# Data-driven Networking Research: models for academic collaboration with industry (a Google point of view)

Jeffrey C. Mogul
Priya Mahadevan
Christophe Diot
Google
network-data-sharing@google.com

John Wilkes
Phillipa Gill
Amin Vahdat
Google

## ABSTRACT

We in Google's various networking teams would like to increase our collaborations with academic researchers related to data-driven networking research. There are some significant constraints on our ability to directly share data, which are not always widely-understood in the academic community; this document provides a brief summary. We describe some models which *can* work – primarily, interns and visiting scientists working temporarily as employees, which simplifies the handling of some confidentiality and privacy issues. We describe some specific areas where we would welcome proposals to work within those models.

## CCS CONCEPTS

• **Networks → Network measurement**; **Network experimentation**;

## KEYWORDS

Network data, research collaboration

## 1 MOTIVATION

Google depends on networking, and we owe a tremendous debt to the academic research community that has provided the networking innovations we use. Many of us are former academic researchers ourselves, and Googlers actively participate in the research community in various modes, including collaborations with people in academia.

We frequently hear from academic friends that they struggle to do realistic networking research at scale, not just because they lack access to large-scale infrastructure (although that is a problem) but also because they lack access to real-world, large-scale data: workloads, traffic traces, failure types and statistics, etc. "Could Google please share more data with us" is a common request.

We collectively (Google + academia) would benefit from better collaborative, data-driven networking research: academic researchers would be able to test their work against real-world conditions, and Google is better off when academic research is grounded in reality, since that makes the research results more applicable to our needs.

## 2 CONSTRAINTS ON DATA SHARING

Unfortunately, we have seldom actually been able to do this collaborative research – not for lack of desire or a lack of interesting problems, but because Google operates under some constraints that might not be widely understood. We briefly explain these constraints, to help potential academic partners understand how to propose successful collaborations. (Google is not the only company that has constraints, but our constraints might be different from those at other companies, and it might be useful to potential partners to understand why.)

These constraints fall into a few general categories:

- **Privacy**: User privacy is Google's primary concern. In fact, while we allow interns and visiting scientists (mostly) full access to our confidential technical information, such as source code, these temporary employees never have direct access to user data. With the introduction of strict privacy regulations such as GDPR, and our active participation in the cloud market, a surprisingly broad set of networking data now counts as "user data." (For example, traffic patterns for a cloud user, or for an Internet prefix, that could be re-identified using data from non-Google sources.)
  As we discuss later, anonymization is unlikely to allow us to release much data for public use.
- **Business concerns**: Many companies, including Google, avoid revealing the details of internal technologies that allow us to compete successfully. Network-related data, even if fully anonymized, can still expose technical information we would prefer not to reveal, or business-growth information that might affect investors (and hence subject to strict government regulation).
- **Scale**: We collect a lot of infrastructural data, but at our scale, the collection infrastructure itself is expensive to build and operate. This makes it hard for us to start collecting novel kinds of information – and sometimes researchers assume that we are already collecting data that we might not actually be collecting.
- **Operational risk**: Creating special-case data-collection mechanisms (for example, to get around scaling problems) can lead to operational risks. It is even more difficult to "just try out an idea" generated in academia; even seemingly-safe changes (e.g., to a congestion-control algorithm) create substantial risk.

- **Staff time**: Even when the other constraints have allowed us to release datasets in the past, we have found that this consumes months of someone's time – working through these constraints while producing high-quality, well-documented datasets.

The first two are "hard constraint"; the others are "softer" in the sense that they can be mitigated through effort.

## 2.1 Anonymization is not a solution by itself

Privacy cannot be solved by data anonymization alone. Anonymization seldom works in practice for networking data [1]; either it is too easy to de-anonymize (at least partially), and/or too much information is lost in the anonymization process, which generally makes it impossible to interpret anonymized data in a meaningful way for research purposes. In addition, the risk of getting anonymization wrong is high. It would be easier to share data if there were a proof of appropriate levels of privacy with no way to reverse engineer. Obviously, delivering such a proof for domain-specific data and individual company "schemas" is somewhere between difficult and impossible by itself. Anonymization ends up being both too labor intensive and too risky (it is hard to predict what the data analysis will reveal) for most companies to take that path. NDAs and other legal tools are sometimes necessary, but never sufficient.

## 3 MODELS FOR COLLABORATIVE DATA-DRIVEN RESEARCH

With the constraints described above in mind, we have found several models for collaborative data-driven research that do work (and we are open to suggestions to other models that respect those constraints).

Our past successes have almost always occurred when an academic collaborator has had a Google badge of some type – usually an intern, sometimes a visiting scientist. Having a collaborator temporarily "in house" side-steps many of the constraints described above.

- Interns and visiting scientists can be given access to anonymized data in some cases, after careful review, because as employees they have clearer legal obligations with respect to privacy and business confidentiality.
- When a collaboratively-written paper is primarily the product of Google employees, even temporary ones, we can deliver valuable insights without the complexities of releasing datasets to the public. (For confidentiality reasons, we cannot always publish the papers we might want to write.)
- When we do release datasets, having interns carry out much of the work can mitigate the staff-time constraint, while still providing a useful learning experience for the interns (and, perhaps, providing the interns with an opportunity to exploit those explicitly-released datasets in their university research, after finishing at Google).

We hire a lot of interns every year, but relatively few visiting scientists. However, because visiting scientists have more experience and often stay with us longer, they can create a much deeper collaboration. We are therefore especially interested in ideas for intern-driven and visiting-scientist collaborations around data-driven research.

We have some examples of a few successful data-driven collaborations, where a visiting scientist or intern was instrumental in either obtaining the data or preparing it for release:

- Borg traces: https://github.com/google/cluster-data provides access to workload traces from the scheduling system at the heart of Google's cluster management software and systems. Each trace set required the full-summer effort of an intern and host to collect and publish, and requires some ongoing effort to answer questions from users. The traces have been cited several hundred times in external papers. They also caused us to write "Obfuscatory obscanturism: making workload traces of commercially-sensitive systems safe to release" [6].
- "Evolve or Die: High-Availability Design Principles Drawn from Google's Network Infrastructure" [3] was written primarily by a visiting scientist, and draws on internal reports on high-impact failures. (Here, the data is not the traditional highly-structured trace data that drives a lot of networking research, but unstructured text that cannot be "anonymized" in any useful way.)

Several other recent papers were co-authored by Googlers and interns and/or visiting scientists, using some data from Google's networks (but not primarily "measurement" papers):

- "Carousel: Scalable Traffic Shaping at End-Hosts;; [7].
- "Sundial: Fault-tolerant Clock Synchronization for Datacenters" [4].
- "Network Error Logging: Client-side measurement of end-to-end web service reliability" [2].

## 4 SUGGESTIONS FOR POTENTIAL COLLABORATIONS

We encourage academic researchers to focus less on "can we obtain network-related data from Google?" and more on "how can we do more collaborative, data-driven networking research with Google?"

We close with some starting-point ideas for future collaborations – think of these as an informal "call for proposal" for research projects that each could involve a mixture of Google funding for university work, and some at-Google work carried out by visiting scientists and/or interns. Such projects would potentially lead to data releases within the context of a "sandboxed" collaboration, where the sandbox boundary is "people with a Google badge of some form, either permanent or temporary." More importantly, we believe that many of these, and ideas like them, could lead to peer-reviewed publications.

Please treat these as examples and suggestions to spur your own thinking.

- With the shift of massive compute resources to support machine learning, how does this affect conventional assumptions about network traffic patterns in datacenter networks? Similar questions might apply to large scale data-analytics.
- Does the research community rely on untested assumptions about where congestion occurs in datacenter networks and what causes it, both at the macro scale, and on packet-level timescales?

- We have seen a lot of interest in developing ML-inspired congestion control and adaptive-bitrate video encoding algorithms (e.g., Pensieve [5]), but also a lot of confusion on what actually works "in the wild" [8]). Can we develop a corpus of useful training and evaluation data, both to inspire a new generation of algorithms, and also to improve their validation over a sufficiently large range of use cases?

  Note that the existing M-Lab collaboration, https://www.measurementlab.net/, might be an existing vehicle for access to such data, with fewer restrictions; are there questions that the M-Lab approach does not support?

- For research into the public Internet, where M-Lab might be the right approach in general, but today is too limited, one could explore extensions of the M-Lab model, such as funding for placement of additional measurement nodes. The *Pantheon of Congestion Control* [9] was a community evaluation platform for academic research on congestion control, but it is no longer maintained; this illustrates the need for significant investments to sustain such platforms.

- While the research community has been focusing on formal methods for network verification, in the real world we often have to rapidly root-cause bugs without the benefit of formal approaches – for example, if a cloud customer reports a performance problem, is the root cause in the provider's network, and if so, where? This problem lies between totally manual approaches and formal ones, and steps towards more automation probably depend on a lot of data from various simultaneous sources.

- Can we spot otherwise undetected problems, before they are reported, via passive measurement at scale that leverage statistical inference, tomography, and other tools?

- Since trace-anonymization *per se* is often infeasible, under the constraint that it must preserve user privacy and proprietary information, if there were an open-source system that could properly mediate access via techniques such as differential privacy or running researcher-written analyses on-premises at Google, we might be able to create an intern-driven project to stand up such an infrastructure at Google. There are many "ifs" before we could approve that, and one interesting research question is whether it is even possible for a differential-privacy approach to comply with regulations such as GDPR, and under what constraints.

People often think of network data as structured data, such as a traffic matrix time series, a topology graph, or other traces. We encourage our collaborators, especially visitors, to broaden their view of "network data" to include architectural principles, and lessons learned from failures and near-misses, which can be shared in publications such as "Evolve or Die" [3].

There are many challenging research problems in networking that Google cannot solve without the help of the academic community, and many problems that the academic community cannot study without access to data from large operators. Data sharing is a challenge, but it is also an opportunity.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Mark Allman and Vern Paxson. 2007. Issues and Etiquette Concerning Use of Shared Measurement Data. In *Proc. Conference on Internet Measurement (IMC) (IMC '07)*. 135–140. https://doi.org/10.1145/1298306.1298327

[2] Sam Burnett, Lily Chen, Douglas A. Creager, Misha Efimov, Ilya Grigorik, Ben Jones, Harsha V. Madhyastha, Pavlos Papageorge, Brian Rogan, Charles Stahl, and Julia Tuttle. 2020. Network Error Logging: Client-side measurement of end-to-end web service reliability. In *Proc. NSDI*. 985–998. https://www.usenix.org/conference/nsdi20/presentation/burnett

[3] Ramesh Govindan, Ina Minei, Mahesh Kallahalla, Bikash Koley, and Amin Vahdat. 2016. Evolve or Die: High-Availability Design Principles Drawn from Google's Network Infrastructure. In *Proc. SIGCOMM*. http://dl.acm.org/authorize.cfm?key=N19254

[4] Yuliang Li, Gautam Kumar, Hema Hariharan, Hassan Wassel, Peter H. Hochschild, Dave Platt, Simon Sabato, Minlan Yu, Nandita Dukkipati, Prashant Chandra, and Amin Vahdat. 2020. Sundial: Fault-tolerant Clock Synchronization for Data-centers. In *Proc. OSDI*. 1171–1186. https://www.usenix.org/conference/osdi20/presentation/li-yuliang

[5] Hongzi Mao, Ravi Netravali, and Mohammad Alizadeh. 2017. Neural Adaptive Video Streaming with Pensieve. In *Proc. SIGCOMM*. 197–210. https://doi.org/10.1145/3098822.3098843

[6] Charles Reiss, John Wilkes, and Joseph L. Hellerstein. 2012. Obfuscatory obscanturism: making workload traces of commercially-sensitive systems safe to release. In *CloudMAN*. Maui, HI, USA. http://www.e-wilkes.com/john/papers/2012.04-obfuscation-paper.pdf

[7] Ahmed Saeed, Nandita Dukkipati, Valas Valancius, Terry Lam, Carlo Contavalli, and Amin Vahdat. 2017. Carousel: Scalable Traffic Shaping at End-Hosts. In *ACM SIGCOMM 2017*. https://research.google/pubs/pub46460/

[8] Francis Y. Yan, Hudson Ayers, Chenzhi Zhu, Sadjad Fouladi, James Hong, Keyi Zhang, Philip Levis, and Keith Winstein. 2020. Learning in situ: a randomized experiment in video streaming. In *Proc. NSDI*. 495–511. https://www.usenix.org/conference/nsdi20/presentation/yan

[9] Francis Y. Yan, Jestin Ma, Greg D. Hill, Deepti Raghavan, Riad S. Wahby, Philip Levis, and Keith Winstein. 2018. Pantheon: The Training Ground for Internet Congestion-Control Research. In *Proc. USENIX Annual Technical Conference (USENIX ATC '18)*. 731–743.