

Visualizing Internet-Measurements Data for Research Purposes: the NeuViz Data Visualization Tool

Original

Visualizing Internet-Measurements Data for Research Purposes: the NeuViz Data Visualization Tool / Futia, Giuseppe; Enrico, Zimuel; Basso, Simone; DE MARTIN, JUAN CARLOS. - (2013). (Intervento presentato al convegno Congresso Nazionale AICA 2013 tenutosi a Fisciano (Italia) nel 18-20 Settembre 2013).

Availability:

This version is available at: 11583/2516321 since:

Publisher:

Published

DOI:

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

Visualizing Internet-Measurements Data for Research Purposes: the NeuViz Data Visualization Tool

Giuseppe Futia¹, Enrico Zimuel², Simone Basso¹, Juan Carlos De Martin¹

¹Nexa Center for Internet & Society, Politecnico di Torino (DAUIN),
Corso Duca degli Abruzzi, 24, 10129 Torino (TO), Italy
{giuseppe.futia, simone.basso, demartin}@polito.it

²R&D Department, Zend Technologies Inc.,
19200 Stevens Creek Blvd, Cupertino, CA 95014, USA
enrico@zend.com

Abstract. *In this paper we present NeuViz, a data processing and visualization architecture for network measurement experiments. NeuViz has been tailored to work on the data produced by Neubot (Net Neutrality Bot), an Internet bot that performs periodic, active network performance tests. We show that NeuViz is an effective tool to navigate Neubot data to identify cases (to be investigated with more specific network tests) in which a protocol seems discriminated. Also, we suggest how the information provided by the NeuViz Web API can help to automatically detect cases in which a protocol seems discriminated, to raise warnings or trigger more specific tests.*

Keywords: data visualization, network performance, big data.

1. Introduction

The Internet is a cornerstone of our societies and has been enabling unprecedented levels of social interaction, content sharing, business creation, as well as innovation in many fields. As Frischmann argues convincingly, one of the main reasons why the Internet is so relevant for us is that the Internet is an *infrastructural resource*, i.e., a shared piece of infrastructure that is typically managed as a *commons* in a non-discriminatory way [Frischmann, 2012].

However, the Internet is not an infrastructural resource as a fact of nature, or because of an immutable, technological law; the current status of the Internet is, instead, the consequence of specific choices, both private and public, that could very well change over time. For example the policy decision of who (the State or the Internet Service Providers) should finance (and under which conditions) the so-called 'Next Generation Networks' (NGNs) has the potential of radically changing the landscape.

In fact, many parties (including the authors of this contribution) believe that, if States allow the Internet Service Providers (ISPs) to implement premium services to collect more money and finance NGNs, the infrastructural-resource characteristics of the Internet may become less relevant, and the Internet may

lose part of its *generativity* (i.e., the property of enabling more and more people to write and distribute software and/or media content [Zittrain, 2009]).

To be fair, there is little empirical evidence supporting most policy positions on both sides of the debate. On the one hand, for instance, it is hard to prove empirically *ex ante* that allowing ISPs to implement premium services will reduce the generativity of the Internet. On the other hand, there is surprisingly little evidence backing the ‘bandwidth hogs’ argument (i.e., the argument that there is a little number of people that consume most bandwidth). The Internet policy debate, in general, is so ill informed by poor data, by missing data, and by data provided by one single stakeholder that – we agree with Palfrey and Zittrain – there is a need for more, better data to anchor the debate to solid foundations and move forward [Palfrey and Zittrain, 2011].

This is indeed starting to happen: more and more network measurement tools and visualizations, in fact, are being developed by researchers and companies worldwide. Many of such tools and visualizations are hosted by Measurement Lab (<http://measurementlab.net/>), an umbrella project run by the Open Technology Institute and the PlanetLab Consortium, and supported by academic partners and companies such as Google.

In this paper, in particular, we propose NeuViz (Neubot Visualizer), an architecture that allows us to process and visualize the data collected by Neubot, the network neutrality bot [Basso et al, 2011a], one of the tools hosted by Measurement Lab. Neubot – a project of the Nexa Center for Internet & Society – is a centrally-coordinated bot that runs the in background on the user computer and periodically runs network-performance tests that currently emulate HTTP and BitTorrent, and, in future, will emulate other protocols.

The purpose of NeuViz is to visualize and navigate Neubot data through its Web user interface, to search for cases (to be investigated with more specific network tests) in which a protocol seems discriminated. Also, NeuViz is designed to help, in the future and with a more advanced Neubot architecture, to automatically detect cases in which a protocol seems discriminated, to raise warnings or trigger more specific tests.

Many existing visualization architectures are based on cloud services and allow to query the data on demand using SQL-like query languages; compared to such visualization tools, NeuViz is much more optimized for the specific purpose of visualizing network measurement data. We designed, in fact, a robust, scalable backend architecture to support special-purpose, complex data analysis, in which the query (or the filtering algorithm) is executed in advance on the network-experiments dataset, and in which the result is stored in one (or more) NoSQL database(s), for fast data access.

We evaluate our work by loading into NeuViz the results of two Neubot network tests (Speedtest and BitTorrent) collected in the January 2012 - May 2013 period. We show that NeuViz helps us to effectively navigate Neubot data to identify cases in which a protocol seems discriminated. Also, we suggest that the information provided by the NeuViz Web API can help to automatically detect cases in which a protocol seems discriminated.

The rest of this paper is organized as follows. In Section 2 we describe related network measurement tools and visualizations. In Section 3 we describe Neubot and the Neubot data that we used in this paper. In Section 4 we

describe the NeuViz architecture. In Section 5 we describe our implementation choices. In Section 6 we describe what we learnt from browsing Neubot data with NeuViz. In Section 7 we draw the conclusions, and we describe future developments.

2. Related Work

In this section we mention the related tools and visualizations. Some of the tools that we mention (including Neubot) are hosted by Measurement Lab (M-Lab) [Dovrolis et al, 2010], a distributed server platform that also provides advanced services (e.g., the possibility of querying the hosted tools data using BigQuery, a RESTful service to query big datasets using an SQL-like query language – <https://developers.google.com/bigquery/>).

We mention four tools similar to Neubot: (i) Glasnost, which is an M-Lab based tool that compares a certain protocol flow (e.g., BitTorrent, Emule) with a *reference flow* to detect traffic shaping and its cause (e.g., the port number, the payload) [Dischinger et al, 2010]; (ii) NDT, the Network Diagnostic Tool, which measures the download and upload speed between the user computer and a Measurement Lab server [Carlson, 2003] and records the TCP state using Web100 [Mathis et al, 2003]; (iii) SpeedTest.net (<http://www.speedtest.net/>), which estimates the broadband speed of the user's connection using parallel HTTP flows; (iv) Grenouille, which measures the round trip delay, the download speed, and the upload speed (<http://grenouille.com/>).

Differently from Glasnost Speedtest.net, and NDT – which run on-demand tests – Neubot and Grenouille run tests in the background; however, Neubot uses diverse protocols, while Grenouille focuses on the performance only.

We mention four visualizations similar to NeuViz: (i) the two visualizations of Glasnost data developed by, respectively, the Syracuse University School of Information Studies (<http://dpi.ischool.syr.edu/MLab-Data.html>), and the Open Knowledge Foundation (<http://netneutralitymap.org/>); (ii) the many BigQuery-based NDT visualizations by the M-Lab team (e.g., <http://goo.gl/m9WbS> and <http://dmadev.com/2012/11/19/>); (iii) the SpeedTest.net visualizations available at NetIndex.com (<http://www.netindex.com/>); (vi) the Grenouille visualizations available at <http://grenouille.com/>.

Similarly to the NDT visualizations NeuViz is based on the world map; however, NeuViz is optimized for complex data analysis and uses precomputed data, while the NDT visualizations are more interactive and fetch the data from BigQuery on demand. Also, the aim of NeuViz is similar to the aim of the Glasnost visualizations; both, in fact, aim at making access networks more transparent by, respectively, showing anomalies and alleged shaping.

3. Neubot and Neubot data

Neubot is a free-software Internet bot that performs active, lightweight network-performance tests [Basso et al., 2011a]. Once installed on the user's computer, Neubot runs in the background and every 30 minutes performs active transmission tests with servers hosted by Measurement Lab. To coordinate the botnet composed of all the Neubot instances worldwide, there is the so-called

Master Server, which suggests each Neubot the next test to run as well as the default test parameters. Currently, the Master Server does not optimize the suggestions returned to each Neubot; however, as we will show, the information returned by NeuViz could help the Master Server to implement more dynamic policies.

Neubot implements three network performance tests: Speedtest, BitTorrent, and RawTest. Speedtest measures the network performance using the HTTP protocol, BitTorrent measures the network performance using the BitTorrent protocol, and the RawTest test measures raw, TCP-level performance (hence the name of the test). In this paper we only describe the Speedtest and the BitTorrent tests, because we are mainly interested to use NeuViz to find cases in which a protocol seems discriminated.

Speedtest is an HTTP-based test – originally inspired to the test of SpeedTest.net, hence the test name – that downloads and uploads data using a single HTTP connection [Basso et al., 2011b]. The test measures the download and the upload speed at the application level. Also, the test estimates the base Round Trip Time (RTT) using as a proxy the time that the connect system call takes to complete (later indicated as *connect time*). The test transfers a number of bytes that guarantees that each phase of the test (download, upload) lasts for about five seconds.

The BitTorrent test is similar to the Speedtest test, except that it uses the BitTorrent protocol (http://www.bittorrent.org/beps/bep_0003.html) instead of the HTTP protocol. However, while Speedtest makes a single GET request for a large-enough amount of data, BitTorrent – to better emulate the BitTorrent protocol – downloads many small chunks in a request-response fashion and, to approximate a continuous transfer, makes many back-to-back requests at the beginning of the test.

Measurement Lab (which hosts Neubot on its servers) periodically collects the Neubot experiments results saved on its servers and publishes such results on the Web (<http://measurementlab.net/data>) under the terms and conditions of the Creative Commons Zero 1.0 Universal license. We mirrored the data provided by Measurement Lab, and we converted such data in CSV format, generating CSV files that contain one month of data each. To prepare this paper, we imported into NeuViz the CSV files from January 2012 to May 2013 (reading 5,383,376 test, from 4,037 Neubot clients worldwide, for a total of 1.5 GB).

Each CSV file contains the following fields (the type is indicated in parentheses): client address (str); connect time, in second (float); download speed, in byte/s (float); Neubot version (str); operating system platform (str); server address (str); test name (str: “speedtest” or “bittorrent”); timestamp of the test, i.e., the number of seconds since 1970-01-01 00:00 UTC (int); upload speed, in byte/s (float); unique identifier of the Neubot instance (str).

4. Description of the NeuViz Architecture

Fig. 1 shows the NeuViz architecture, which is a pipeline that processes data provided by *Producers*, and which organizes the data such that *Consumers* can visualize (or further process) such data. The pipeline is composed of a *Backend*

and a *Frontend*: the Backend receives data from many Producers and processes such data to allow for efficient visualization; the Frontend is a Web interface that visualizes the data. In the middle there is a Web API.

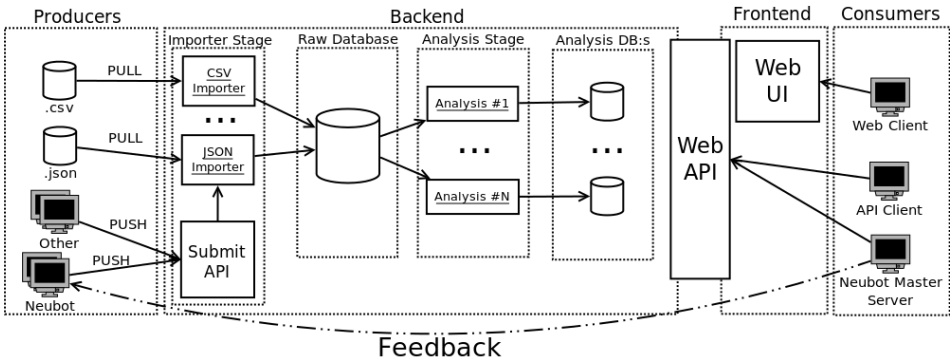


Figure 1: NeuViz architecture

4.1 Producers

As a first approximation a Producer is a static dataset. For example, in this paper we used Neubot data expressed in CSV format and in the future we may want to import datasets from other projects (e.g., SpeedTest.net) and encoded in many formats (e.g., JSON).

NeuViz also includes a Submit API, which allows network-experiment tools (e.g., Neubot and possibly other network-measurement tools) to push the result of their experiments just after the experiments are run. We added the API because we want to create a feedback loop in which data posted by Neubot is processed by NeuViz and consumed by the Master Server to provide better suggestions to Neubot instances.

4.2 The Backend

The Backend is composed of two processing stages (the *Importer Stage* and the *Analysis Stage*), each followed by a database stage (respectively, the *Raw Database* and the *Analysis Databases*).

The Importer Stage is composed of many modules: each module imports data from a specific network-measurement tool (e.g., Neubot, Grenouille) and format (e.g., CSV, JSON). Also, network-measurement tools can push their data via the Submit API. The problem of whether (and how) to authenticate the tool submitting its data is not discussed in this paper.

The Raw Database receives heterogeneous data organized in a uniform format (e.g., JSON) by the Importer Stage. It is not practical to reduce all the input data to the same schema, because each network experiment saves different metrics; therefore, NoSQL seems the best technology to implement the Raw Database. Also, the Raw Database shall be scalable-enough to handle continuous streams of data posted on the Submit API by Neubot (and possibly by other network measurement tools).

The Analysis Stage is a collection of modules that periodically fetch data from the Raw Database and process it to produce the aggregate data needed for the visualizations. We already implemented a world-map-based visualization and we are working on other visualizations. Of course, the Analysis Stage needs to be scalable, because we need to process multiple times the raw data stored into the Raw Database. Also, the Analysis Stage should minimize the computational cost of adding the results of new experiments to NeuViz.

The Analysis Databases are a number of conceptually-separated databases that store data which is ready to be visualized on the NeuViz Frontend with minimal computational cost. We want, in fact, to allow the user to visualize and browse the data as seamlessly as possible.

4.3 The Web API, the Frontend, and Consumers

The Web API connects the Backend and the Frontend (i.e., the Web interface of NeuViz). A Web browser that uses the NeuViz Web interface is the default Consumer; however, other Consumers are possible. For example, as mentioned above, as part of our future work we plan to extend the Master Server to fetch data from NeuViz to provide better suggestions to the Neubot instances.

In this paper we don't discuss whether and how the access to the Web API and to the Frontend should be restricted. This will possibly be the subject of a future work.

5. Implementation Choices

In this section we describe the implementation of the first NeuViz prototype (<https://github.com/neubot/neuviz>), and we explain our implementation choices.

5.1 Backend

We implemented the Importer Stage and the Analysis Stage in Python. The Importer Stage converts the input into a JSON document and adds geolocation information, if needed. The Analysis Stage prepares the data for a world-map-based visualization that allows the user to navigate the spatial level (countries, and cities) and the temporal level (hour of the day). The output of the Analysis Stage is a set of JSON documents.

For both the Raw Database and the Analysis Databases we use MongoDB, a NoSQL database that indexes JSON documents and that is very often deployed in big data scenarios [Moniruzzaman and Akhter, 2013]. To improve performance, we exploit the *indexes* feature of MongoDB to speed up the query execution. We imported 1.5 GiB of Neubot CSV files, from January 2012 to May 2013, and we stored such data into a MongoDB database, processing about 5.3M samples in less than 60 minutes. We run the Importer on a laptop with an Intel Core i7 CPU at 2.0 Ghz, with 8 GB of RAM, and a 256-GB SSD, running GNU/Linux 3.5.0.

We geolocated the incoming data using the GeoLite Free Geolocation Database (<http://dev.maxmind.com/geoip/legacy/geolite/>). As explained in the GeoLite website, when the database is not up-to-date, the geolocation loses

1.5% of accuracy each month because IP addresses are re-assigned. To minimize the damages caused by out-of-date GeoLite databases, we never used databases older than two months when we geolocated Neubot results.

In the Analysis Stage we computed the median number of tests, the median number of Neubot instances, the median connect time, the median download speed, and the median upload speed, for both Speedtest and BitTorrent, for each hour of the day, for each month, and along the geographical dimension (i.e., country, city), and the business dimension (i.e., ISP). We decided to use the median, which is a common index used to analyze network traffic, to avoid the risk that few outliers could dominate our indexes.

For future scalability we designed the code in a way that potentially allows us to use MapReduce techniques on cloud services (e.g., Amazon Elastic MapReduce, <http://aws.amazon.com/elasticmapreduce/>). To this end we divided the Importer and the Analysis code into a *map* step and a *reduce* step.

5.2 The Web API and the Frontend

To access the NeuViz API, the user sends the following HTTP/1.1 request: GET /neuviz/1.0/<viz>/<params>, where <viz> is the name of the visualization, and <params> is a placeholder for (possibly-empty) parameters. The returned JSON contains a recursive set of dictionaries that represent the geographical dimension (country, city) and the business dimension (ISP). The leaves are dictionaries that contain the following hour-wide median statistics for the Speedtest and the BitTorrent tests: download speed, upload speed, connection time, number of Neubot instances, number of tests.

The Web Interface, written using D3.js (<http://d3js.org/>), allows the user to explore different network measurement performances at different geographic dimensions (country, cities, and ISPs). For simplicity, and since it does not seem to cause any performance issue, we currently use the Web interface to compute some statistics, e.g., the difference between the median Speedtest download speed and the median BitTorrent download speed that we use in Section 6.2 to compare the performance of BitTorrent and Speedtest.

6. Results

In this section we report what we learnt from using NeuViz to browse Neubot data, both in terms of number of tests and in terms of performance.

6.1 Number of Neubot Tests

Fig. 2 shows the visualization of the number of tests per country and per hour. The alpha channel of the country color indicates the median number of tests per country. The visualization, in particular, shows the median number of tests performed between 9:00 PM and 10:00 PM (local time) in April 2013. The selected country is Canada, in which the median number of tests performed is indicated by the number in the bottom right corner (1084).

By selecting other countries in the visualization, we have seen that the countries with more median tests per hour between 9:00 PM and 10:00 PM in

April 2013 are: the US (4223); Italy (2866); Germany (2285); and Canada (1084). Other countries have less tests per hour.

The availability of the number of tests per country is interesting because, by knowing the number of tests per country, the Master Server could maximize the test coverage; e.g., it can increment the frequency of testing on countries where there are few Neubot users.

6.2 Comparison of Speedtest and BitTorrent performance

Before studying the visualization that shows the difference between the Speedtest and the BitTorrent test download and upload speeds, we checked that the Speedtest and the BitTorrent connect times were 'comparable'. To this end we arbitrarily define 'comparable' two median connect times whose difference is smaller than five millisecond, in our experience a reasonable threshold for this kind of analyses.

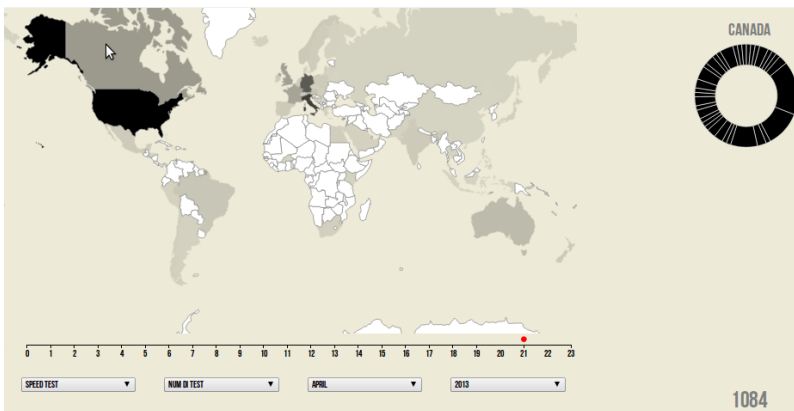


Figure 2: NeuViz interface of the worldwide map for Neubot data of April 2013

The visualization of the difference between the median BitTorrent connect time and the median Speedtest connect time shows, surprisingly, that in Italy such difference is always positive and often greater than five millisecond (i.e., the Speedtest connect time is typically lower). Italy is the only country in which, for 2013 data, we noticed this behavior.

Also we noticed interesting things from the comparison of the median upload speed in countries in which the median connect times are comparable. We noticed, in fact, that in 2013 the median upload difference between Speedtest and BitTorrent in Canada was very often positive, while the same difference was very often negative in the US (see Fig. 3).

Moreover, when comparing the download speeds in countries in which the connect times are comparable, we also noticed that the US Speedtest download speed is always lower (in median) than the BitTorrent one for every hour of the day and for every month of 2013. Interestingly, instead, the download speeds are comparable in Italy, in which – as we have seen – there is a connect time bias in favor of Speedtest.

The above observations lead us to speculate that: (a) BitTorrent is slightly faster than Speedtest; (b) in Italy the two tests are comparable because of the connect-time bias that we observed; (c) the BitTorrent upload speed seems to be discriminated in Canada. Of course, these are only hypotheses that need to be verified (or contradicted) by more detailed experiments.

6.3 Concluding Remarks

Despite being still in beta stage, NeuViz allowed us to discover the three diverse network anomalies we described in 6.2. In the future, a more advanced Master Server could learn, from the NeuViz API, about similar anomalies and ask Neubot instances that are near the anomalies to gather more information needed to investigate the anomalies (e.g., by capturing packets to gather RTT samples useful to understand whether there is a connect-time bias).

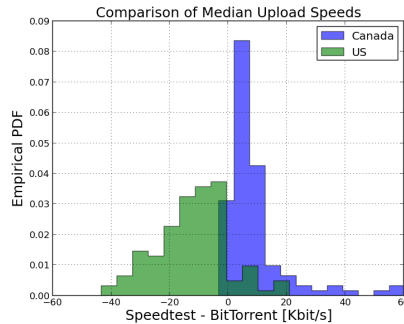


Figure 3: the Empirical Probability Density Function (PDF) of the difference of the median upload speed of US and Canada

7. Conclusion and Future Work

In this paper we described NeuViz, an architecture that allows us to process and visualize the data collected by Neubot, the active, network-measurement tool developed by the Nexa Center for Internet & Society. The purpose of NeuViz is to visualize and navigate Neubot data through its Web user interface, to search for cases (to be investigated with more specific network tests) in which a protocol seems discriminated.

Differently from other visualization architectures NeuViz is much less flexible and much more optimized, on purpose. NeuViz, in fact, executes the queries in advance and the result is stored into one or more NoSQL databases (using MongoDB), for fast data access. The Backend of NeuViz, written in Python, is structured to ease the task of porting it to a cloud-based MapReduce solution, for future scalability. The Web interface Frontend of NeuViz shows a world-map-based visualization of Neubot results implemented using the D3.js library.

To evaluate NeuViz we loaded one-year-and-a-half records collected by two network tests periodically run by Neubot, called Speedtest (based on HTTP) and BitTorrent. We showed that NeuViz effectively helped us to identify cases (to be investigated with more specific network tests) in which a protocol seems

discriminated. In our discussion we also suggested how the Web API of NeuViz can help to automatically detect cases in which a protocol seems discriminated, to raise warnings or trigger more specific tests (by cooperating with the Master Server of Neubot). As part of our future work we plan to extend NeuViz to automatically raise warnings and to cooperate with the Master Server of Neubot to trigger more-specific network experiments.

Acknowledgments

The first prototype of the NeuViz project has been developed as final project of the BigDive course 2013 (<http://www.bigdive.eu>). We would like to thank Christian Racca of the TOP-IX Consortium and all the staff and teachers of the BigDive course for their support during the development of this project.

References

- [Basso et al, 2011a] Basso S., Servetti A., De Martin J. C., The network neutrality bot architecture: A preliminary approach for self-monitoring of Internet access QoS, in Proc. of the Sixteenth IEEE Symposium on Computers and Communications, Corfu, Greece, 2011.
- [Basso et al, 2011b] Basso S., Servetti A., De Martin J. C., The hitchhiker's guide to the Network Neutrality Bot test methodology, in Proc. of Congresso Nazionale AICA 2011, Torino, 2011.
- [Carlson, 2003] Carlson R., Developing the Web100 Based Network Diagnostic Tool (NDT), In Proc of the Passive and Active Measurement Conference, 2003.
- [Dischinger et al, 2010] Dischinger M., Marcon M., Guha S., Gummadi K. P., Mahajan R., Saroiu S., Glasnost: Enabling End Users to Detect Traffic Differentiation, in Proc. of USENIX Symposium on Networked Systems Design and Implementation, 2010.
- [Dovrolis et al, 2010] Dovrolis C., Gummadi K. P., Kuzmanovic A., Meinrath S., Measurement Lab: Overview and an Invitation to the Research Community, ACM SIGCOMM Computer Communication Review, 40, 3, 2010, 53–56.
- [Frischmann, 2012] Frischmann B. M., Infrastructure: The Social Value of Shared Resources, Oxford University Press, 2012.
- [Mathis et al, 2003] Mathis M., Heffner J., Reddy R., Web100: Extended TCP Instrumentation for Research, Education and Diagnosis, ACM SIGCOMM Computer Communication Review, 33, 3, 2003, 69–79.
- [Moniruzzaman and Akhter, 2013] Moniruzzaman A. B. M., Akhter H. S., NoSQL Database: New Era of Databases for Big data Analytics - Classification, Characteristics and Comparison, International Journal of Database Theory and Application, Vol. 6, No.4, 2013.
- [Palfrey and Zittrain, 2011] Palfrey J., Zittrain J., Better Data for a Better Internet, Science, 334, 6060, 2011, 1210-1211.
- [Zittrain, 2009] Zittrain J., The future of the Internet--and how to stop it., Yale University Press, 2009.