden (about 0.5 million observations per year coming from 217 end users). Indeed, in the studies commissioned to SamKnows by the European Union, only partial results are published for Sweden (SamKnows, 2013, 2014).<sup>15</sup>

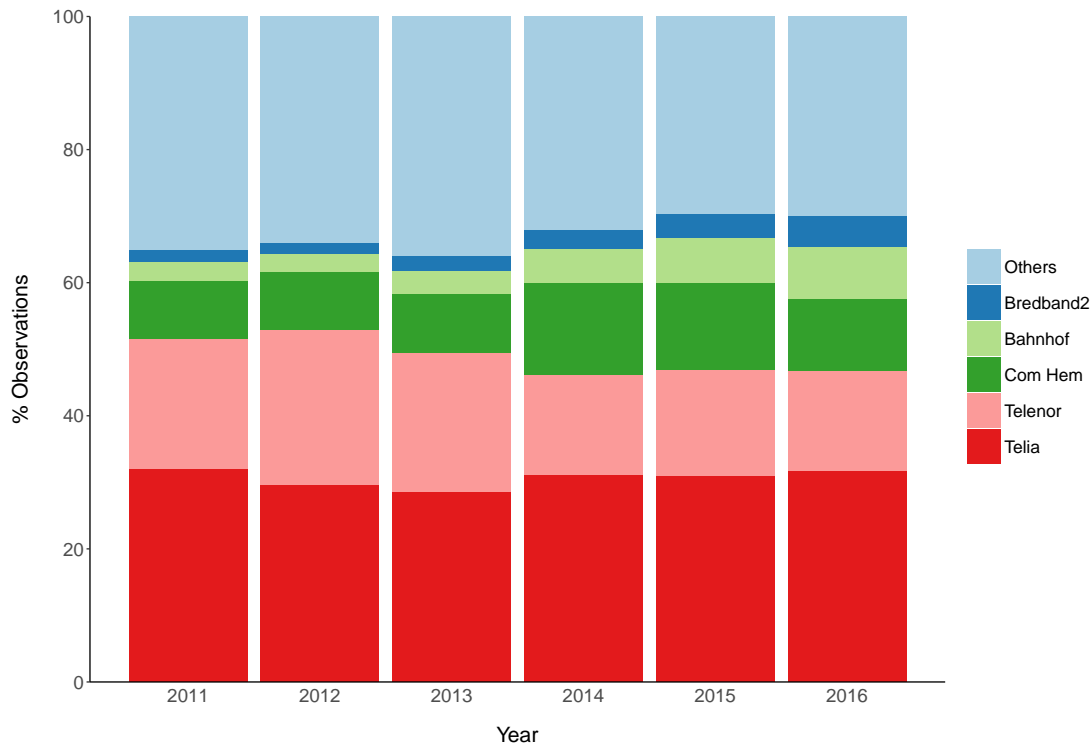
Nevertheless, SamKnows remains an important reference given that it is the only benchmark for Sweden in which possible bottlenecks within the home network (e.g. concurrent use of a single connection, poor WiFi performance) are discounted. SamKnows average speed is well above that measured by *Bredbandskollen*, which may indicate that congestion within the home network (path from 1 to 2 in Figure 3) is a relevant bottleneck in Sweden.

Figure 10 shows the breakdown of observations per operator.

---

<sup>15</sup>For 2012, SamKnows published the average download speeds for FTTx and xDSL connections in Sweden, but not for cable. For 2013, only the average download speed of FTTx connections was published. In Figure 9, the average for SamKnows is estimated by taking for cable connections in Sweden the average speed of cable connections in Europe in 2012. This is a conservative estimate, given that average speeds in Sweden tend to be above the European average.



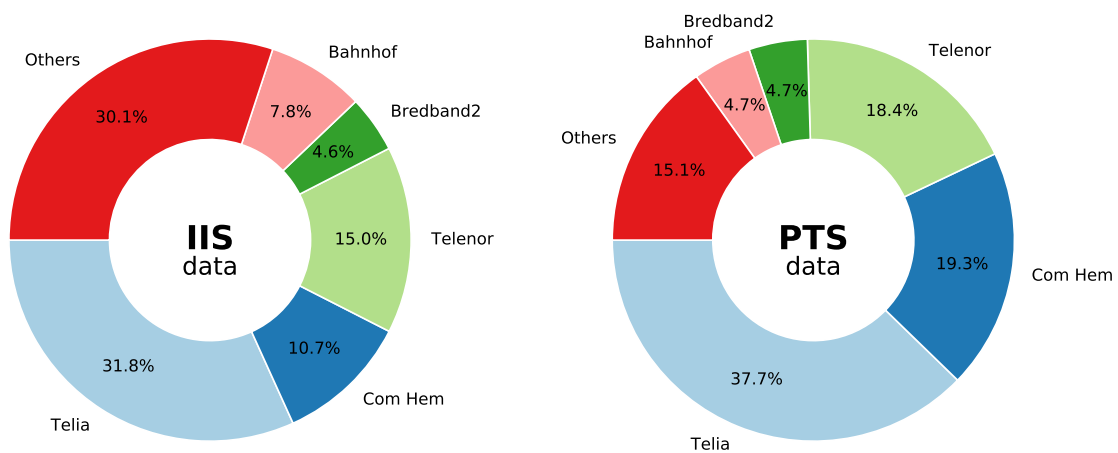


**Figure 10:** Share of *Bredbandskollen* observations per operator.

When compared to the actual market shares (Swedish Post and Telecom Authority, 2016), we note that there are some significant differences (Figure 11). Indeed, the incumbent operator, Telia, is underrepresented, as well as Com Hem. On the other hand, Bahnhof and smaller operators included in the category others are largely overrepresented.

If the population represented in *Bredbandskollen* corresponded to the whole population of fixed-broadband subscriber in Sweden, the number of tests per operator should follow a multinomial distribution with the probability of each operator being equal to the market share. That is, if there were no selection bias, each test would be drawn (with replacement) from this multinomial distribution.

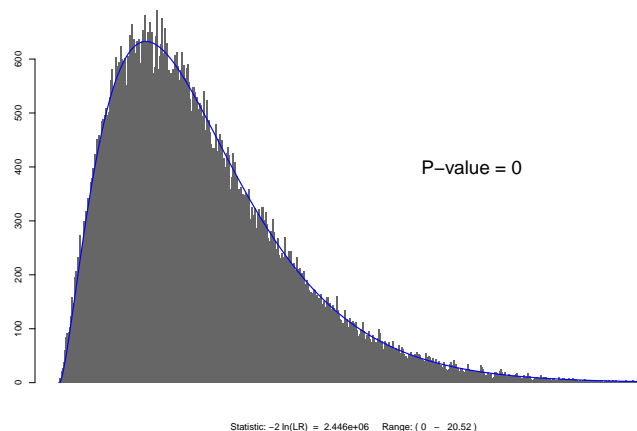
This formulation allows us to test the null hypothesis that *Bredbandskollen*'s population is the same as that of all fixed-broadband subscribers in Sweden. If confirmed, this will be an indication of no (or mild) selection bias. If rejected, there will be evidence of an important selection bias.



**Figure 11:** Share of *Bredbandskollen* observations per operator (left) and actual market shares (right), 2016.

**Source:** Swedish Post and Telecom Authority (PTS) for market shares.

Figure 12 shows the results of a chi-squared test, which strongly rejects the null. As could be expected given the large number of observations, the deviations between *Bredbandskollen*'s share of observations per operator and the actual market shares indicate a significant selection bias.

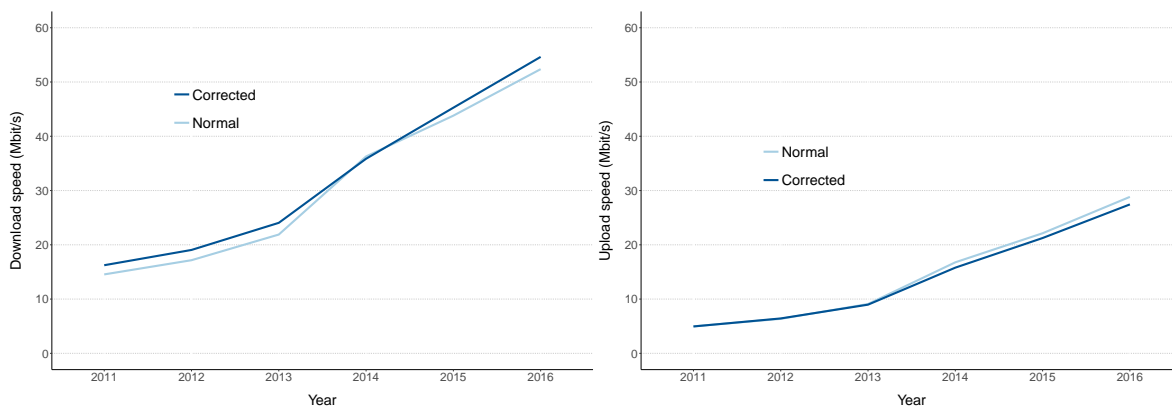


**Figure 12:** Chi-squared test: share of *Bredbandskollen* observations per operator and actual market shares.

**Source:** Swedish Post and Telecom Authority (PTS) for market shares.

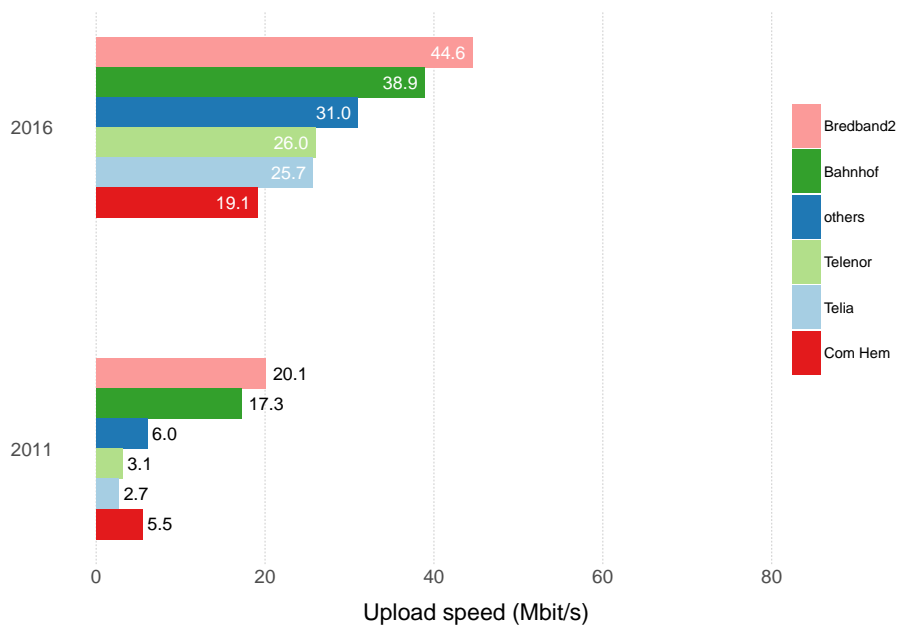
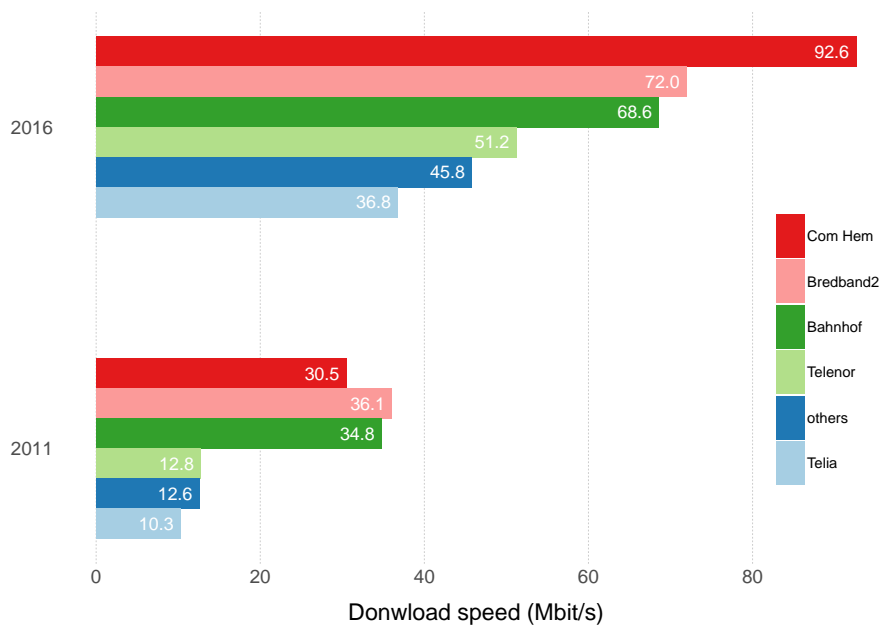
Post-stratification methods have been shown to eliminate selection bias in some non-representative samples (Zagheni and Weber, 2015). In particular, "post-stratification eliminates the bias due to selection or coverage problems if, within each adjustment cell, the probability that each case completes the survey is unrelated to the case value on the survey variable of interest" (Baker et al., 2013). In the case of broadband speed tests, this would mean that the end users of an under- or over-represented operator that decided not to perform the test did not do so because their performance would have been different from those end users from the same operator that performed the test. This is a strong assumption in our case.

Nevertheless, we test the effects of adjusting the weights given to each *Bredbandskollen*'s observation so that the final distribution corresponds to actual market shares (Figure 13). The correction has, however, a minimal effect on the aggregate values.



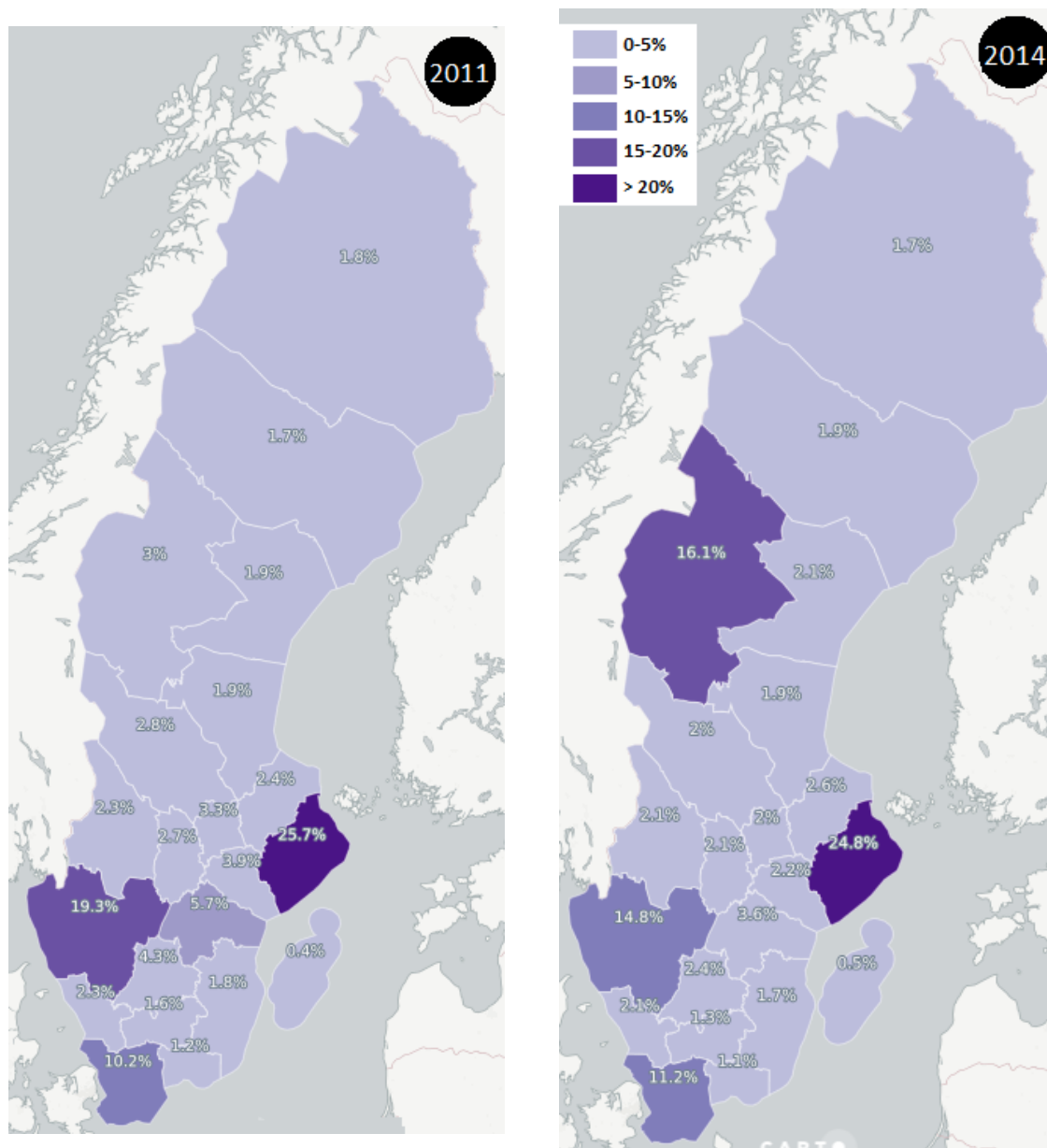
**Figure 13:** Fixed-broadband download speeds (left) and upload speeds (right), Sweden, 2011-2016 – *Bredbandskollen* adjusted test results.

The analysis of *Bredbandskollen*'s speed tests broken down by operator shows that although major differences exist between average speeds of different operators, they have opposite effects on the adjusted total average (Figure 14). That is, small operators ("others") and Bahnhof are over-represented and therefore their weight is adjusted downwards in the re-weighted average. Conversely, Telia and Com Hem are under-represented and their weight is adjusted upwards. The resulting adjusted average, however, barely changes.



**Figure 14:** Fixed-broadband download speeds (top) and upload speeds (bottom), by operator, Sweden 2011 and 2016 – *Bredbandskollen* test results.

Figure 15 and Figure 16 show the geographical dimension of the data collected by *Bredbandskollen*.

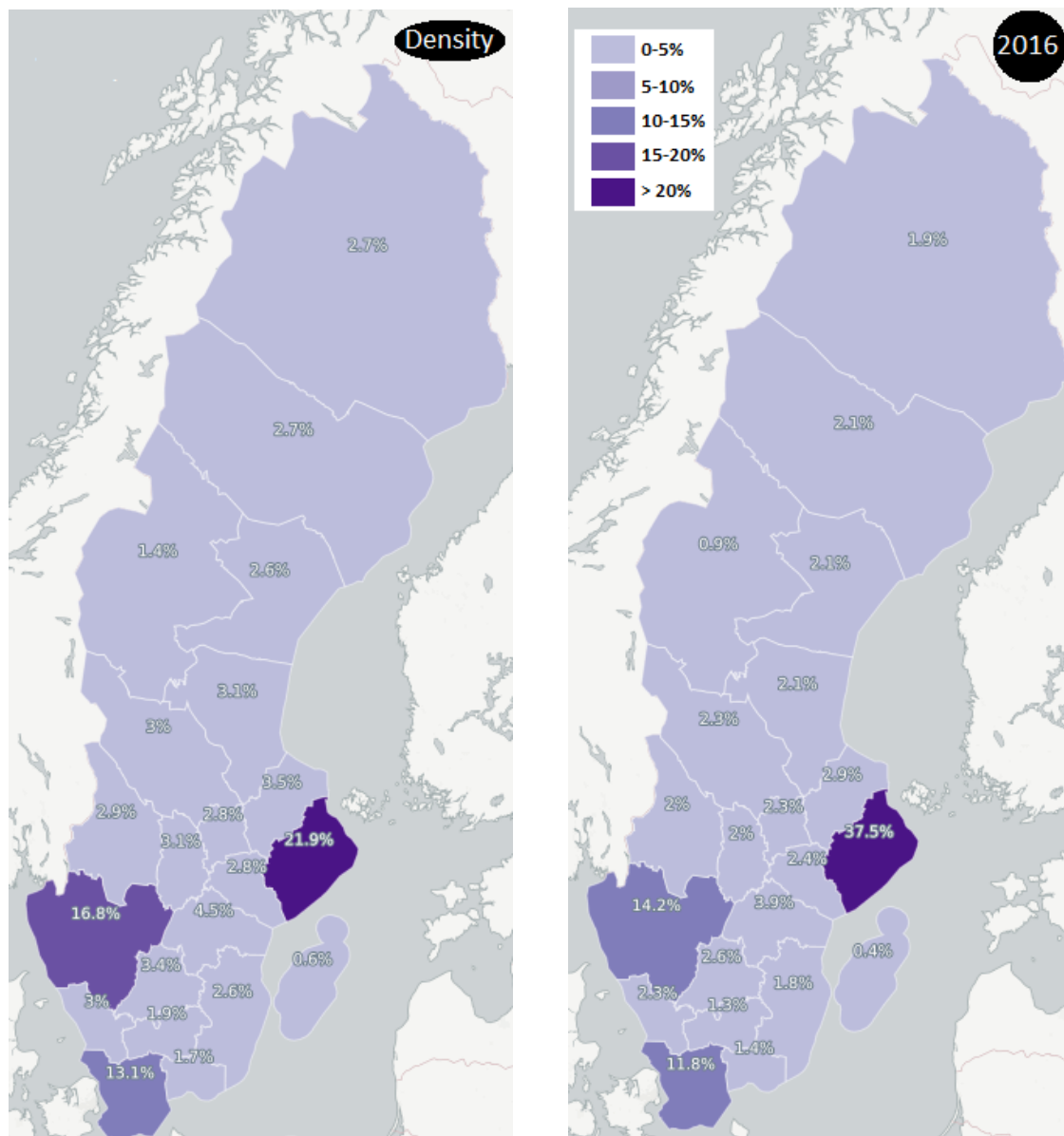


**Figure 15:** Share of *Bredbandskollen* observations per region, 2011 and 2014.

As it could be expected, the region of Stockholm concentrates most of the observations, followed by Västra Götlands region<sup>16</sup> and Skåne region.<sup>17</sup> This is in line with the population

<sup>16</sup>The second most populated region in Sweden after Stokholm. The largest city in Västra Götlands's is Gothenburg.

<sup>17</sup>The third most populated region in Sweden. Its largest city is Malmö.



**Figure 16:** Household density and share of *Bredbandskollen* observations per region, 2016.

**Note:** Household density is calculated as the number of households covered with fixed Internet service in a region over the total number of households covered by fixed Internet services in the whole country. Data on broadband coverage are sourced from the [Swedish Post and Telecom Authority](#).

density of these three regions, which concentrate most of the households with access to fixed-broadband services in Sweden.

An unexpected finding is the large number of observations from Jämtland region in 2013 and 2014. Jämtland is a sparsely populated, rural region in the middle of Sweden (see Figure 15, right). For some unknown reason, it is largely over-represented in *Bredbandskollen*'s sample in 2013 and 2014, with more than 15 per cent of total observation coming from that region in these two years. This is in stark contrast with 2015, where less than 1 per cent of *Bredbandskollen*'s observations were located in Jämtland region.

Stockholm region could also be over-represented in 2016, because it accumulates 37.5 per cent of all observations in that year whereas it only represents 21.9 per cent of the households with potential access to a fixed-broadband connection. In order to know precisely whether its representation in the sample is correct or not, we would need data on the number of fixed-broadband subscriptions broken down by region, which unfortunately PTS does not collect.

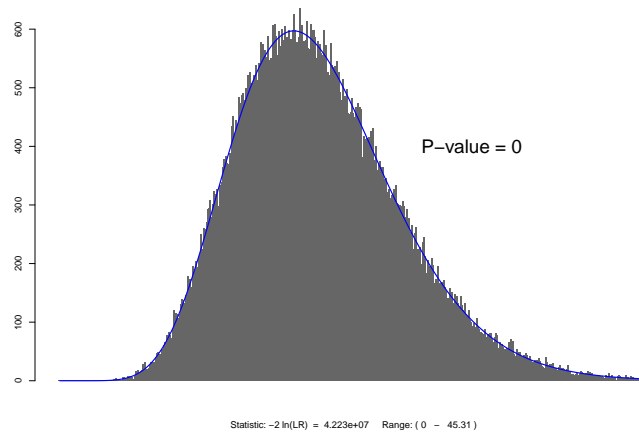
If we run a chi-squared test to check whether the sample distribution per region corresponds to the household density, the results strongly reject the null hypothesis that both come from the same distribution (Figure 17).

This is not as conclusive a sign of selection bias as that obtained with the chi-squared test ran on the observations per operator. Indeed, in this case we are assuming that the share of subscriptions per region corresponds to the share of potential customers per region. This is often not the case, as dense urban areas tend to concentrate most of the connections, whereas sparse rural areas are less connected, even proportionally to their population (see chapter 6 in International Telecommunication Union, 2016).

Nevertheless, what is definitely a sign of selection bias or measurement error (e.g. failure to geolocate appropriately a large number of observations) is the oscillation observed in the share of observations per region. This applies to both Jämtland and Stockholm regions. Fixed-broadband subscriptions cannot be expected to increase sharply given the high penetration

levels in Sweden.<sup>18</sup> Therefore, the share of observations corresponding to these regions should not change abruptly from one year to another.

The oscillations observed could be explained either by an error in geolocating the IPs of a large number of end users running the test in these regions or by a structural change in the selection bias (e.g. some active users in the Jämtland region performed many tests on *Bredbandskollen* in 2013 and 2014, but in 2015 and 2016 they stopped the tests or started using another speed measurement platform). *Bredbandskollen*'s data on unique end users could shed some light on this latter hypothesis.



**Figure 17:** Chi-squared test: share of *Bredbandskollen* observations and household density per region.

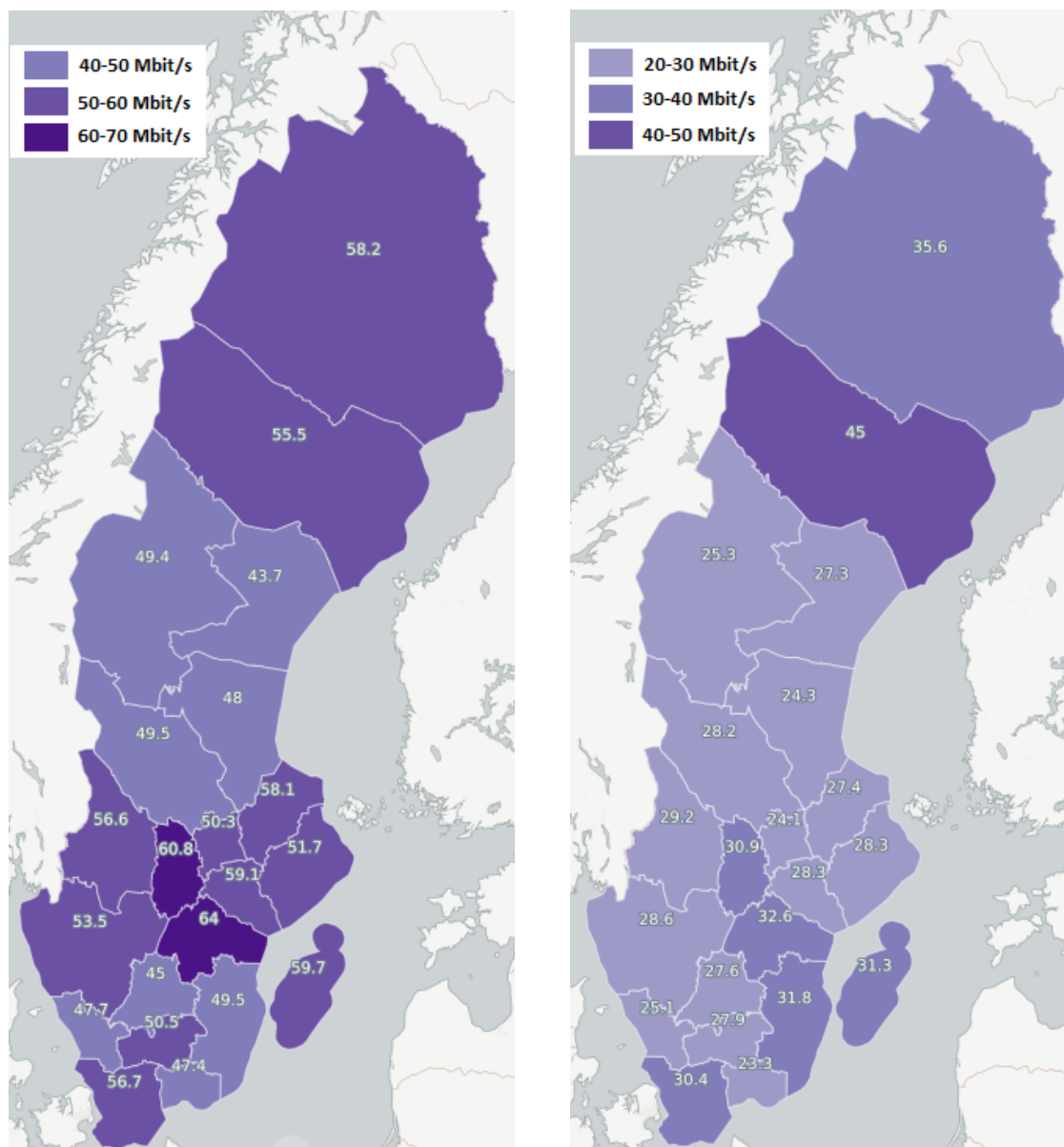
**Note:** Household density is calculated as the number of households covered with fixed Internet service in a region over the total number of households covered by fixed Internet services in the whole country. Data on broadband coverage are sourced from the [Swedish Post and Telecom Authority](#).

---

<sup>18</sup>There were 32.0 fixed-broadband subscriptions per 100 inhabitants in Sweden in 2011 and 36.1 in 2016. This is a significant progress given the high starting penetration level in 2011. However, there has not been a drastic change in the trend that could explain the oscillations observed in *Bredbandskollen*'s share of observations per region. Fixed-broadband figures sourced from [ITU](#).



Lastly, Figure 18 shows the average broadband speeds recorded by *Bredbandskollen* in each region in 2016.



**Figure 18:** Fixed-broadband download speeds (left) and upload speeds (right), by region, Sweden, 2016 – *Bredbandskollen* test results.

A striking finding from Figure 18 is that the regions with the highest population densities (i.e. Stockholm, Västra Götlands and Skåne) are not the ones with the highest average speeds. On the contrary, the highest average upload speeds (and also some of the highest download speeds) are recorded in Norrbotten and Västerbotten, two sparsely inhabited regions at the uppermost part of Sweden.

There is nothing in the analysis of *Bredbandskollen*'s observations indicating an anomaly in the representation of these regions in the sample. Therefore, this surprising result seems to be correct and would need to be cross-validated against some qualitative evidence.

Another message that emerges from the map of average speeds in Sweden is the relative equality of speeds in all regions. Moreover, this equality is grounded on a very high baseline: 40 Mbit/s download and 20 Mbit/s upload, as the minimum averages per region. This is uncommon and suggests that Sweden has had much more success than most other countries in ensuring that the benefits of high-speed broadband reach rural and sparsely populated areas.

## 5 Conclusions

The analysis of the data from *Bredbandskollen* carried out in this project has provided strong evidence of a selection bias. The lack of sufficient ground truth data on broadband speeds in Sweden does not allow to correct *Bredbandskollen*'s sample in a way to ensure its representativity for the whole population of Sweden.

If *Bredbandskollen* could share part of the information it has on the number of unique users running the tests, some post-stratification adjustments could be envisaged to improve the representativeness of the sample. Even if the information on unique test users was shared at an aggregate level to avoid any privacy concerns (e.g. per region), that would enable a wider range of tests and adjustments to tackle the issue of the selection bias.

Another piece of information that would help in the correction of the bias is the breakdown of observations per technology (i.e. xDSL, fiber and coaxial cable). *Bredbandskollen* obtains information on the technology directly from network operators, in parallel to the test. This

information is not contained in the log records shared with ITU, but would be very helpful. For instance, it would allow the linking of *Bredbandskollen* dataset with the statistics from the Swedish regulator on fixed-broadband technologies, thus adding another dimension from which to test and correct the selection bias.

The analysis has also found an unstable sample composition in the geographical distribution of *Bredbandskollen*'s observations. This may be explained by a technical issue in the geolocation of the IPs of a large number of test users. Otherwise, it may be an important bias introduced by the voluntary online sampling approach that underpins the test.

Notwithstanding the issues highlighted above, the geographical analysis of *Bredbandskollen*'s results suggest that there does not seem to be an urban/rural broadband speed divide in Sweden. Indeed, there is a relative equality of speeds in all regions, supported by a very high baseline: 40 Mbit/s download and 20 Mbit/s upload.

Even if the accuracy of the average speed measurement cannot be guaranteed given the evidence of a selection bias, the result that there are small differences between regions can still hold under the assumption that the selection bias is constant or evolves in parallel with time for all regions. Except for Jämtland and Stockholm, for which there is some evidence of instability in their sample composition, the assumption of a common evolving selection bias for the other regions is plausible.

This would confirm that Sweden is an outstanding country in terms of inclusive high-speed broadband deployments.

## References

- Baker, R., Brick, J., Bates, N., Battaglia, M., Couper, M., Denver, J., Gile, K., and Tourangeau, R. (2013). Non-probability sampling. <http://www.aapor.org/Education-Resources/Reports/Non-Probability-Sampling.aspx>. Report of the AAPOR Task Force, American Association for Public Opinion Research, Boston, MA.
- Bauer, S., Clark, D. D., and Lehr, W. (2010). Understanding broadband speed measurements. [https://groups.csail.mit.edu/ana/Publications/Understanding\\_broadband\\_speed\\_measurements\\_bauer\\_clark\\_lehr\\_TPRC\\_2010.pdf](https://groups.csail.mit.edu/ana/Publications/Understanding_broadband_speed_measurements_bauer_clark_lehr_TPRC_2010.pdf).
- Bauer, S., Lehr, W., and Mou, M. (2016). Improving the measurement and analysis of gigabit broadband networks. In *TPRC44, September 2016, Alexandria, VA*.
- Boletín Oficial del Estado, number 156, Friday 27 June (2014). <https://www.boe.es/buscar/doc.php?id=BOE-A-2014-6729>. Pages 49561-49583.
- Boyd, D. and Crawford, K. (2012). Critical questions for big data: Provocations for a cultural, technological, and scholarly phenomenon. *Information, communication & society*, 15(5):662–679.
- Canadi, I., Barford, P., and Sommers, J. (2012). Revisiting broadband performance. In *Proceedings of the 2012 ACM conference on Internet measurement*, pages 273–286. ACM.
- Chetty, M., Sundaresan, S., Muckaden, S., Feamster, N., and Calandro, E. (2013). Measuring broadband performance in south africa. In *Proceedings of the 4th Annual Symposium on Computing for Development*, page 1. ACM.
- Couper, M. (2013). Is the Sky Falling? New Technology, Changing Media, and the Future of Surveys. *Survey Research Methods*, 7(3):145–156.
- Davidsson, P. (2017). Bredbandskollen Surfshastighet i Sverige 2008-2016. Technical report, Internetstiftelsen i Sverige (IIS).
- Davies, R. (2016). Broadband as a universal service. Technical report, European Parliamentary Research Service. PE 581.977.

- Economides, N. (2008). *Public Policy in Network Industries*, chapter Handbook of Antitrust Economics. MIT Press, Cambridge, Mass.
- European Commission (2010). A Digital Agenda for Europe. <http://eur-lex.europa.eu/legal-content/EN/ALL/?uri=CELEX:52010DC0245R%2801%29>. Communication from the Commission to the European Parliament, the Council, the European Economic and Social Committee and the Committee of the Regions. COM(2010) 245 final/2.
- European Telecommunications Standards Institute (2005). Speech Processing, Transmission and Quality Aspects (STQ); User related QoS parameter definitions and measurements; Part 4: Internet access. [http://www.etsi.org/deliver/etsi\\_eg/202000\\_202099/20205704/01.01.01\\_60/eg\\_20205704v010101p.pdf](http://www.etsi.org/deliver/etsi_eg/202000_202099/20205704/01.01.01_60/eg_20205704v010101p.pdf). ETSI EG 202 057-4 V1.1.1 (2005-10).
- Federal Communications Commission (2015a). 2015 Broadband Progress Report and Notice of Inquiry on Immediate Action to Accelerate Deployment. [https://apps.fcc.gov/edocs\\_public/attachmatch/FCC-15-10A1\\_Rcd.pdf](https://apps.fcc.gov/edocs_public/attachmatch/FCC-15-10A1_Rcd.pdf).
- Federal Communications Commission (2015b). 2015 Measuring Broadband America Fixed Report. <https://www.fcc.gov/reports-research/reports/measuring-broadband-america/measuring-broadband-america-2015>.
- International Telecommunication Union (2011). Handbook for the Collection of Administrative Data on Telecommunications/ICT 2011. <http://www.itu.int/en/ITU-D/Statistics/Pages/publications/handbook.aspx>.
- International Telecommunication Union (2014). Manual for Measuring ICT Access and Use by Households and Individuals, 2014. <http://www.itu.int/en/ITU-D/Statistics/Pages/publications/manual2014.aspx>.
- International Telecommunication Union (2016). Measuring the Information Society Report 2016. <http://www.itu.int/en/ITU-D/Statistics/Pages/publications/mis2016.aspx>.
- Lehr, W., Bauer, S., and Clark, D. D. (2013). Measuring Performance when Broadband is the New PSTN. *Journal of Information Policy*, 3:411–441.

- M-Lab Research Team and others (2014). ISP interconnection and its impact on consumer Internet performance – a measurement lab consortium technical report. [https://www.measurementlab.net/publications/M-Lab\\_Interconnection\\_Study\\_US.pdf](https://www.measurementlab.net/publications/M-Lab_Interconnection_Study_US.pdf).
- Ministerio de Energía, Turismo y Agenda Digital (2017). Informe de seguimiento de los niveles de calidad de servicio. Parámetros específicos para el servicio de acceso a Internet: Velocidad de transmisión de datos conseguida. Primer trimestre de 2017. [http://www.minetad.gob.es/telecomunicaciones/es-ES/Servicios/CalidadServicio/informes/Documents/Seguimiento\\_SAI\\_T1\\_17.pdf](http://www.minetad.gob.es/telecomunicaciones/es-ES/Servicios/CalidadServicio/informes/Documents/Seguimiento_SAI_T1_17.pdf).
- Ofcom (2017). UK fixed-line broadband performance, November 2016: Research Report. <https://www.ofcom.org.uk/research-and-data/telecoms-research/broadband-research/uk-home-broadband-performance-2016>.
- Prasad, T., Shekhawat, D., Guha, S., Goveas, N., and Deshpande, B. (2016). Analysis of impartial quality measurements on indian broadband connections. In *Communication (NCC), 2016 Twenty Second National Conference on*, pages 1–6. IEEE.
- Riddlesden, D. and Singleton, A. D. (2014). Broadband speed equity: A new digital divide? *Applied Geography*, 52:25–33.
- Rohman, I. K. and Bohlin, E. (2012). Does broadband speed really matter as a driver of economic growth? investigating oecd countries. *International Journal of Management and Network Economics*, 2(4):336–356. PMID: 51888.
- Rood, H., Yoshikawa, D., Wevers, R., Post, R., Tetteroo, A.-J., Kuipers, A., and Schurmann, R. (2012). Usability of broadband performance measurements for statistical surveys. In *2012 TRPC*.
- SamKnows (2013). Quality of Broadband Services in the EU (March 2012). Technical report, European Commission.
- SamKnows (2014). Quality of Broadband Services in the EU (October 2013). Technical report, European Commission.

- Sundaresan, S., Burnett, S., Feamster, N., and De Donato, W. (2014). BISmark: A Testbed for Deploying Measurements and Applications in Broadband Access Networks. In *USENIX Annual Technical Conference*, pages 383–394.
- Sundaresan, S., De Donato, W., Feamster, N., Teixeira, R., Crawford, S., and Pescapè, A. (2012). Measuring home broadband performance. *Communications of the ACM*, 55(11):100–109.
- Sundaresan, S., Feamster, N., and Teixeira, R. (2016). Home Network or Access Link? Locating Last-Mile Downstream Throughput Bottlenecks. In *PAM 2016 - Passive and Active Measurement Conference*, pages 111–123, Heraklion, Greece.
- Swedish Post and Telecom Authority (2016). The Swedish Telecommunications Market – First Half-year 2016. <https://www.pts.se/upload/Rapporter/Tele/2016/Swedish-Telekommunikations-Market-first-half-2016.pdf>.
- Telecom Regulatory Authority of India (2006). Regulation quality of service standards for broadband services.
- Telecom Regulatory Authority of India (2014). Quality of service of broadband service (second amendment) regulations. [http://www.trai.gov.in/sites/default/files/201406250411117962118scan0004\\_0.pdf](http://www.trai.gov.in/sites/default/files/201406250411117962118scan0004_0.pdf).
- Tiru, M. (2016). Big Data for Measuring the Information Society. In *14th World Telecommunication/ICT Indicators Symposium*, Botswana.
- TNS Opinion & Social (2014). Special Eurobarometer 414 “E-Communications and Telecom Single Market Household Survey”. Technical report, European Commission.
- United Nations Economic and Social Council (2016). Statistical Commission, Forty-seventh session, Report of the Inter-Agency and Expert Group on Sustainable Development Goal Indicators. <https://unstats.un.org/unsd/statcom/47th-session/documents/2016-2-IAEG-SDGs-Rev1-E.pdf>. E/CN.3/2016/2/Rev.1 (8-11 March 2016).

- Wattegama, C. and Kapugama, N. (2011). Volunteer computing model prospects in performance data gathering for broadband policy formulation. *Communications and Strategies*, (81):153–174.
- Whitacre, B., Gallardo, R., and Strover, S. (2014). Broadband’s contribution to economic growth in rural areas: Moving towards a causal relationship. *Telecommunications Policy*, 38(11):1011–1023.
- Zagheni, E. and Weber, I. (2015). Demographic research with non-representative internet data. *International Journal of Manpower*, 36(1):13–25.
- Zuhyle, S. and Mirandilla-Santos, G. (2015). Measuring broadband performance: Lessons learnt, challenges faced. In *Proceedings of CPRSouth 2015, Taiwan*.