
What is Normal?

A Big Data Observational Science Model of Anonymized Internet Traffic

Jeremy Kepner

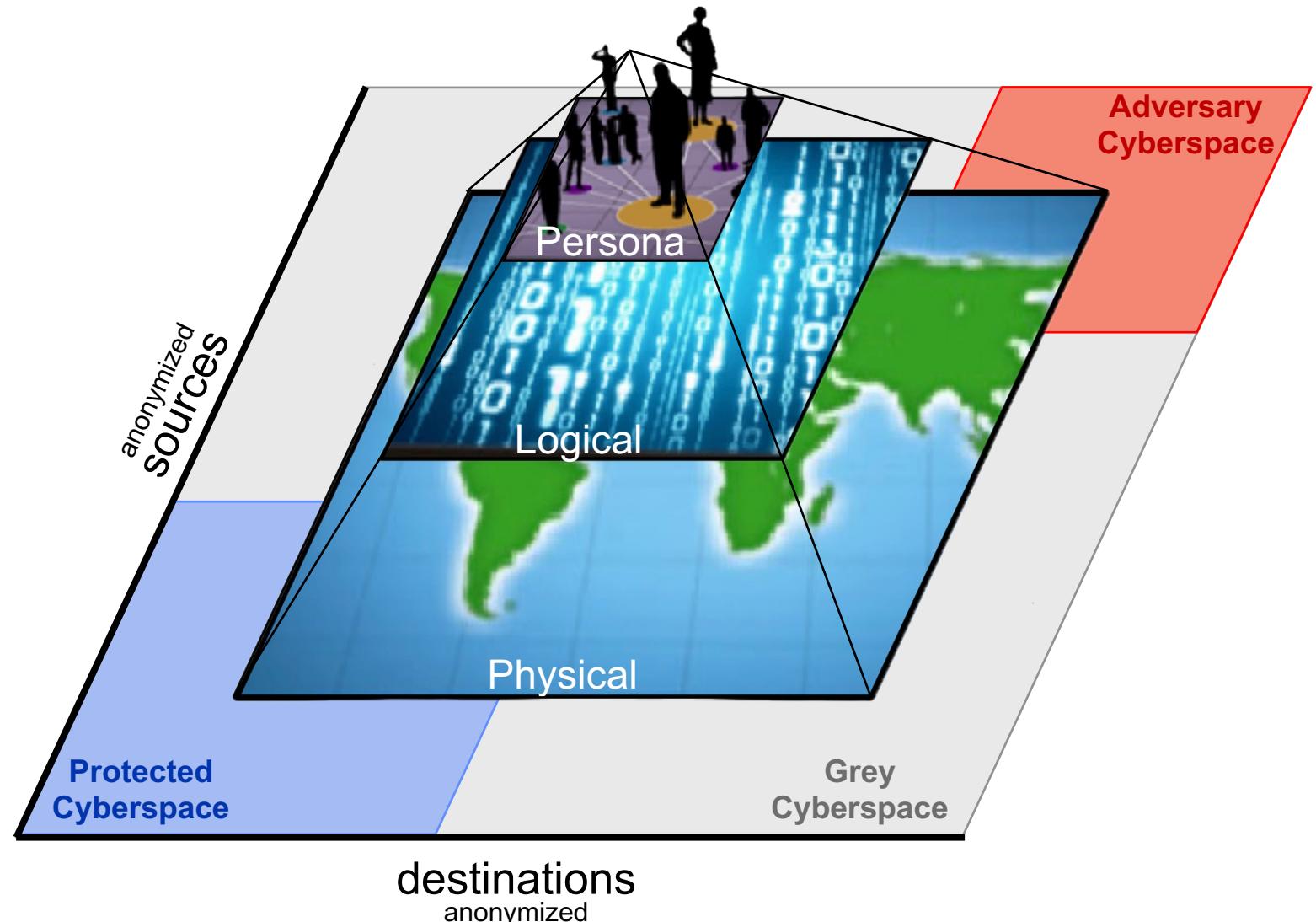
November 2025



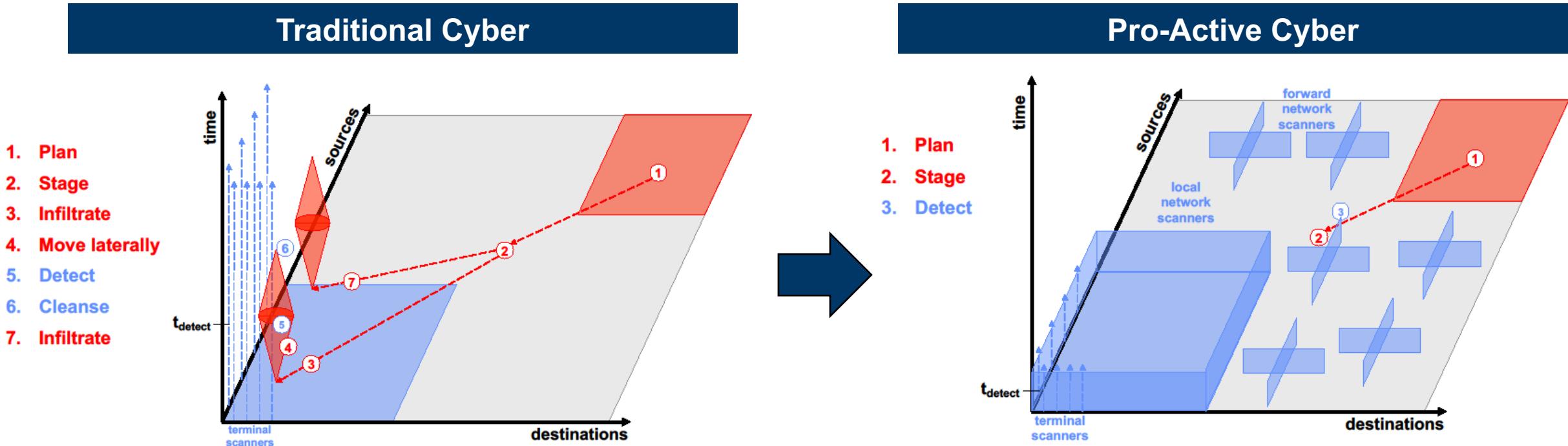
Research was sponsored by the Department of the Air Force Artificial Intelligence Accelerator and was accomplished under Cooperative Agreement Number FA8750-19-2-1000. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Department of the Air Force or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

Problem We Are Trying To Solve

- Problem
 - Need better sensors and AI algorithms to improve networks
- Approach
 - New cyber sensors
 - Novel privacy-preserving hierarchical AI techniques
 - Assembled the largest corpus of publicly available network data to support this research (100+ trillion events)



Expected Impact: Pro-Active Cyber



- Move sensors to outside of protected enclave
- Provide lead time to cyber defenders



Traffic Matrix Tutorial

Which attack stage is most relevant to the observed traffic? Hint:
youtu.be/_HDplx2MxyA or

stage

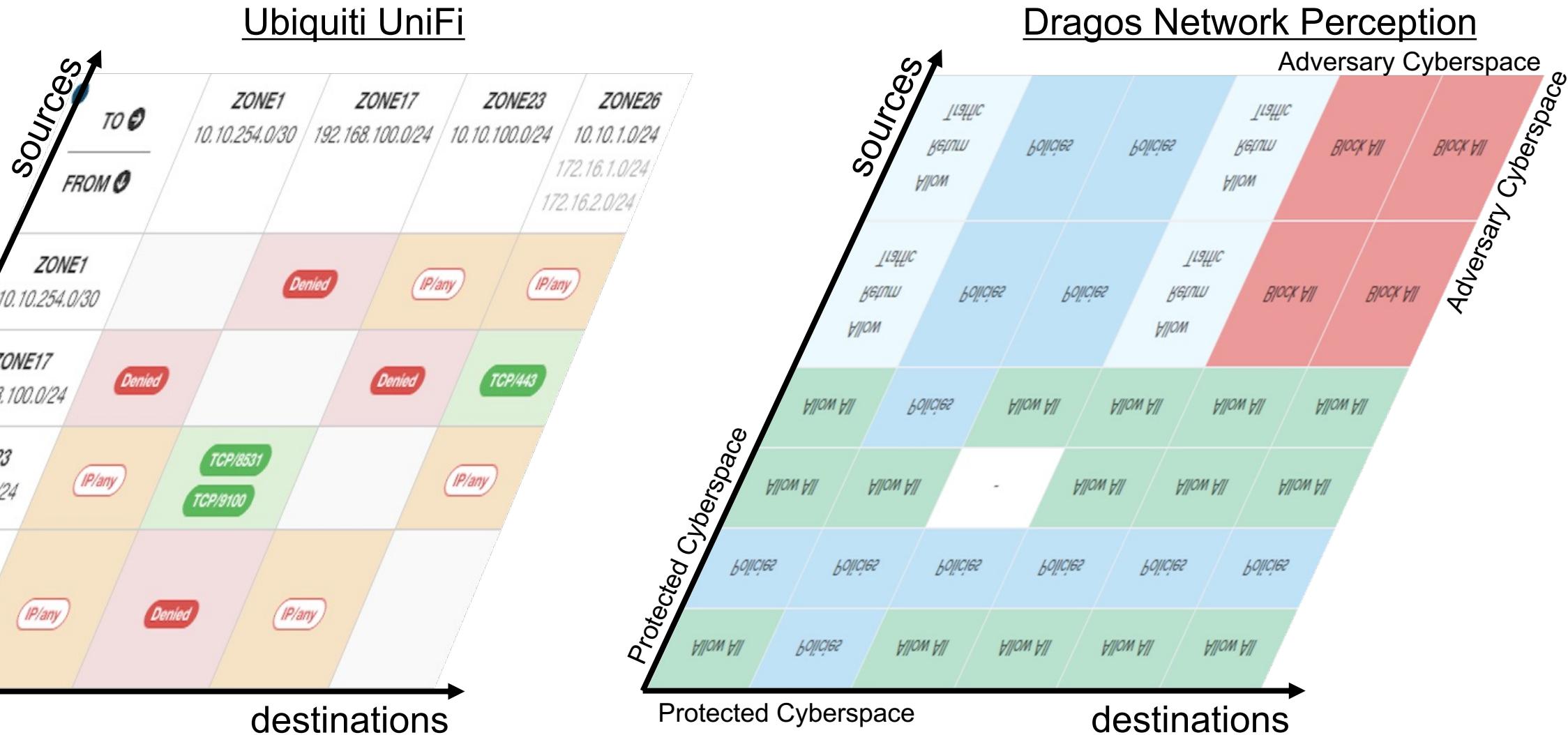
infiltrate

move laterally

Spacebar - Switch View
Q/E - Rotate screen



Industry Firewall Zone Matrix Adoption





Anomaly Detection and Observational Science

- Core challenge for anomaly detection systems is adequate models of normal

The concept of normality. It is one of the main steps to build a solution to detect network anomalies. The question ``how to create a precise idea of normality?'' is what has driven most researchers into creating different solutions through the years. This can be considered as the main challenge related to anomaly detection and has not been entirely solved yet.¹

- Reproducible observations of cyberspace have been recommended as a core foundation for the science of cyber-security

The highest priority should be assigned to establishing research protocols to enable reproducible [observations]²

Other domains (land, sea, undersea, air, and space) rely on detailed observational science models of their environment to understand what is normal

- Network operators/owners collect, analyze, and share data only to improve network operations/security
 - Data recipients inherit that intent
- CAIDA & MIT have pioneered new approaches that approve appropriately anonymized data with clear who/what/where/why/how agreements
 - MIT now has the largest collections of data because of our unique ability to analyze these data at scale while protecting privacy



CAIDA Master Acceptable Use Agreement (AUA)

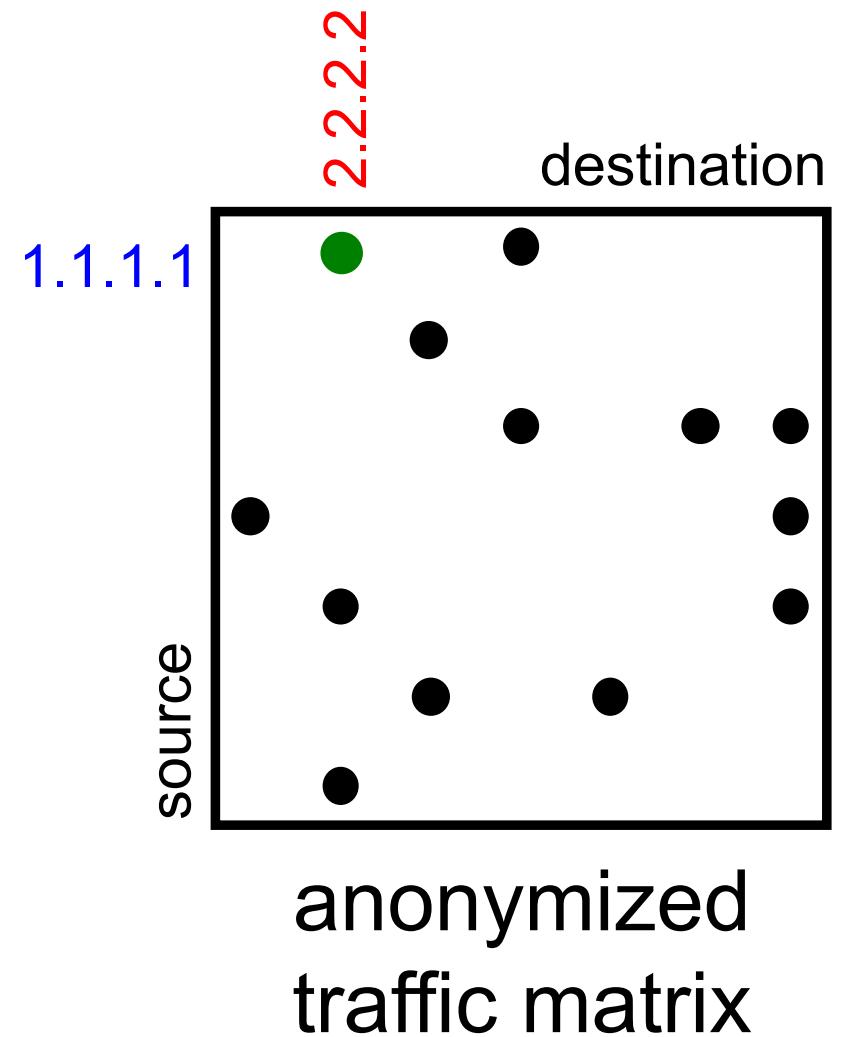
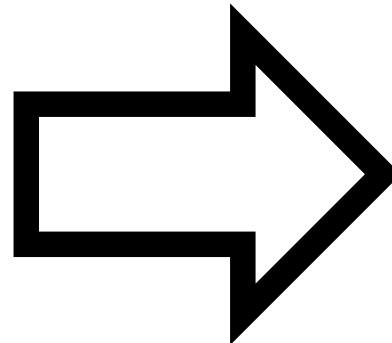
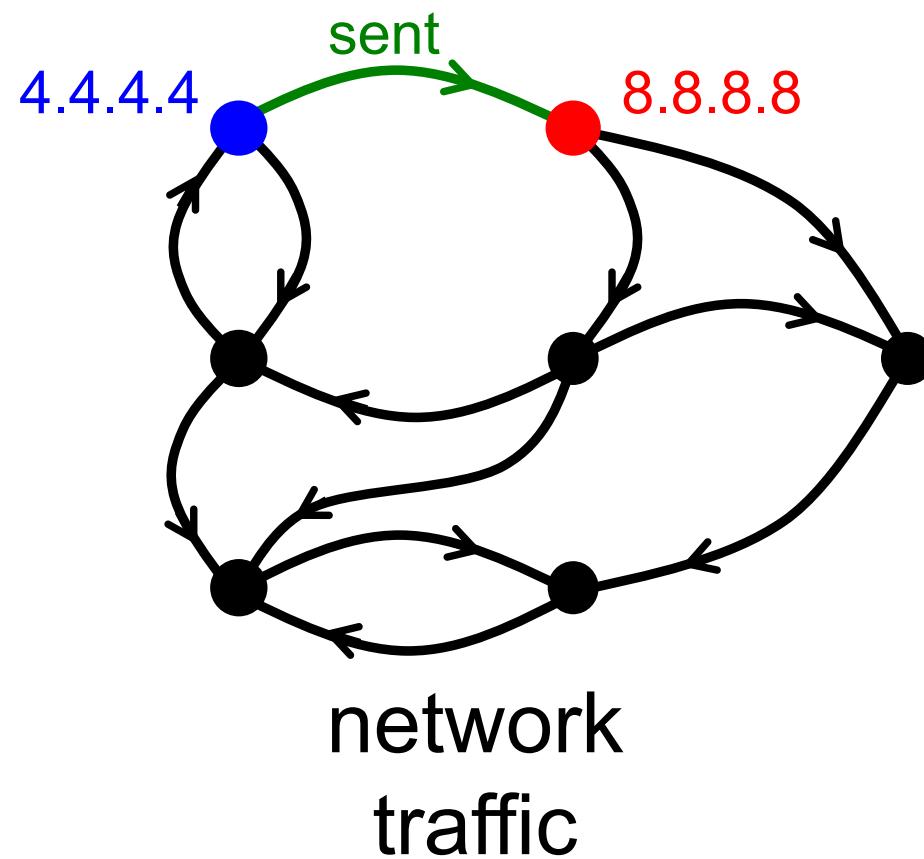
The CAIDA Acceptable Use Agreement (AUA) applies to the majority of the CAIDA datasets, with the exception of datasets covered by the [Acceptable Use Agreement for Publicly Accessible Datasets \(Public-AUA\)](#). The relevant dataset agreement is shown as part of the data request process. See available CAIDA Datasets.



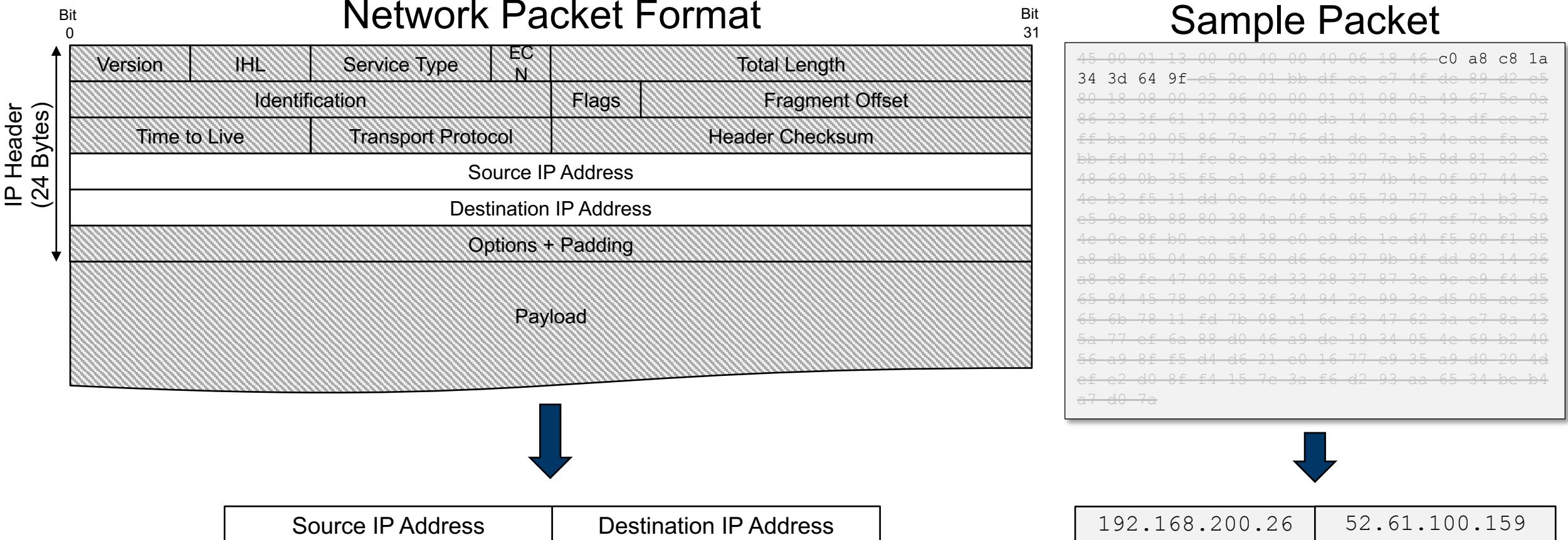
IS&T Data Request Form

Overview:

Please fill out the form below to make a request for enterprise administrative data managed by IS&T. The time to receive data depends on the complexity of the request and availability of our staff. Please give us as much advanced notice as possible if you are working with a deadline.



Keep Only Source/Destination IPs



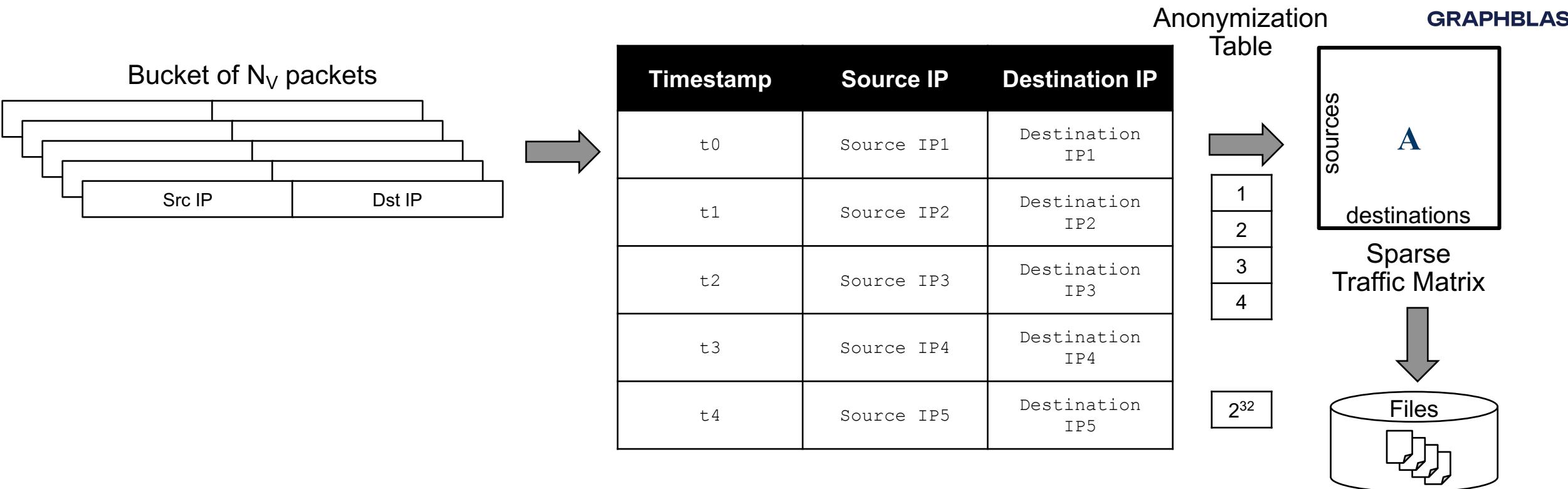
ECN: Explicit Congestion Notification

IHL: Internet Header Length

Construct Anonymized Traffic Matrix

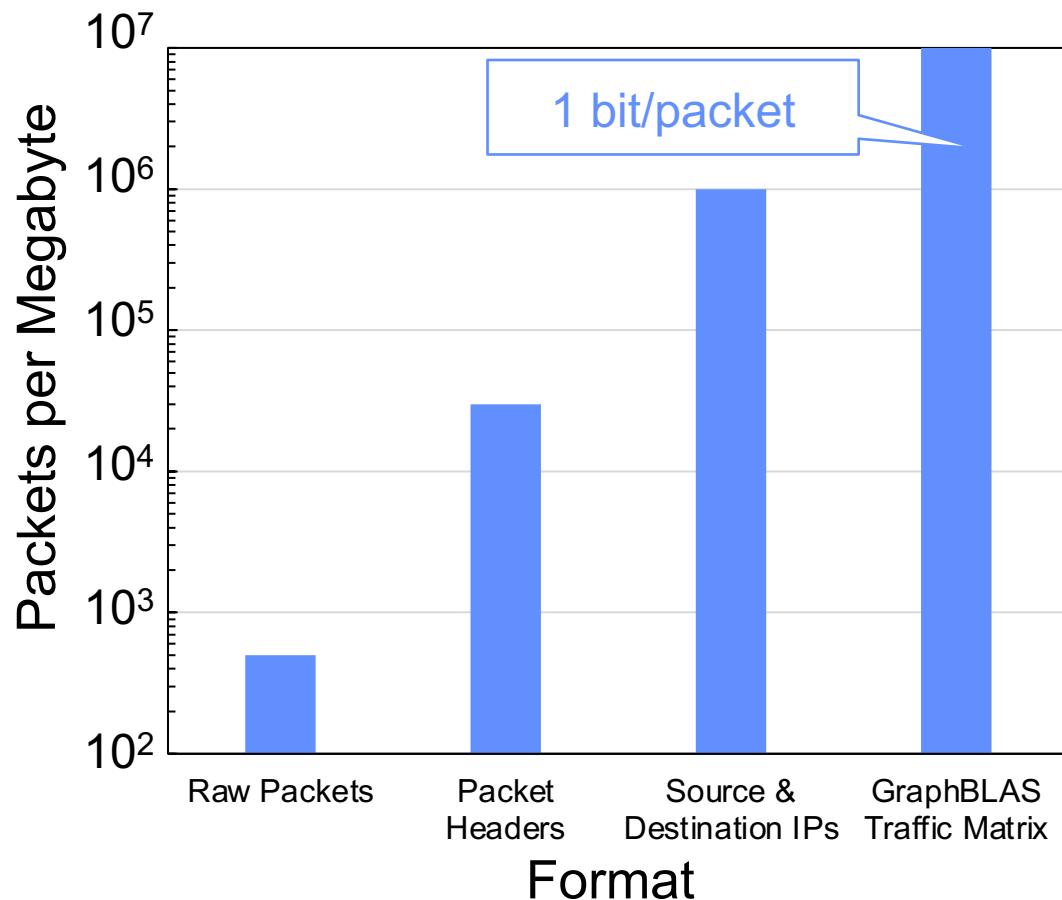


GRAPHBLAS



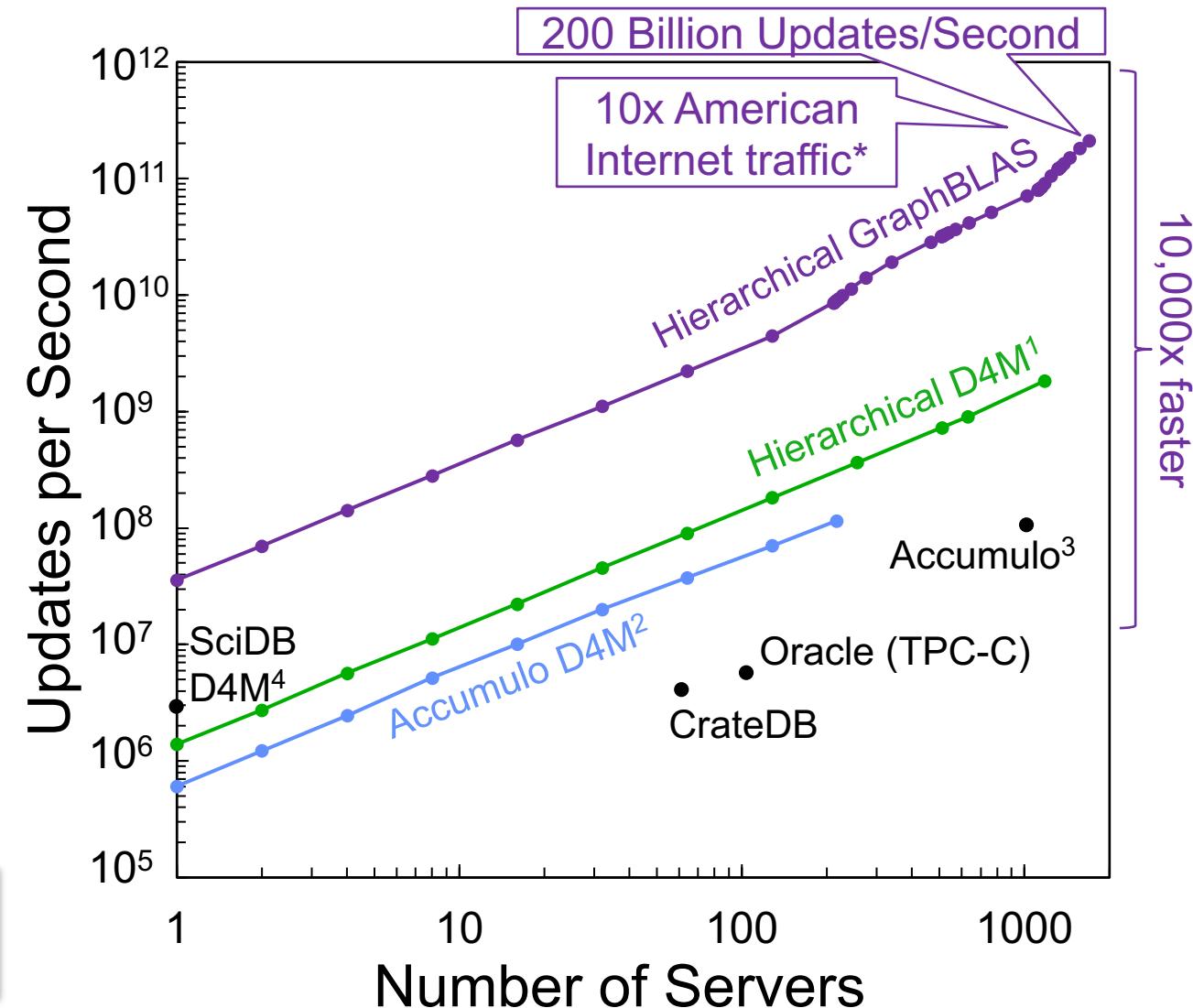
Performance Breakthroughs

Compression

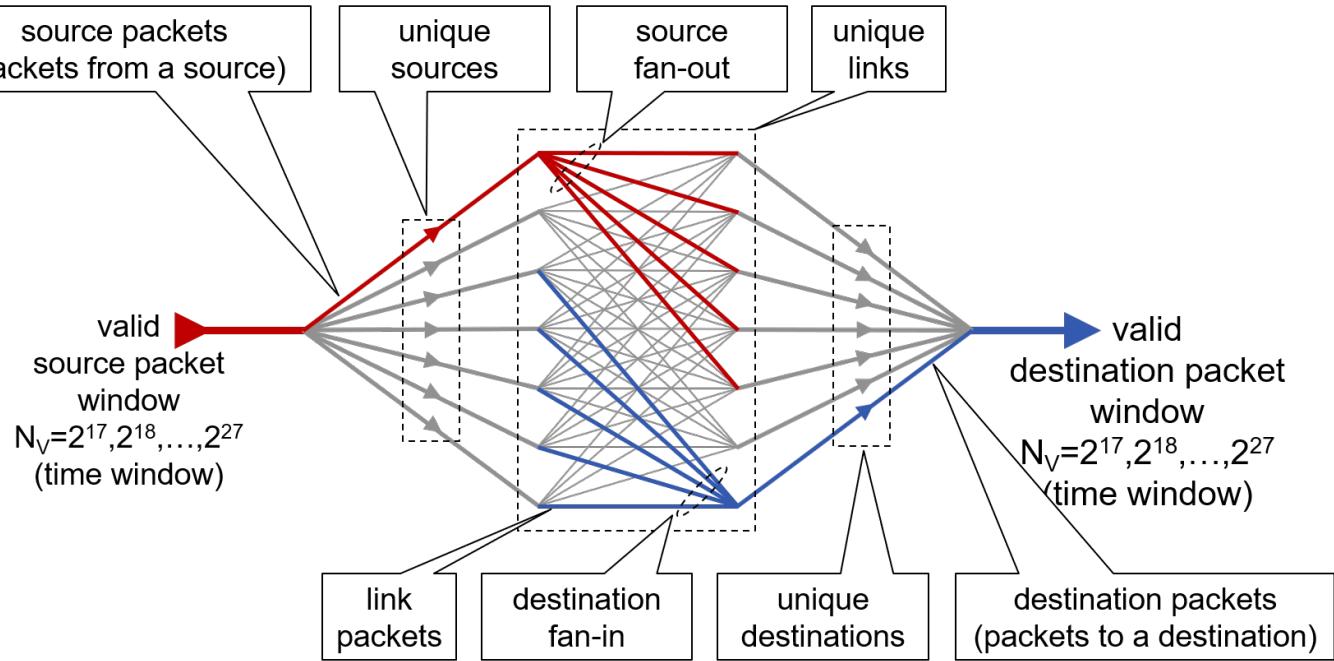


GraphBLAS Traffic Matrices
3,000x compression & 10,000x faster

Speed



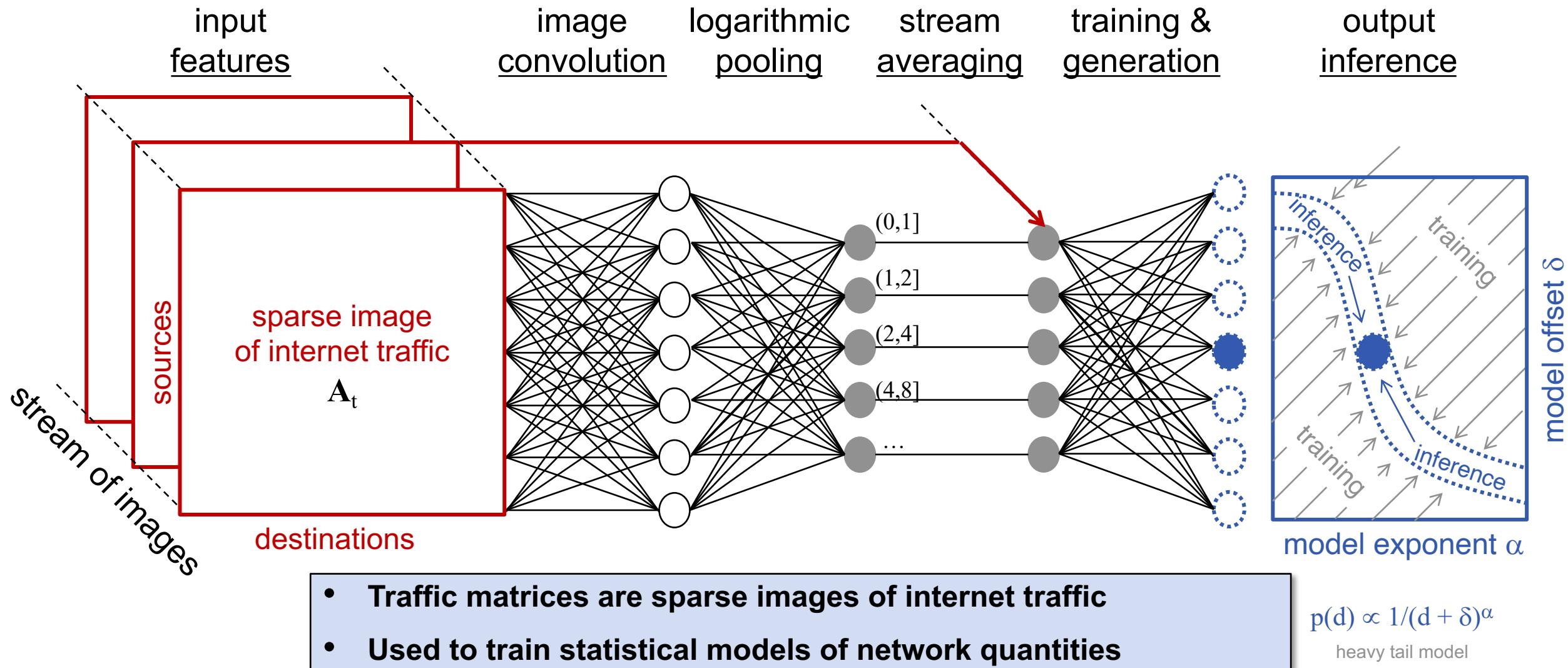
GraphBLAS Hypersparse Math



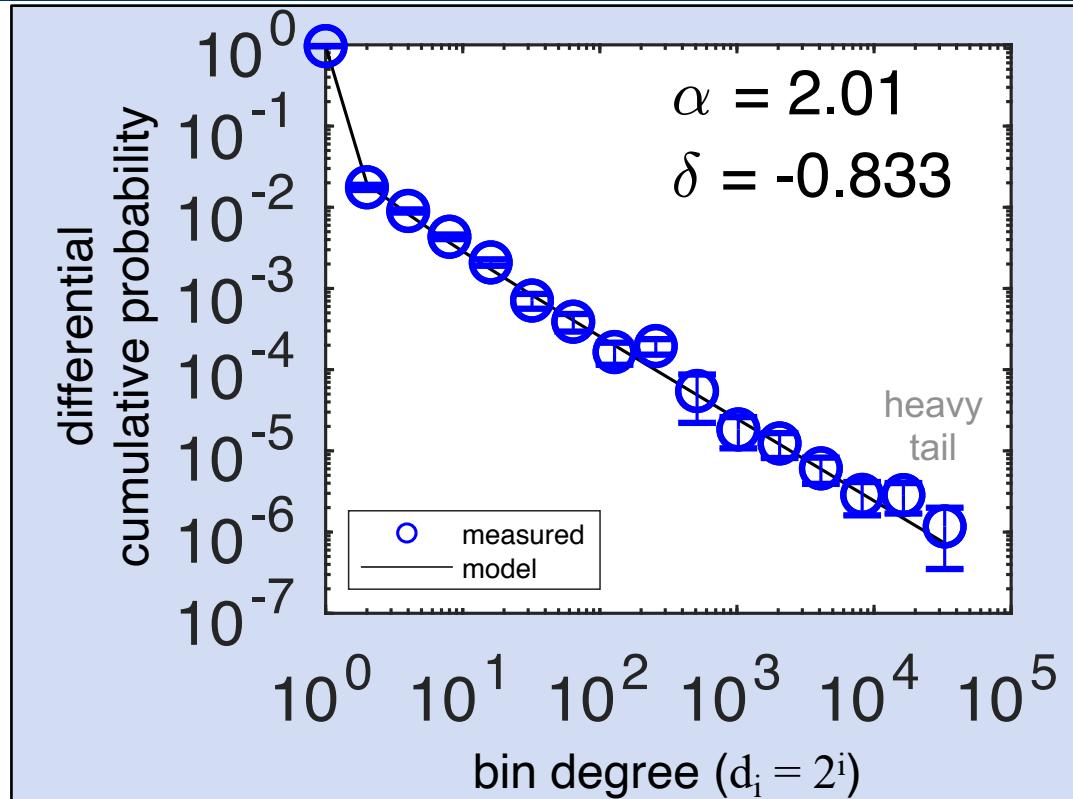
Aggregate Property	Summation Notation	Matrix Notation
Valid packets N_V	$\sum_i \sum_j \mathbf{A}_t(i, j)$	$\mathbf{1}^\top \mathbf{A}_t \mathbf{1}$
Unique links	$\sum_i \sum_j \mathbf{A}_t(i, j) _0$	$\mathbf{1}^\top \mathbf{A}_t _0 \mathbf{1}$
Link packets from i to j	$\mathbf{A}_t(i, j)$	\mathbf{A}_t
Max link packets (d_{\max})	$\max_{ij} \mathbf{A}_t(i, j)$	$\max(\mathbf{A}_t)$
Unique sources	$\sum_i \sum_j \mathbf{A}_t(i, j) _0$	$\mathbf{1}^\top \mathbf{A}_t \mathbf{1} _0$
Packets from source i	$\sum_j \mathbf{A}_t(i, j)$	$\mathbf{A}_t \mathbf{1}$
Max source packets (d_{\max})	$\max_i \sum_j \mathbf{A}_t(i, j)$	$\max(\mathbf{A}_t \mathbf{1})$
Source fan-out from i	$\sum_j \mathbf{A}_t(i, j) _0$	$ \mathbf{A}_t _0 \mathbf{1}$
Max source fan-out (d_{\max})	$\max_i \sum_j \mathbf{A}_t(i, j) _0$	$\max(\mathbf{A}_t _0 \mathbf{1})$
Unique destinations	$\sum_j \sum_i \mathbf{A}_t(i, j) _0$	$ \mathbf{1}^\top \mathbf{A}_t _0 \mathbf{1}$
Destination packets to j	$\sum_i \mathbf{A}_t(i, j)$	$\mathbf{1}^\top \mathbf{A}_t _0$
Max destination packets (d_{\max})	$\max_j \sum_i \mathbf{A}_t(i, j)$	$\max(\mathbf{1}^\top \mathbf{A}_t _0)$
Destination fan-in to j	$\sum_i \mathbf{A}_t(i, j) _0$	$\mathbf{1}^\top \mathbf{A}_t$
Max destination fan-in (d_{\max})	$\max_j \sum_i \mathbf{A}_t(i, j) _0$	$\max(\mathbf{1}^\top \mathbf{A}_t)$

GraphBLAS hypersparse traffic images enable efficient computation of network quantities in C/C++/Python/Julia/Matlab/Octave

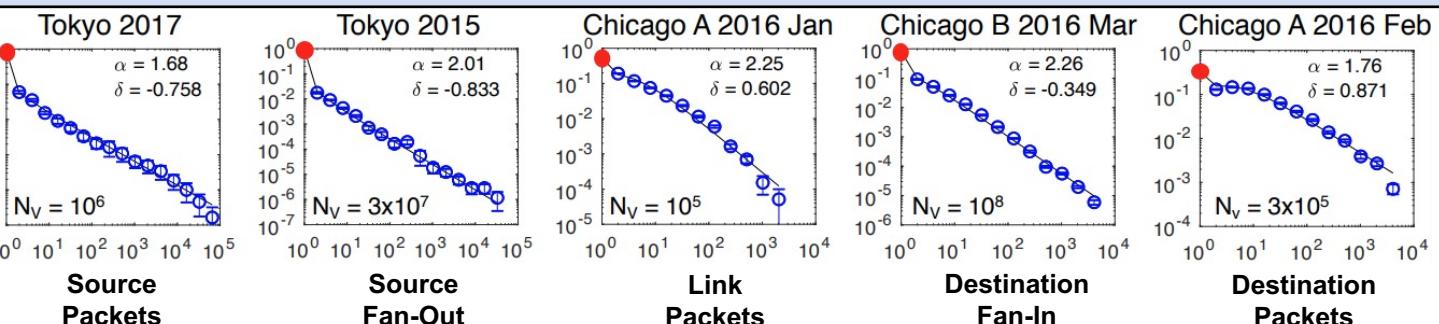
Graph Neural Network Analysis of Large-Scale Internet Traffic



High Precision Heavy Tail Internet Background Model



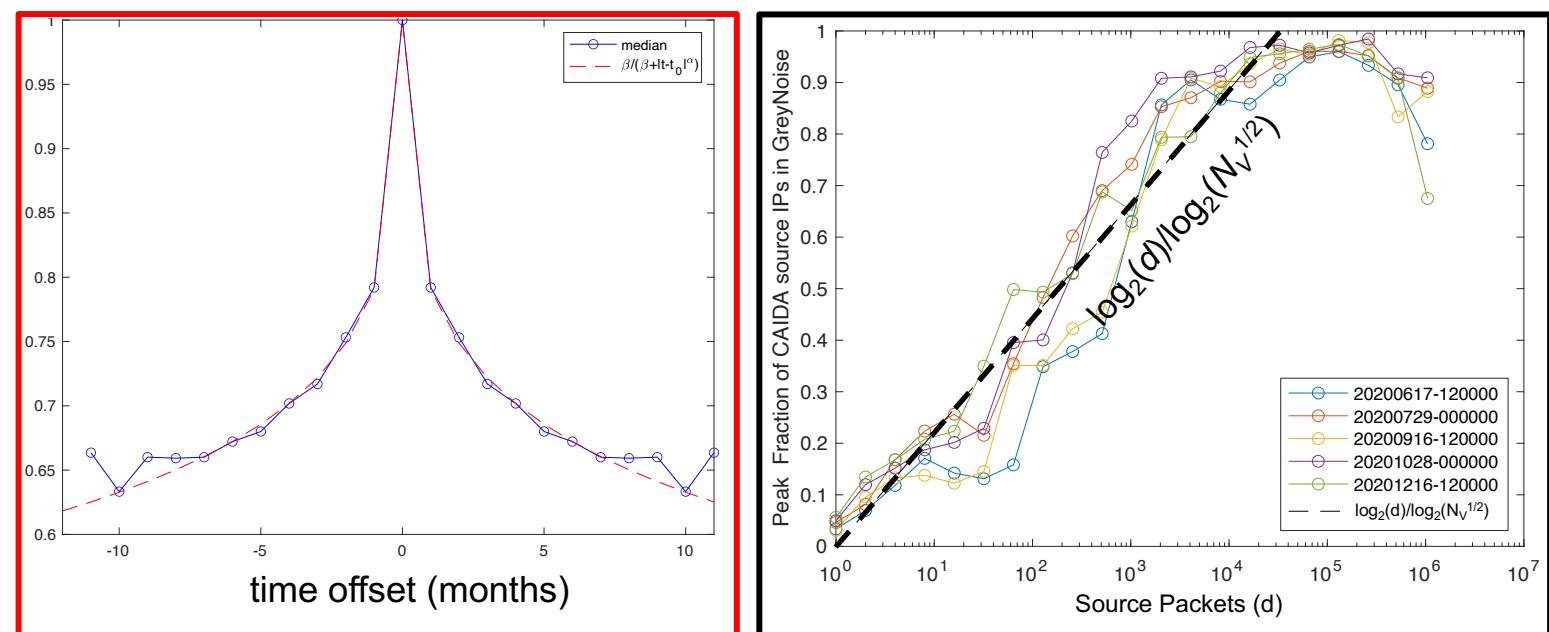
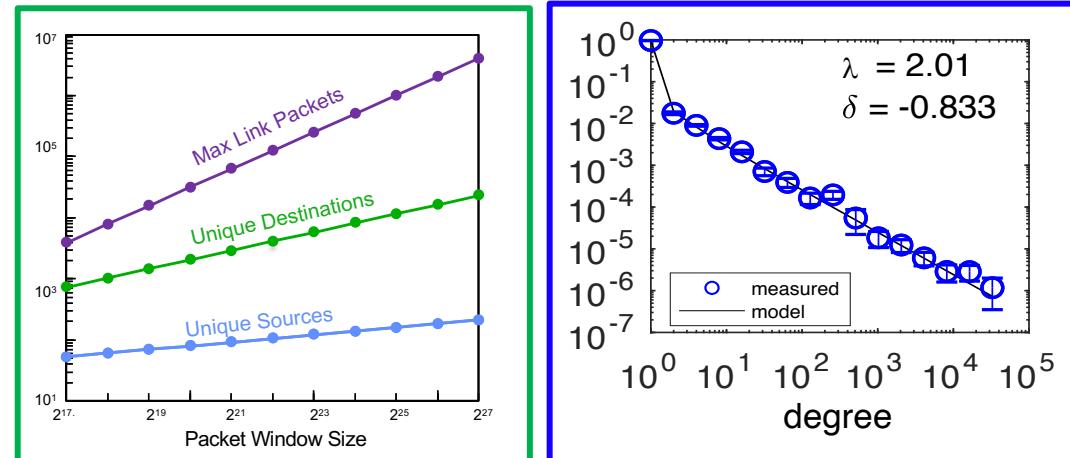
- Characterizes each distribution with two parameters



Holistic Internet Background Model

- Given a window with N_V incoming packets, the probability of seeing a source emitting d packets being observed again at time t by a different observer is proportional to

$$\begin{aligned}
 & N_V^a \\
 & 1/(d + \delta)^\lambda \\
 & \beta/(\beta + t^\alpha) \\
 & \log_2(d)/\log_2(N_V^{1/2})
 \end{aligned}$$



- where $a, \lambda, \delta, \alpha$, and β , are site specific model parameters that vary slowly in time

AI Ready Anonymized Network Sensing Challenge

- IEEE High Performance Extreme Computing (HPEC) Conference -

The screenshot shows the GraphChallenge website interface. At the top, there are logos for IEEE HPEC, MIT, Amazon Web Services, and NVIDIA. Below the header is a navigation bar with links: Motivation, Challenges, Data Sets, Scenarios, Submit, News, Champions, and Contact. The main content area has a "Home" link. The page is divided into four main sections:

- Specification:** Shows a network traffic graph with nodes 4.4.4.4 and 8.8.8.8, and an anonymized traffic matrix.
- Code:** Displays a GitHub-style code listing for `/src/pcap/pcap2grb.c` by msjoneso, updated 5 months ago. It includes 601 lines (522 loc) and 20.2 KB. Below the code is a snippet of C code:


```
1 #define _GNU_SOURCE 1
2
3 #include <GraphBLAS.h>
```
- RESOURCE CATALOG:** Features a "caida" logo and a "Anonymized Network Sensing Graph Challenge Dataset" entry. The dataset is identified by the identifier `dataset:2024_ieee_anonymized_nsg_challenge`. It includes links for overview, related, cite, access, and README.txt.
- Data:** Shows a large orange arrow pointing right, indicating the flow from specification to data.

New! Anonymized Network Sensing Graph Challenge: This challenge constructs and analyzes anonymized traffic matrices from network packet capture (PCAP) data to enable open community-based approaches to protecting networks.

- Specification: slides, [paper](#), [example serial parse code](#), [example serial analyze code](#), [example data sets](#)



OneSparse



RPI



Northeastern University

ILLINOIS INSTITUTE OF TECHNOLOGY



FABRIC