# Lawrence Berkeley National Laboratory

## LBL Publications

Title
The NetSage measurement and analysis framework in practice

Authors
Schopf, Jennifer M
Turner, Katrina
Doyle, Dan
et al.

# The NetSage measurement and analysis framework in practice

Jennifer M. Schopf[1] · Katrina Turner[2] · Dan Doyle[1] · Andrew Lake[3] · Jason Leigh[2] · Brian L. Tierney[1]

## Abstract

Data sharing is required for research collaborations, but effective data transfer performance continues to be difficult to achieve. The NetSage Measurement and Analysis Framework can assist in understanding research data movement. It collects a broad set of monitoring data and builds performance Dashboards to visualize the data. Each Dashboard is specifically designed to address a well-defined analysis need of the stakeholders. This paper describes the design methodology, the resulting architecture, the development approach and lessons learned, and a set of discoveries that NetSage Dashboards made possible.

## 1 Introduction

Scientific investigation is highly collaborative and requires the ability to seamlessly share data between institutions to enable scientific discovery. However, effective data sharing, especially for large data sets, can be challenging. For example, a common astronomy workflow involves a telescope producing data sets of 100 TBytes of data every day, which are then sent to multiple international sites for analysis. The data transfers and individual site processing

✉ Jennifer M. Schopf
  jmschopf@indiana.edu

  Katrina Turner
  khj@hawaii.edu

  Dan Doyle
  daldoyle@globalnoc.iu.edu

  Andrew Lake
  andy@es.net

  Jason Leigh
  leighj@hawaii.edu

[1]  Indiana University, Bloomington, Indiana, USA

[2]  University of Hawaii, Honolulu, Hawaii, USA

[3]  Lawrence Berkeley National Laboratory, Berkeley, California, USA

must complete before the next data collection window in order to re-calibrate the telescope. If the data is not received in a timely fashion, the telescope cannot get re-focused and a viewing window may be lost. In another example, it took over three months to transfer data from a set of climate science experiments to a central location for analysis [16]. It is not uncommon for researchers to resort to shipping disks instead of using the network for data delivery due to network performance issues, such as is done for the Event Horizon Telescope black hole analysis [19]. Its also not unusual for a research collaboration to poorly estimate the time it may take to transfer a data file, such as when the Arecibo Telescope data center planned to back their data up in the cloud, and assumed it would take weeks and not the multiple years required by their initial approach.

The ability to measure and interpret network behavior is critical to understanding data transfer performance. Information about the end-to-end data path makes it possible to identify and address problems that are found. Our experience has shown that without this data, researchers often misunderstand the true performance of their data transfers, which can significantly impact their workflows and time to science.

This paper details the NetSage Measurement and Analysis Framework [56, 81], which is used to understand data transfer performance. We describe the stakeholder-focused methodology used to design performance Dashboards, each of which is focused on specific user questions.

We detail the software architecture and its implementation that uses multiple data sources and takes advantage of related approaches. We then walk through several use cases to show the types of analyses and discoveries that the NetSage Framework enables.

This paper describes both the international and US domestic deployments of NetSage Portals as of April 2021. Previous publications about the NetSage Framework, for example [93], only addressed the International Portal and a minimal set of components and use cases. Other publications, such as [29, 30], focused on how the project worked with large-scale data and some very preliminary deployments prior to wide-scale collection of flow data. Presentations by the NetSage Project at numerous venues are available on the NetSage Project Website [56].

## 2 NetSage overview

The NetSage Framework was created around a set of three design principles:

– NetSage *unifies* a broad set of data into a cohesive visualization, thereby enabling additional discoveries not possible with only a single data source.
– NetSage is *open* in that the data sets collected are meant to be widely accessible, with performance Dashboards open to the public. No other monitoring approach shares as much data in a public setting. None of the data shared is considered private or public, but making it widely available increases the usability and the audience. The NetSage is also available under open source license and available via GitHub [52].
– NetSage is *privacy-aware* and was developed with privacy concerns in mind. It contains no personally identifiable information (PII) about individual hosts or users of the network. Information about small flows, such as common with email or web searches, is also discarded. Also, if required by the data provider, the Dashboards can be secured by password or Shibboleth [42]. *The Dashboards described in this paper are all open to the public. We encourage readers to visit the URLs cited to see the full Dashboards in detail.*

The innovative aspect of NetSage is not in the individual pieces but rather in the integration of data sources to support objective performance observations as a whole. We refer to the *NetSage Framework* as an inclusive term for the methodology, architecture, and approach. *NetSage deployments* can collect data from routers, switches, active testing sites, and science data archives. A NetSage deployment uses a combination of passive and active measurements to provide longitudinal performance visualizations via *performance Dashboards*. We refer to a set

of performance Dashboards for a project or set of resources as the *NetSage Portal* for that project or set of resources. The Dashboards can be viewed by resource collection, institutions or projects to identify changes of behaviors for data transfers using visualizations of data over time periods, as described in Sect. 6.

The NetSage Framework was developed for use by a broad user community, not by a single Network Operations Center (NOC). It was developed not only for network engineers, but resource owners, research technology departments, application scientists, CIOs, and funding agencies. The design approach was to meet a set of end user questions through innovative Dashboards, not just to supply measurement data to a NOC, which is the primary use case for most of the projects cited in the Related Work section.

The main use cases for the NetSage Framework have included:

– Understand the data movement patterns across a suite of resources;
– Identify the main sources and destinations for large data transfers, or flows;
– Visualize information about different research projects and science domains that are moving data;
– Display patterns of behaviors for data movement between organizations.

NetSage deployments now encompasses 14 regional deployments, listed in Table 1 and shown in Fig. 1. The NetSage Project was originally funded as part of the National Science Foundation (NSF) International Research and education Network Connections (IRNC) [37] program to develop and deploy advanced measurement services to measure how the science and engineering community was taking advantage of NSF-funded research networks and exchange points. The original deployment included working with, and gathering data from, the seven funded IRNC projects: Atlantic Wave Software Defined Exchange [32], America's Lightpath Express and Protect (AmLight ExP) [33], Networks for European, American, and African Research (NEAAR) [82], Pacific Islands Research and Education Networks (PIREN) [39], Pacific Wave [24], StarLight [43], and TransPAC4 [80], in addition to a set of science data archives. The international projects share the NetSage International Portal Bandwidth Dashboard [53], as shown in Fig. 2. Other individual resource sets, such as the Great Plains Network (GPN), iLight (the Indiana state network), KINBER, and so on, have a Portal specific to their data. In addition, the NetSage *All Data Portal* [57] includes the complete set of data collected for all of the other Portals, which is useful to analyze traffic more broadly, for example, by a science project that spans many networks globally. Table 1 lists the different data sources

**Table 1** The list of the US domestic and international deployments for NetSage as of March 2020 and the number and type data sources collected from each for use in NetSage deployments

| Project/resource set | SNMP | Flow | Science Archive |
|---|---|---|---|
| GPN | 2 | 2 | – |
| iLight | – | 5 | – |
| KINBER | – | 2 | – |
| FRGP/NCAR | – | 1 | 1 |
| SoX | – | 3 | – |
| TACC | – | 1 | 4 |
| PIREN, UH Astro | 5 | 2 | 1 |
| NERSC | – | – | 10 |
| ANA | 7 | – | – |
| Pac Wave/ CENIC | 17 | 7 | – |
| NEAAR/ NEA3R | 3 | 3 | – |
| TP4/ TP5 | 2 | 3 | – |
| AmPath/ AmLight | 11 | 6 | – |
| StarLight | 2 | – | – |

collected for each deployment. Note that perfSONAR test points are only available for the international resources and are not listed separately.

The data sources for 2020 included over 2.5 Billion flow data records collected from almost 60 sources, as well as SNMP data from over 60 unique nodes across the US and internationally. During 2020, over 3,600 unique users in 101 countries visited NetSage Dashboards.

## 3 Design methodology

The NetSage team adapted the Immersive Empathic Design Methodology (IEDM) [11] for developing visualizations. Visualization experts have commonly used this process, or similar techniques such as Design Thinking [9], to successfully produce effective visualizations for many decades. This methodology has eight steps:

1. Create profiles for representative stakeholders to understand their visualization needs.
2. Sketch storyboards to characterize the type of visualization to answer the stakeholder identified needs.
3. Present storyboards to stakeholders for feedback, which is often accomplished by recording storyboard presentations for stakeholders to view and comment on.
4. Update the storyboards based on the feedback from Step 3, and reiterate, as time and resources allow.
5. Develop prototypes based on the storyboards.
6. Give early working prototypes to stakeholders for them to try out in their own workflows.

7. Elicit feedback from the stakeholders.
8. Iterative development using the feedback to produce a successively better system, as well as to introduce additional requested features.

For example, for the initial NSF International Networks NetSage deployment, the profiles for representative stakeholders were defined by identifying the set of end users for the NetSage Dashboards and the types of questions they might ask of the data. The initial Dashboard users included:

– Network resource owners and operators who wanted to know the status of the resources;
– Collaborative research teams trying to understand resource use and how their data transfers would behave;
– Engineering staff to ensure effective resource use; and
– Staff members at funding agencies who needed additional insight into their investments.

After discussions with representatives from each audience, sets of use case questions were identified, including:

– What is the present state of the NSF-funded international network resources?
– What are the top sources or destinations for data flows using the NSF-funded international network resources?
– What are the top science domains that use the NSF-funded international network resources?
– What is the maximum, minimum, and average duration of large data transfers?
– Which countries are sharing data using the NSF-funded international network resources?
– Are there patterns of behaviors that can be identified about how the NSF-funded international network resources are used?
– Which sources or destinations have transfers that are not effectively using the NSF-funded international network resources?

These questions were used by the NetSage development team to design Dashboards with visualizations to provide the answers. A series of hand-drawn graphical storyboards were produced to describe the proposed Dashboards. The feedback during Step 3 enabled the NetSage development team to identify commonalities across the stakeholders and to adapt the Dashboard designs accordingly. This approach not only verified that user goals were being addressed, but also that each Dashboard was focused on addressing the response to a particular question.

As the use of the NetSage Framework expanded beyond the NSF-funded international network resources, we continued to deploy this design methodology and to incorporate feedback accordingly, leading to additional features for all end user stakeholders.

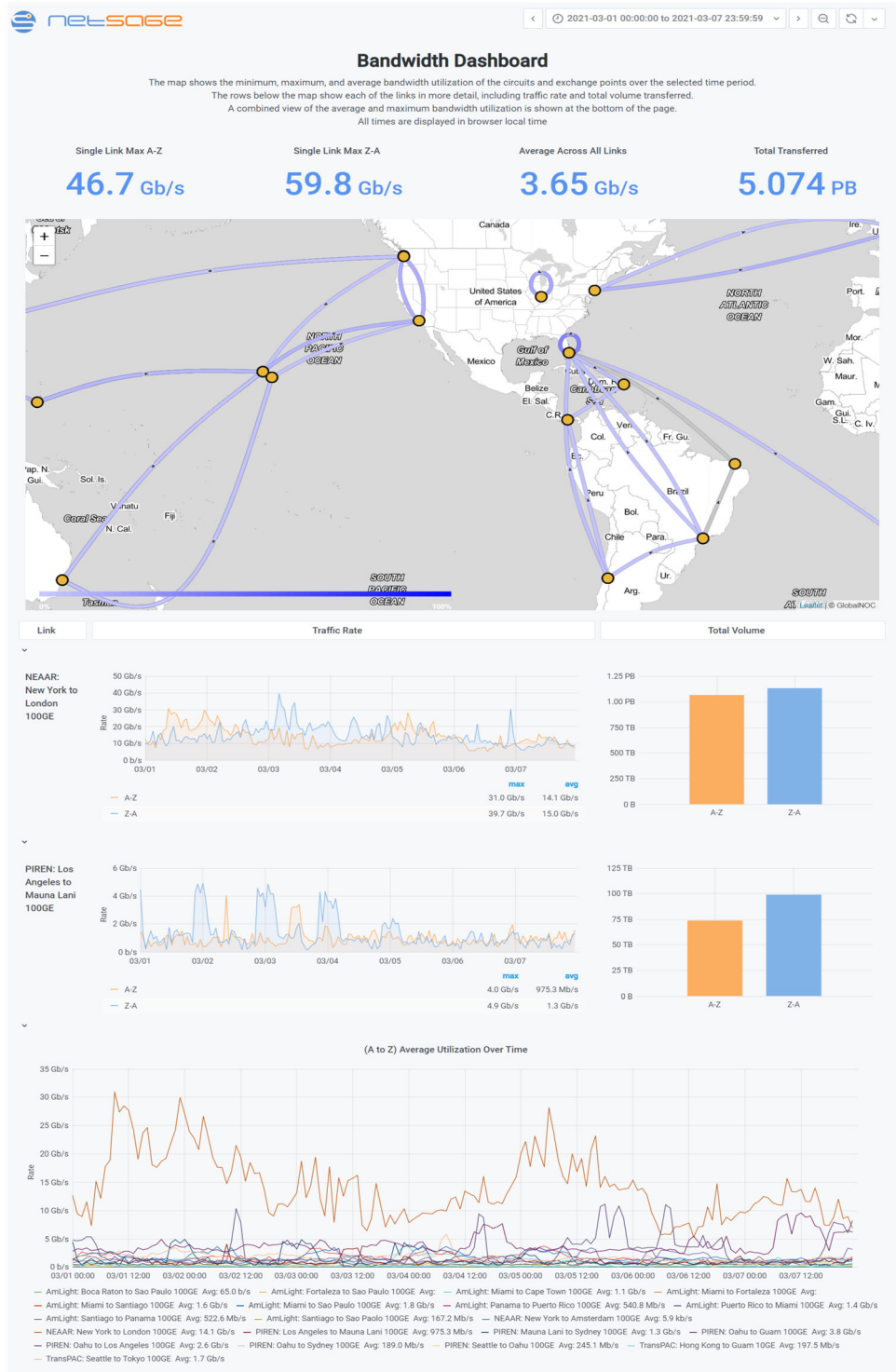**Fig. 1** A map showing the US domestic and international deployments for NetSage as of March 2021
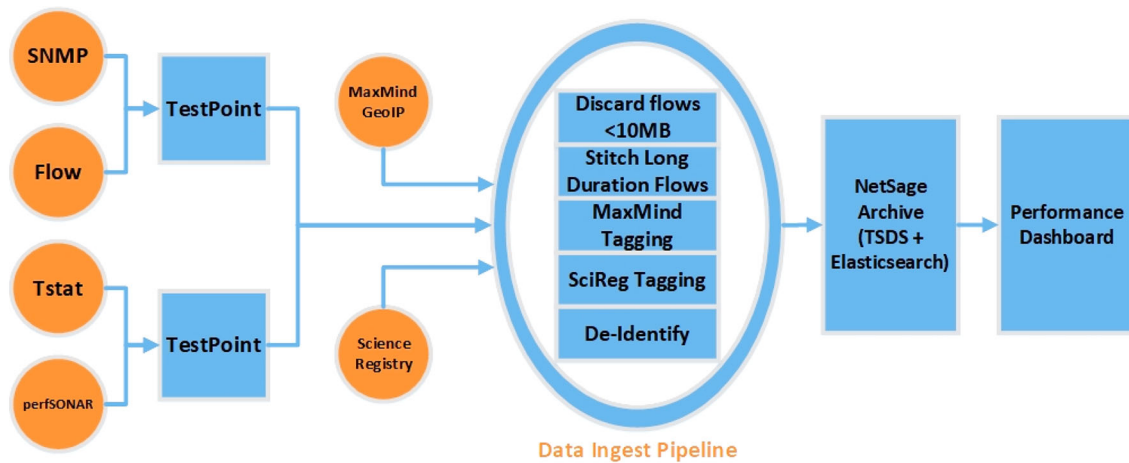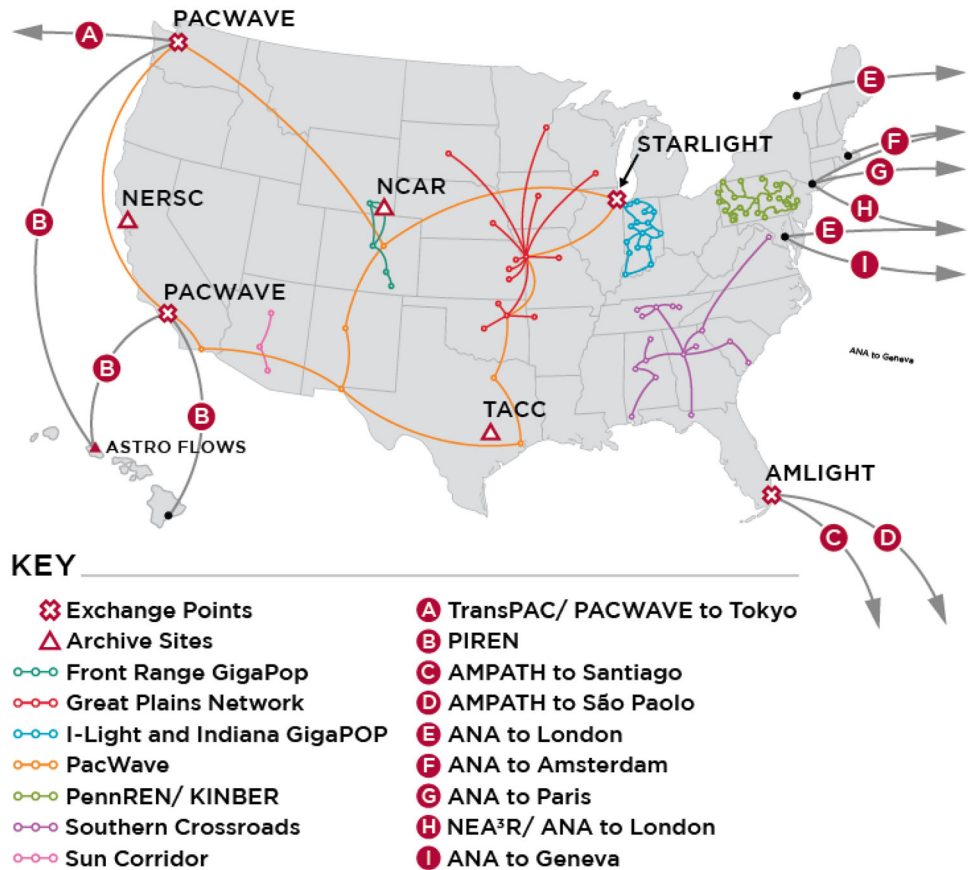
**Fig. 1** continued



**Fig. 2** A logical diagram of the NetSage Architecture, consisting of Data sources, the Data Ingest Pipeline, the NetSage Archive, and sets of performance Dashboards

## 4 NetSage architecture

The NetSage Software consists of a set of open source tools that follows a basic monitoring tool architecture, as shown in Fig. 2. *NetSage TestPoints* are a collection of software and hardware components that gather active and passive data into records that are sent to the *Data Ingest Pipeline*. The five-stage Pipeline filters those records and adds additional tags before de-identifying the data. The records are then stored in the *NetSage Archive*, a centralized storage framework consisting of two different databases, a Time Series Data System (TSDS) archive [89] and an

Elasticsearch archive [20]. *Performance Dashboards*, built using the open source Grafana [31] analysis and visualization engine, access the records from the NetSage Archive and present visualizations to answer the questions identified by the stakeholders.

## 4.1 NetSage data collection TestPoint

The core of the data collection process for a NetSage deployment is the set of hardware and software that make up the logical NetSage TestPoint. TestPoints use both active and passive measurement techniques to gather data for a broader understanding of network behavior.

TestPoints passively collect data from routers or switches using the Simple Network Management Protocol (SNMP) [10], an application–layer protocol for information about managed devices on IP networks, generally with a polling rate of once every 60 seconds. This data set includes the interface name, the number of input and output bits, any errors or discards, and data about unicast or multicast use.

The second type of passive data collected by the Test-Point is flow data from routers using tools such as NetFlow [51], sFlow [78], or IPFIX [12]. Flow data is typically sampled by a router at between 1:100 and 1:1000 packets. This data set includes information for sampled flows including the source and destination, the number of bits and packets transferred, the duration of the flow, the flow type, and the protocol and port used.

The third type of passive data collected by the TestPoint comes from packet header inspection tools running on science data archives using Tstat [45, 90], which was developed as part of the European Union (EU) Measurement Plane (mplane) FP7 project [88]. Tstat examines the packet headers for data flowing in and out of instrumented science archives and reports TCP statistics for each flow, including the congestion window size, the number of packets re-transmitted, the source and destination of the flow, the number of bits and packets transferred, the duration of the flow, the flow type, and the protocol and port used. Unlike similar data collected using a standard router flow tool, Tstat data is not sampled.

The fourth data set collected by the TestPoint is from active measurements using perfSONAR [76, 87], an open source network measurement suite designed to provide end-to-end performance metrics. There are currently over 2000 publicly registered perfSONAR nodes deployed worldwide [77]. NetSage deployments may include perfSONAR data sets for active measurements of throughput, latency, and loss.

Each of these data sets can be used in multiple ways. In the examples and figures that follow, we highlight which data sets are the source of the information given.

## 4.2 Data Ingest Pipeline

Records from the NetSage TestPoints are sent to the Data Ingest Pipeline, which consists of five stages, as shown in Fig. 2.

In the first stage, records for flows smaller than a threshold are discarded. The threshold is generally set to the size of 10 MBytes over 5 min, but can be adapted on a per-installation basis. The primary goal of the NetSage Framework is to understand large-scale data transfers, so records related to small flows do not need to be retained. With a threshold at 10 MBytes, the number of records collected is reduced by 90%, however we retain data for approximately 90% of the volume of the data transferred. This filtering has several additional benefits. It reduces both the CPU needed to run the Pipeline software as well as the overall storage space required by the NetSage Archive. It also increases the level of privacy, as flows related to emails or web page downloads are smaller than the threshold and therefore discarded.

In the second stage of the pipeline, we ensure that each record represents all of the data for a single flow by stitching together the existing multiple records for longer flows. Most flow collection techniques share data at specified time intervals, generally 5 min. If a single data transfer occurs over several time intervals, multiple partial records will be created, one for each time frame. We stitch these together to ensure each record represents a full transfer.

At the third stage, tags are added to the record to map the source and destination of the data transfer to their Autonomous System (AS) Numbers and Names [4]. Previously, the MaxMind GeoIP database [44] was used, but in December 2020, the NetSage Data Ingest Pipeline shifted to using the Center for Applied Internet Data Analysis (CAIDA) AS to Organization Mapping Data set [8], which maintains better human-readable naming. To identify institutions that do not have their own ASN, the Pipeline uses information from the Shared Whois Project (SWIP) [84]. The SWIP data is used to formally document cases where subsets of IP space need to be identified. For example, a regional network may allocate a piece of its own IP space to a member institution. SWIP enables the traffic to that IP-subset to be properly attributed to the institution and not the regional network, which in turn allows the NetSage Dashboards to list the sources and destinations more accurately.

In the fourth stage, tags are added to the record to identify the science domain and project information using the NetSage Science Registry [55]. One of the primary use cases for NetSage is to better understand which science domains and projects are using a set of networked

resources. In order to do this, we need a mapping of IP address spaces to specific science projects and domains. We created the Science Registry to do this mapping. The system supports collaborative and crowd-sourced data entry. With this system, we can add additional tags to each Flow record as needed, including the research project name, the science domain, the universities or institutions involved with the project, geo-location data, or other related data. As of April 2021, the NetSage Science Registry contains over 430 entries for over 370 Organizations and over 250 Science Projects. The Registry is continually updated by the NetSage team and its collaborators.

In the fifth and final stage, the low order bits of the IP addresses are stripped off to de-identify the data. One of the system design goals of the NetSage Framework was to avoid storing personally identifiable information (PII), and this stage of the pipeline addresses that requirement. For IPv4, the 8 low-order bits are removed. For IPv6, the 64 low-order bits are removed. Full details are given in the NetSage Data Privacy Policy [54], which was developed to balance the need for user privacy with the practical value of the data. Note that since the NetSage Archive does not include full IP addresses, there is no reference data related to a person or at a personal level, so this approach is compliant with the European Union General Data Protection Regulation (GDPR) [25] as well.

Currently, the Data Ingest Pipeline can be run in two ways—at a collection point at Indiana University or in a container on a system owned by the resource owner. The data is simply pointed to the correct collector. By using the containerized deployment, a resource owner can further control the data sharing process. Current NetSage deployments are split roughly 50–50 in how they deploy the Data Ingest Pipeline.

### 4.3 Security of NetSage Ingest Pipeline

As with any network monitoring system, the overall entire system must follow security best practices to ensure that data is not tampered with and to ensure that only the intended data is ever published. The NetSage deployment at Indiana University follows the Defense in Depth [18] philosophy. All components require two-factor authentication to access. All data transfers are via SSH connections. When the NetSage team works with other sites to access monitoring data, a security review is undertaken.

Also, as mentioned above, only public data sets that are not considered sensitive are published in public Dashboards. In cases where there are privacy concerns for the data, the NetSage Dashboards are secured by password or Shibboleth.

### 4.4 NetSage Archive

After passing through the Data Ingest Pipeline, the data sets are stored in the centralized NetSage Archive. The Archive consists of a Time Series Data System (TSDS) archive and an Elasticsearch archive. The Indiana University OmniSOC [75] hosts the NetSage Archive and provides production-quality security and support for the data resource.

TSDS was developed by the Indiana University GlobalNOC [27] to provide an efficient way to store data with consistent time intervals. It provides well-structured and high-performance storage and retrieval of time series data and metadata, in our case, SNMP and perfSONAR data.

The Elasticsearch, Logstash, and Kibana (ELK) Stack is open source software that forms a scalable system used to flexibly ingest, store, and analyze sporadic event data. The Elasticsearch archive stores data as JSON documents and indexes it for quick searching and retrieval. The NetSage Archive uses the Elasticsearch archive to store flow data and data from Tstat. One of the features of Elasticsearch is that it is designed to be horizontally scalable, meaning that both the performance and capacity of the Archive can be increased by adding more nodes to the cluster running the NetSage Archive.

The NetSage archive is structured so that if other databases of information were needed, they could also be included in the archive. In this way, each type of data can be stored so that access to the data itself is done efficiently.

### 4.5 Dashboard Visualization Components

NetSage Dashboards are used to visualize the answers to the stakeholder questions that were identified as part of the design methodology. A set of NetSage Dashboards for a particular suite of resources is referred to as a *NetSage Portal* for those resources. The Dashboards are built using the open source Grafana analysis and visualization engine. Each Dashboard consists of a set of *Visualization Components* that show different aspects of the data in response to a query. In cases where Grafana did not have a ready-made Visualization Component, new ones were developed using D3.js [15], and contributed back to Grafana when possible. In some cases, the following Figures include only partial Dashboards. Each caption includes a reference to the exact URL to see the full Dashboard if needed.

The NetSage Visualization Components range from simple to complex, each included as part of a given Dashboard to tell its part of the story to meet a design goal. On the simple end of the scale, basic line and bar charts are used to answer stakeholder questions such as "what is the

present state of a resource?" in several Dashboards, such as shown in Fig. 2.

We use *maps* to show geographical relationships for different data sets. For example, the Advanced North Atlantic (ANA) Network [1] uses the ANA Portal Bandwidth Dashboard in Fig. 3 to enable stakeholders to easily visualize and compare the volume of traffic between physical end points. In this case, it can also be used to easily see if load balancing across the links is taking place based on the collected SNMP data.

We use *Heatmaps* to show changes in values over time, where the x-axis is time and a darker color indicates a larger value. Heatmaps can answer stakeholder questions such as "Are there patterns of behaviors for the use of the resources?". This type of display can accentuate changes of behavior over time, as seen in Fig. 4, which shows the International Portal Latency Pattern Dashboard for the NSF-funded NEAAR link using perfSONAR active measurement latency data.

*Sankey Graphs* [79] are used to show relationships between items using a ribbon graphic, where the width of the ribbon indicates the quantity proportionately. We use Sankey Graphs as a visual way of answering stakeholder questions such as "Which institutions are sharing different types of science data between them?", as shown in Fig. 5 for a set of science domains using the iLight Portal Science Discipline Dashboard.

A *Bump Chart* [92] is a special form of a line plot that is well-suited for exploring changes in rank over time. We use Bump Charts to compare the rankings of Top Talkers over time, as shown in Fig. 6. This allows end users to easily see multiple observations with respect to each other, rather than the actual values itself. For example, the figure shows that for the Great Plains Network (GPN), energy science data transfers for the CMS experiments, with data sources at UNL and FermiLab, dominate the list consistently.

Finally, a *Slope Graph* [91] allows the direct visualization of the relationship between two variables. For NetSage Dashboards, Slope Graphs are used to show relationships between Sources and Destinations of flows. For example, Fig. 7 shows the Slope Graph for the FRGP Portal Flow Data Dashboard and visualizes the relationships between sources and destinations. While the bar charts (also included in the Dashboard) indicate that NOAA and NCAR/UCAR share more data, their pervasiveness as sources is more clearly described with the Slope Graph Visualization Component.

## 4.6 Other uses of the NetSage framework

### 4.6.1 ESnet peering decisions

The NetSage System is flexible enough to be applied in several other use cases. At the Energy Sciences Network (ESnet) [23], there was a need to replace a legacy flow collection system to serve the network engineers in tasks such as general capacity planning, peering analysis, and determining when to establish a direct connection to a network and at what capacity. In order to meet this need, it was decided to augment the NetSage Data Ingest Pipeline and NetSage Archive to also include information from the Border Gateway Protocol (BGP) [6]. This data set includes information about all the networks involved in delivering packets in a flow to their destination once they leave the local network. This is significant because it includes not just the source and destination networks, but intermediate networks as well.

The ESnet development team is using the standard NetSage codebase in GitHub [52] for their enhancements, which will enable each deployment to leverage the common components for their separate use cases.

### 4.6.2 DDoS detection on international backbones

The International Networks at Indiana University (IN@IU) [35] team has been using components from NetSage with slight alterations to detect Distributed Denial of Service (DDoS) attacks. This type of attack generally consists of an attacker compromising a number of hosts that are directed to send large amounts of data to a target system in an effort to overwhelm the target. Many types of DDOS attacks have traffic patterns that are easy to identify through the examination of packet header samples via flow data.

Unlike the current Data Ingest Pipeline that removes data about all flows under 10MB, this use case requires all data and also full IP addresses for the flows. Because of this, the Dashboards must also be locked down so that the data is kept secure. The NetSage team worked with members of IN@IU to adapt these components, and a new Dashboard was developed as well. The Dashboard displays the top sources, destinations, and flow pairs by highest number of flows under a defined size to identify the possibly suspect flows. It also includes a panel that shows the top single destination of these small flows along with all of the associated sources, aiding in quick identification of a potential DDoS target.
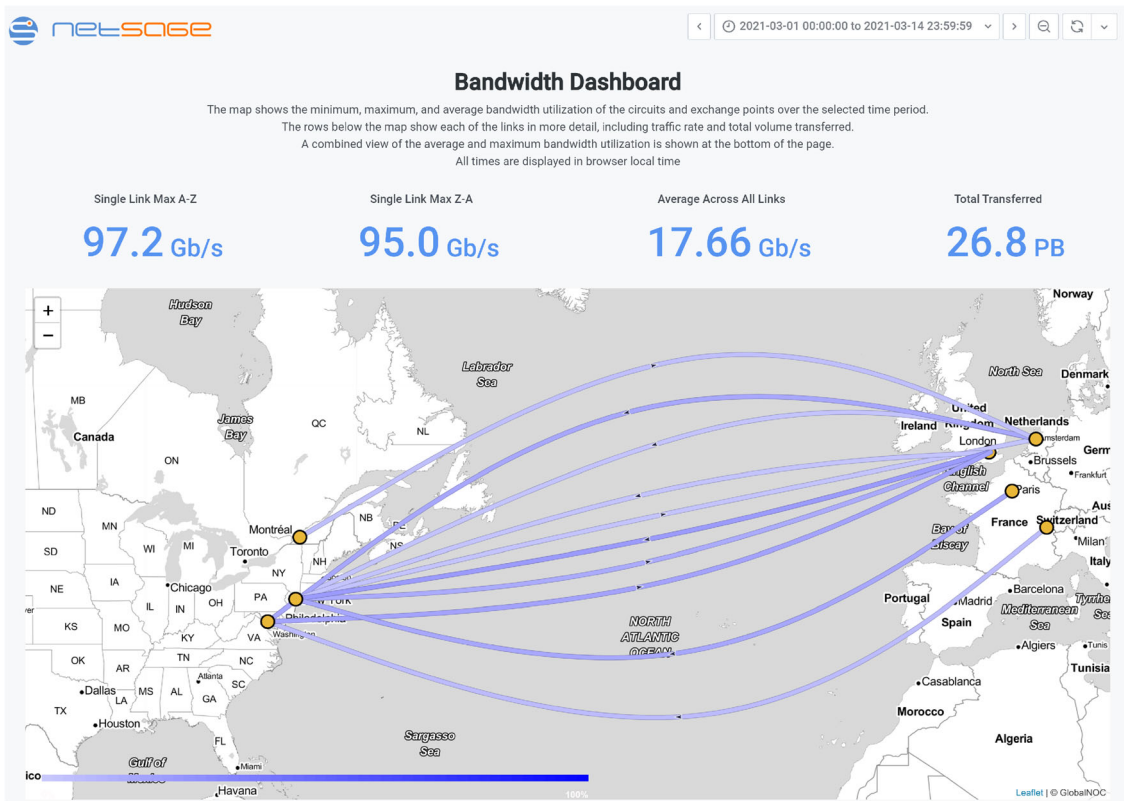
**Fig. 3** The ANA Portal Bandwidth Dashboard [59] shows SNMP data for the ten circuits that are part of the ANA consortium using a map that allows easy inspection
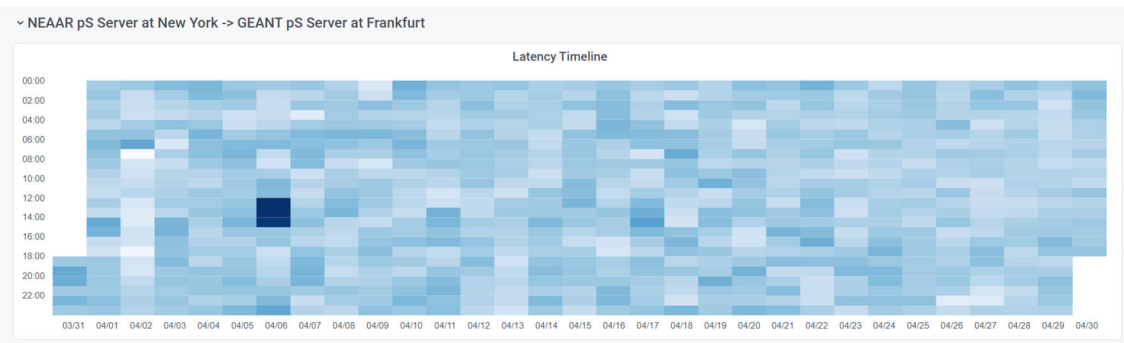


**Fig. 4** A Heatmap showing perfSONAR latency data, which is part of the International Portal Latency Pattern Dashboard for tests between the ManLan and GEANT Open London exchange points over the NEAAR link [65]. The x-axis is the day of the month and the y-axis is the time of day. Darker colors indicate larger amounts of latency
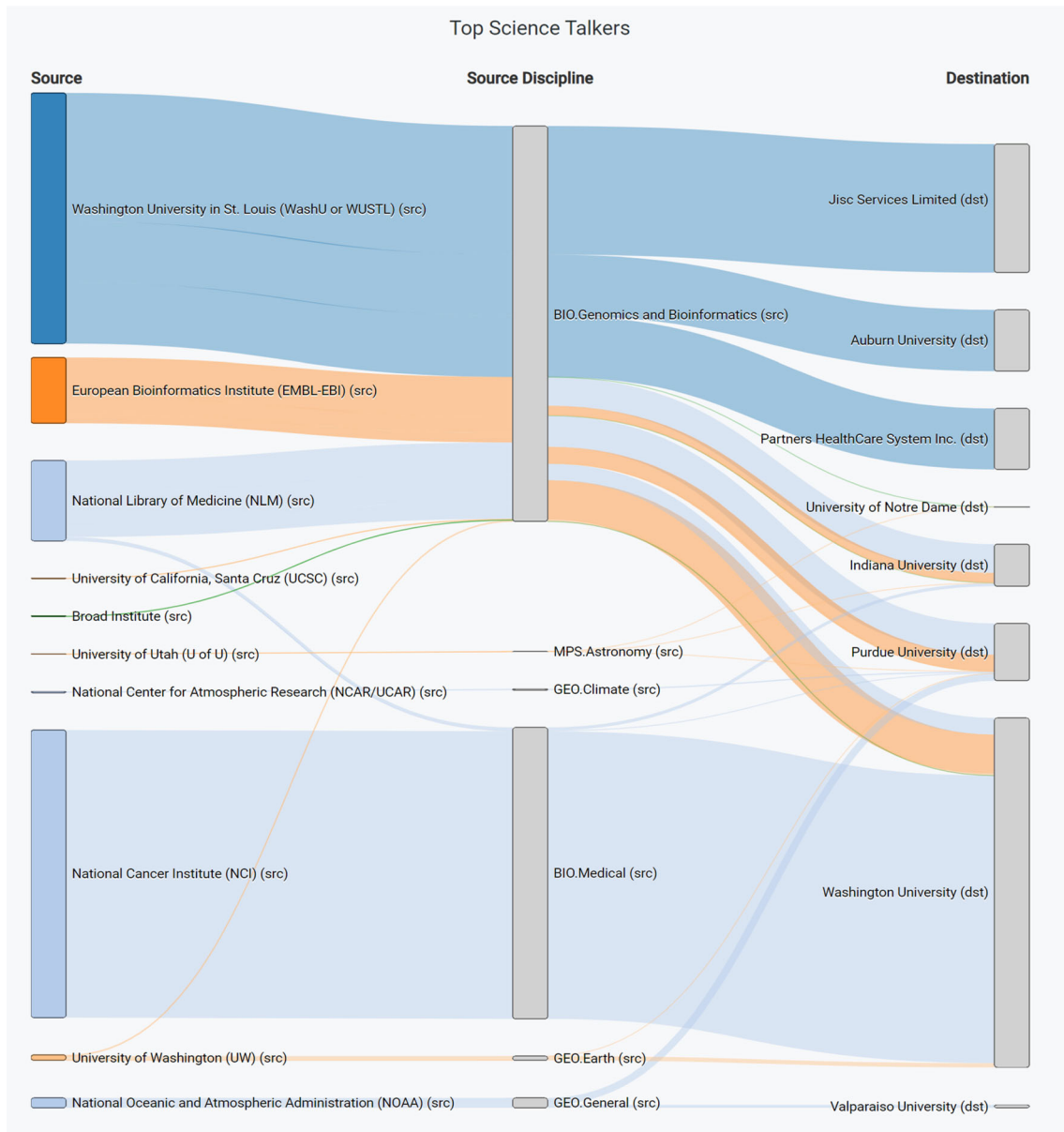
## 5 Design lessons learned

In the course of developing any large-scale pragmatic software framework, plans change and lessons are learned. Three of the major lessons as we have experienced while developing and deploying the NetSage Framework have been to adapt when necessary, to leverage other people's work as much as possible, and that what users request will change as soon as they have a prototype to work with

(sometimes referred to as "No plan survives contact with the enemy"—Helmuth von Moltke).

### 5.1 Lesson 1: Adapt when necessary

Sometimes, a change of plans can result in an improved approach. The NetSage development team had originally planned to collect both sampled and unsampled flow data from routers (not science data archives) by using a packet-header inspection software tool such as ARGUS [3], Zeek

**Fig. 5** A Sankey graph showing which institutions are sharing data for identified science domains on the iLight network [63]. The width of the ribbon is proportionate to the volume of flow data sent

(formerly called Bro) [94], tcptrace [86], or Tstat. This approach would have enabled a comparison between the different measurement approaches. However, part of how these types of packet-header inspection tools function is that they track both sides of the "conversation" between a source and a destination for each data transfer. In order to track both sides of the conversation, data sent from the source to the destination must use the same path through the network as when data is sent from the destination to the source. It was discovered that many, if not most, international data transfers experience asymmetric routing. In other words, the network path from the source to the destination was not the same as from the destination to the source. Because of this, none of this class of packet-header inspection tools could be used in the middle of the path at a router.

However, packet-header inspection tools would be able to gather data when they were placed at the end of a path, that is, at the actual source or destination. We evaluated the common sources and destination of data sets in the US and identified several major data centers. We then worked with them to deploy Tstat and a TestPoint to be able to pull the data into the NetSage Archive. We are now collecting data about transfers to and from science archives at the National Energy Research Scientific Computing Center (NERSC) [50], the Texas Advanced Computing Center (TAAC) [85],

**Fig. 6** A Bump chart showing the GPN Portal over three months [62]
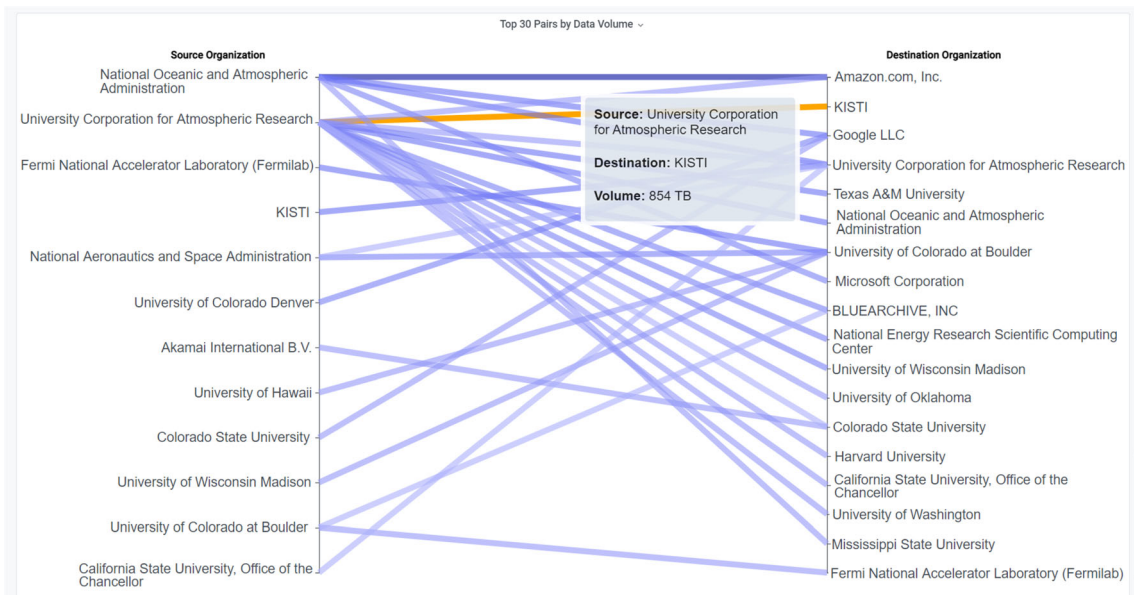


**Fig. 7** A Slope Graph from the FRGP Portal Flow Data Dashboard showing the relationship between sources and destinations [60]

the National Center for Atmospheric Research (NCAR) [49], and the University of Hawaii's Institute for Astronomy (IFA) [34].

## 5.2 Lesson 2: Leverage when possible

The NetSage development team always planned to leverage other open source projects as strongly as possible in order to maximize project resources, in part also looking toward other NSF-funded projects it could utilize. For

example, the initial NetSage Archive implementation used the existing TSDS database. When the data collection expanded to include flow data, the NetSage Archive was updated to include the Elasticsearch archive as well, as opposed to building a tailor-made one in house. Our initial implementation of the Data Ingest Pipeline used NF Dump [73] and custom scripts, which over time have been transitioned to taking advantage of logstash [40] instead. We also shifted from the commercial MaxMind database we were using to map IP addresses to organizations to use the NSF-funded CAIDA AS to Organization Mapping Dataset. Similarly, the initial NetSage NSF-funded international network Portal Dashboards and Visualization Components were written using custom software, because when the project started there was no clear best toolkit approach to building them. In Year 3, we shifted to using Grafana, which has saved countless hours of development and decreased our support burden.

### 5.3 Lesson 3: Changing requests

The NetSage development team, like most builders of pragmatic software, has also discovered there are successes and opportunities when working directly with an active user base. For each new Dashboard that has been storyboarded, designed, and deployed, once the stated requirements were met, the stakeholders will often take the opportunity to request additional functionality. An ongoing challenge has been to keep the Dashboards focused and simple enough for use by a wide audience, while adding the requested additional data visualizations. We continue to expand the use cases we address and the visualizations used to meet stakeholder needs.

## 6 Discoveries made using the NetSage framework

The NetSage Framework is used daily to interpret a variety of networking and data transfer behaviors by resource owners, science teams, and network engineers to better understand the performance and patterns involved in a wide area data movement. End users can work with NetSage Dashboards to gain insight into the data movement at the system, institution, or project level, and to see longitudinal changes in behavior in a broad set of situations.

### 6.1 How a system of resources is used

The use case that the original International Portal was developed to address was to better understand how NSF's multi-million dollar investment in international networks was being used by the US R&E community. The NetSage

team initially developed two, high-level Dashboards to give basic details about the use of a resource, in visual graphs and in summary statistics. The Bandwidth Dashboard, shown in Fig. 2 for the International Portal, uses SNMP data to generate and display a map of the resources, details about the use of each circuit, and summary line graphs for the average and maximum bandwidth utilization. In addition, the Summary Statistics Dashboard, shown in Fig. 8, highlights basic numerical data about the resource using SNMP, Flow, Tstat, and Science Registry data. It gives a high level overview of all the resources in a Portal in this way.

However, NetSage Dashboards can also take advantage of other components to answer several different questions about usage. For example, Fig. 9 shows a Heatmap from the Pacific Wave Portal Flow Data Dashboard that uses flow data to measure data transfers to and from the Zoom video conferencing hosting site during the time frame where R&E institutional use of Zoom changed radically. In March 2020, when universities started to responding to the COVID-19 pandemic-related restrictions, network resource owners needed to know how their systems were responding to the change of use. The Heatmap shows data volumes starting in February that increase on/around March 12 when many US universities declared that researchers could not travel. This was followed 10 days later by a decrease, likely caused by a combination of institutions shifting to Spring Break, institutions issuing work from home directives (so the traffic shifted to home networks not R&E networks), and Zoom shifting some of its hosting to use cloud services, rather than their own IP space.

Another example of changes at a system level that took place during the March 2020 time frame can be seen in Fig. 10. Its common for regional networks to want to track the use of their network over time, which is why NetSage has a Top Talkers over Time Dashboard. In this example, the FRGP Portal Top Talkers Dashboard for January through June, 2020, shows that while the three highest Top Talkers stay the same, a noticeable shift happens in early March, at the time that US educational institutions implemented pandemic restrictions. Note the large amount of change during March.

### 6.2 How the resources of an individual institution can change

Often, an individual institution will want to examine the traffic it sends and receives to better interpret how the institution's researchers are working with collaborators. We created the Flow Data by Institutions Dashboard to answer the core questions asked from that point of view. For example, Fig. 11 shows data for Emory University as part of the SoX Portal Flow Data by Institution Dashboard.

**Other Flow Statistics**

This dashboard provides a summary of flow information for science registry tagging, organizations, countries, protocols, ports, as well as volume and counts of flows.
All times are displayed in browser local time.

Flow Sensor Health (last 3d)

Large Flows Observed

7,803,394

| | | | |
|---|---|---|---|
| 6,350 | Source Organizations | 14,611 | Destination Organizations |
| 7,241 | Source ASNs | 15,978 | Destination ASNs |
| 199 | Source Countries | 236 | Destination Coutries |
| 729,147 | Flows with only one end tagged in Science Registry | 168,690 | Flows with both ends tagged in Science Registry |

Sensor Flow Stats

| Sensor | Total Volume | Largest Flow | Fastest Flow | # Flows |
|---|---|---|---|---|
| UNET UCAR Tstat | 606.9 TB | 32.6 TB | 16.0 Gb/s | 1.4 Mil |
| NEAAR NY-London | 579.3 TB | 560.4 GB | 147.8 Gb/s | 513.0 K |
| Hawaii LA netflow | 212.0 TB | 857.9 GB | 28.8 Gb/s | 189.4 K |
| AMPATH sFlow MCT02 | 139.2 TB | 522.0 GB | 6.8 Gb/s | 2.1 Mil |
| AMPATH sFlow BCT-MI3-... | 113.3 TB | 374.1 GB | 25.9 Gb/s | 401.6 K |
| TransPAC Seattle sFlow | 112.6 TB | 394.9 GB | 21.1 Gb/s | 414.7 K |
| AMPATH sFlow MLXE | 89.7 TB | 197.3 GB | 3.0 Gb/s | 2.1 Mil |
| dtn07.nersc.gov | 85.3 TB | 1.4 TB | 5.4 Gb/s | 34.8 K |
| dtn08.nersc.gov | 83.4 TB | 2.0 TB | 5.1 Gb/s | 24.8 K |

Count of Flows by Sensor (Stacked)

**Fig. 8** A portion of the International Portal Flow Statistics Dashboard for the first week of March 2021 [67] showing the overall statistics for the use of the NSF-funded international networking resources

Specifically, we selected the filters for this Dashboard to show how Emory interacted with international institutions over a 3 month period. During this time, for example, we can see that Emory receives almost three times more data
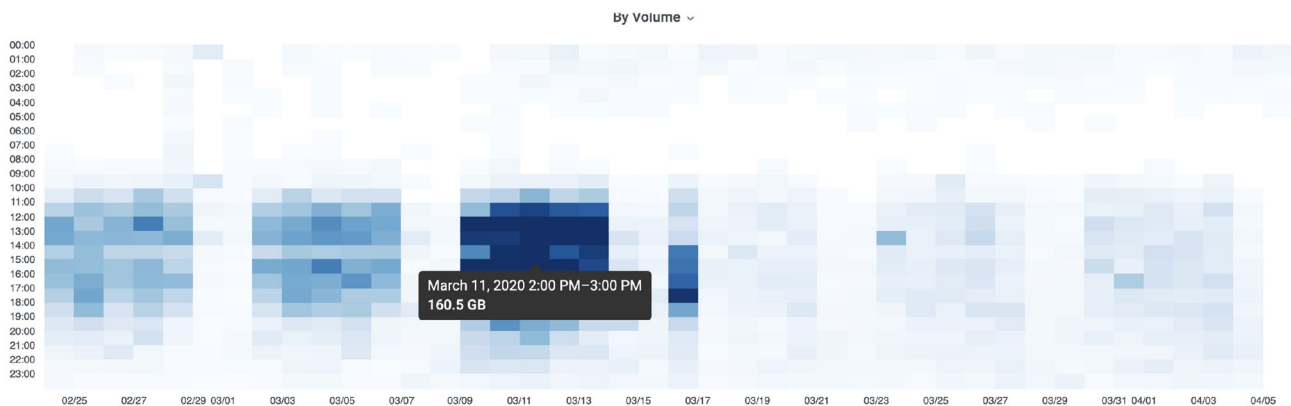
than it sends, and that the researchers are working with a wide variety of international collaboration sites. However, for both sending and receiving data, the main collaborators are in the United Kingdom, as indicated by the listing for JISC, the national R&E network for the UK.

In another example, the National Library of Medicine (NLM) hosts many bioinformatics related data sets, of particular interest to pandemic research in 2020. Figure 12 shows a Heatmap of data transfers with one end at NLM using flow data as part of the KINBER Portal for Individual Flows for the first five months of 2020, and notable are the large transfers in January and late April.

## 6.3 How the resource use of a project can change

NetSage Dashboards can also be used by specific science projects to explore how their data transfers are performing and who is accessing their science data archives. For example, the Engagement and Performance Operations Center (EPOC) [22] often uses the US domestic NetSage deployments to evaluate behaviors of data transfers as part of its Roadside Assistance process [22]. The EPOC team was part of a joint collaboration with partners at TACC, University of Southern California, and Globus, to transfer massive astronomy data sets from the Arecibo Observatory [2] in Puerto Rico to the TACC science archives [14]. The TACC Portal Individual Flows Dashboard enabled engineers to visualize the progress and speed of the the transfers by visualizing the flows between the two institutions, as shown in Fig. 13.

The PAN-STARRS collaboration has used the All Data Portal Project Dashboard for PAN-STARRS to examine the associations between various R&E entities, some expected, others unexpected. The table of sensors in this



By Volume ⌄

March 11, 2020 2:00 PM–3:00 PM
160.5 GB

**Fig. 9** A Heatmap from the Pacific Wave Portal Flow Data Dashboard that displays data related to transfers with one end point at the Zoom video conferencing facility that crossed the Pacific Wave Exchange Point to an academic institution source or destination during February–April 2020 [69]. The x-axis is the day of the month and the y-axis is the time of day

**Fig. 10** The FRGP Portal Top Talkers Dashboard for January through June, 2020 [61], showing how the Top Talkers on R&E networks changed radically during the shift in workspace when COVID-19 related restrictions were put in place in early March
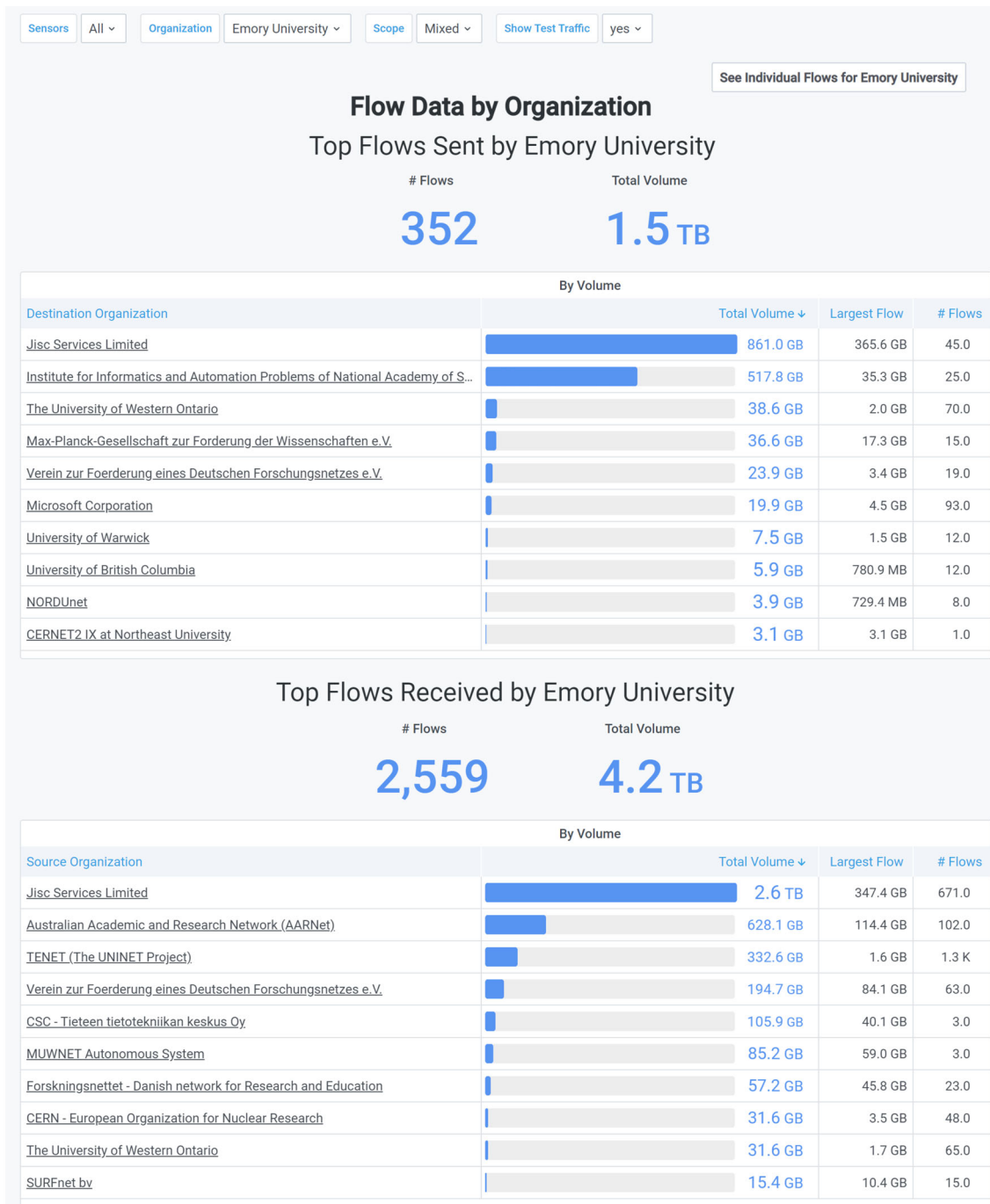
Dashboard, shown in Fig. 14, indicates where the data for the Dashboard was collected and gives a sense of the network paths that are involved when sharing PAN-STARRS data. Scanning that table shows that there is relevant traffic on AMPATH, a collector in Miami that generally shows data traveling between the US and South America. A quick check of the map shows that the South American contact is part of Associação Redes de Interconexión Universitaria (ARIU), the national R&E network for Argentina, rather than one of several Chilean astronomy sites, as might be first assumed.

### 6.4 Identifying changes or unexpected traffic

Resource owners use their NetSage Portal to track typical behavior and to identify when behavior changes occur. When unexpected changes are observed, this may lead to further investigation. For example, Fig. 15 shows part of the Analysis Dashboard for the International Portal for the

NEAAR project, the NSF-funded network collaboration between the US and Europe. End users can utilize an Analysis Dashboard, such as the one shown in Fig. 15, to relate SNMP data changes with the associated flow data to identify what transfers may be causing the shifts in behaviors. A significant change in behavior was experienced by this link in early 2018 because the US Department of Energy network operators added the NEAAR circuit to the set of network resources that support data transfers related to the Large Hadron Collider (LHC). In this particular case, NetSage Dashboards were able to exhibit this change even before the email notification was sent to the owner of the NEAAR circuit.

Another example of unexpected traffic patterns was seen on TransPAC4. TransPAC4 is the NSF-funded international network project that supports connections between the US and Asia. Engineers for TransPAC4 used the International Portal to identify an unusual pattern of behavior where every 10-12 days there was a significant
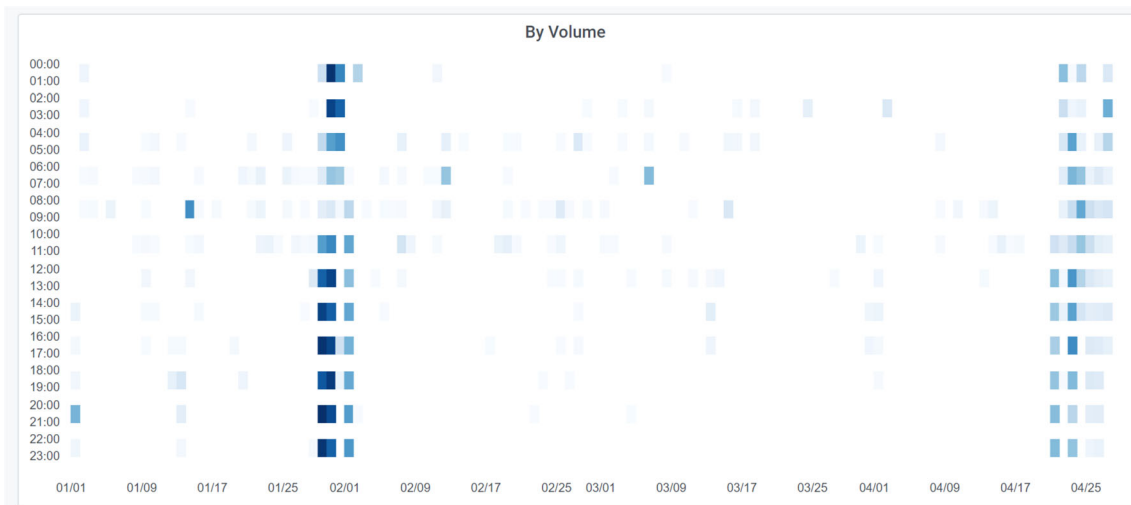
**Fig. 11** SoX Portal Flow Data by Institution Dashboard [70] showing which international institutions send or receive data from Emory University for January–March 2021

increase in the data volume over the network resource, as shown in the partial screen shot of the Analysis Dashboard in Fig. 16. Investigation indicated that the periodic data transfers were taking place between the Instituto di Radioastronomia, in Italy, and the Kashima Space Technology Center in Kashima, Japan, and that the traffic was related to an astronomy very-long-baseline interferometry

(VLBI) project. The workflow for VLBI applications involves several geographically distributed radio telescopes that are all aligned on the same celestial object, in this case, all located in Italy, sending their data to a collector site, in this case in Japan. Identifying large scale data movement patterns such as this can also give resource owners an opportunity to work with end-user scientists to identify

**Fig. 12** A Heatmap showing data transfers from the the National Library of Medicine in early 2020, as part of the KINBER Portal Individual Flows Dashboard. [68]

potential performance options. For example, in this case, additional investigation took place to ensure that this path was the most efficient for the collaboration since there were several different routing options.

## 6.5 Identifying possible routing errors

NetSage is one of several tools that can also be used to identify when data transfers are not takeing the part or route that is most efficient. For example, when additional network capacity comes on line, it is common to see routes change, often in unexpected ways. One example of this was seen when the TransPAC4 project added a new 20G connection between Guam and Hong Kong. The partial Dashboard, shown in Fig. 17, indicates that the top pair of organizations transferring data over the new connection was the Korea Institute of Science and Technology Information (KISTI) in South Korea and the Chinese Academy of Science (identified as Beijing Primezone Technologies, Inc. due to its IP space). Traffic between these locations clearly should not be going over links to the USA. A deeper investigation showed the traffic was routed across the Pacific Ocean twice, resulting in significantly decreased performance. A discussion with the network engineers overseeing different parts of the path determined that the routing preferences were incorrect, and the problem was resolved.

## 7 Related work

Over the past 20 years, the R&E network community has developed numerous monitoring portals similar to Net-Sage. These include:

– The *ESnet portal*, my.es.net [47], includes dashboards for SNMP and flow data, and breaks down data based on both links and DOE sites, but only works with US DOE sites.
– The *GÉANT Tools Portal* [26] is a collection of public and restricted tools to view a wide variety of network monitoring data, including from Cacti and Looking Glass [41].
– The *Gloriad Portal* for the Gloriad network [13], active from 2006 to 2016, included SNMP and some flow data for a small subset of international R&E links, and also had a prototype of the Science Registry.
– *MONIT*, CERN's monitoring portal MONIT for the Worldwide LHC Computing Grid (WLCG) [5, 46], provides tools for monitoring hosts and services in the CERN Data Centre and associated Computing Facilities, as well as experiment activities on the LHC Grid such as data transfers, job execution and site availability.
– *WorldView* [28], developed by IU Global NOC, uses SNMP data to show a map for many R&E networks worldwide.

Each of these portals were developed for use by a network operations center to provide a view for a single network provider or Science Experiment. NetSage is unique in that it was developed for a broader set of users over multiple networks, and to be able to analyze network performance related to data integrated from multiple networks and resources. NetSage Dashboards are designed primarily to understand performance issues and performance degradation, not outages.

In addition to the R&E community portals, there are a number of more general network measurement and
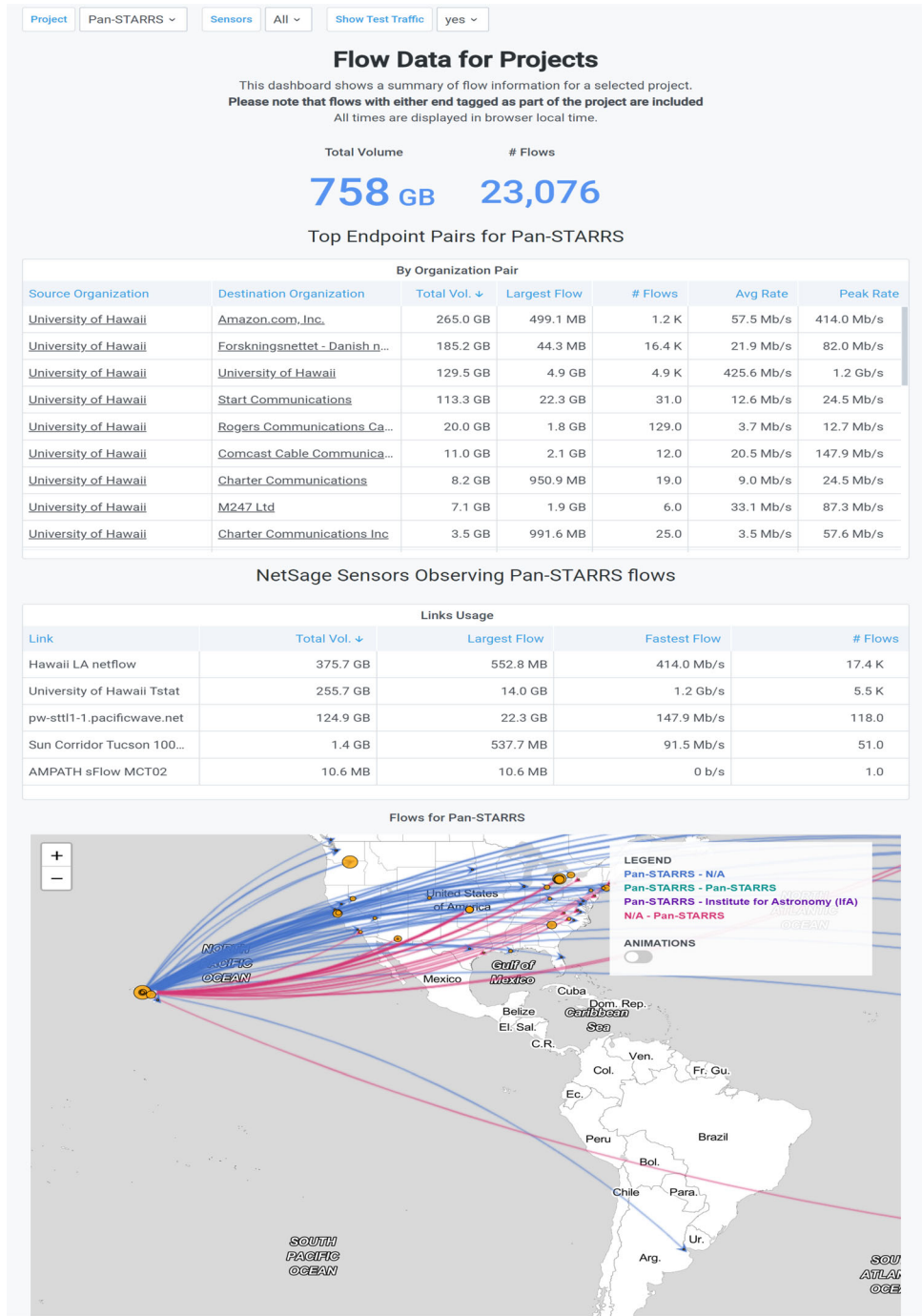
**Fig. 13** The TACC Portal Individual Flows Dashboard showing data transfers from the Arecibo Observatory in Puerto Rico to the TACC science data archive [71]

monitoring tools developed by commercial (or semi-commercial) groups, again also to primarily support network operations center staff for a single network. The most relevant of these include:

– *Argus* [3] is a set of open source tools to collect and analyze flow data.
– *Cacti* [7] is open and supports SNMP out of the box but requires extensions for other data sources. The visualization options are also limited to the default set of graphs, unlike the library of plugins available to NetSage.

– *Deep Field* [17], primarily known in the R&E networking space for its use by Internet2, it uses flow data but has strict login and access control.
– *Elastiflow* [21] is similar to NetSage in approach, but not open source or privacy aware and is deployed only in a limited setting.
– *InMon* [36] is primarily focused on sFlow supported devices and does not support the range of protocols found in NetSage.
– *Kentik* [38] has several similar Dashboards to NetSage, but is deployed on a per-resource set approach, is not

**Fig. 14** The All Data Portal Project Dashboard for PAN-STARRs during February 2021, showing how collaborators are accessing the data sites primarily in Hawaii [58]



open, nor built for collecting network data across multiple organizations.

– *Nagios* [48] is an alarming service designed to do basic evaluation of data sets and notify operators. It is not designed to store time-series data at the scale of the network metrics collected by NetSage and assumes those are managed externally.

– *NFSEN* [73] is tool for displaying line graphs of Netflow data. Its data source and visualization options are much more narrowly focused than NetSage.

– *ntop* [74] is an open source set of tools that collect network metrics such as Netflow and SNMP data. The lack of perfSONAR support and Scientific metadata limit its ability to meet all the NetSage use cases.

**Fig. 15** Traffic volume graph using SNMP data to show the increase in network traffic on the NEAAR link between New York and London for January–February 2018 [66]. Note the two colors indicate traffic in different directions on the circuit

– *SolarWinds* [83] is a commercial system for monitoring many of the same statistics as NetSage but is not open nor built for collecting network data across multiple organizations.

The NetSage Framework was influenced by, or leverages prior work from, some of these portals, such as the flow analysis capabilities of the ESnet portal and the Science Project database used in the now-defunct Gloriad Portal. None of these include all of the features supported by the NetSage Framework, for example, the ability to identify a science resource. Only Kentik and SolarWinds use both flow data and perfSONAR data, similar to the NetSage Framework, but neither of these are open source. No other monitoring system includes data such as Tstat for archival resources. In addition, none of them were created to openly share the level of detail that NetSage does to the general public, and across multiple networks and resources. A summary of all these tools is shown in Table 2. The most common feature supported by these tools that is not

**Fig. 16** A partial flow analysis [72] Dashboard showing SNMP and flow data that identifies recurring data transfers that were part of a VLBI astronomy research project between Italy and Japan over the TransPAC4 link between Seattle and Tokyo in October 2018

included in the NetSage Framework is alerting, which is planned as part of the next year's development cycle.

## 8 Conclusion and future work

In this paper, we have described the NetSage design methodology and architecture as well as a broad set of use cases and discoveries. Within the space of R&E networks, we believe it is the most comprehensive open source approach to date that enables insight into underlying resource behaviors, and as such, differs significantly from other approaches which have been developed for NOC use

only. NetSage enables insights across networks and resources by a broad set of end users, and each performance Dashboard is specifically developed to respond to a stakeholder question.

Future work will continue to be driven by stakeholder requests using our design methodology. In the short term we are developing additional visualizations as requested. For example, to identify the largest flows with the poorest performance among a set of resources, something no other monitoring system described above is capable of at this time. We continue to add more data to the Science Registry to be able to better reveal network use patterns of scientific applications, and to adapt the Project Dashboards in

| Source | Destination | Total Vol. ↓ | Largest Flow | # Flows |
|---|---|---|---|---|
| KISTI | Beijing Primezone Technologies Inc. | 35.0 TB | 72.3 GB | 1.2 K |
| University of Hawaii | Indiana University | 27.6 TB | 14.4 GB | 3.2 K |
| Indiana University | University of Hawaii | 18.7 TB | 12.9 GB | 2.1 K |
| Chinese University of Hong Kong (T… | Jisc Services Limited | 12.2 TB | 36.5 GB | 9.1 K |
| University of Pennsylvania | The University of Hong Kong | 10.9 TB | 45.5 GB | 2.6 K |
| Jisc Services Limited | The University of Hong Kong | 10.8 TB | 18.4 GB | 4.9 K |

**Fig. 17** A partial listing, from collected flow data, of the top pairs of organizations transferring data over the TransPAC4 connectivity between Guam and Hong Kong for September-December 2018 [64]. This table shows that over 35 TB of data was incorrectly moving from KISTI in South Korea to China

**Table 2** Comparison of related work

| | Functionality/ Attribute | Graphical User Interface | SNMP Data | PerfSONAR Data | Flow Data | Science Registry Data | Alerting | Open Source |
|---|---|---|---|---|---|---|---|---|
| Tool | NetSage | Y | Y | Y | Y | Y | N | Y |
| | Argus | Y | N | N | Y | N | Y | Y |
| | Cacti | Y | Y | N | N | N | N | Y |
| | Deep Field | Y | N | Y | N | N | N | N |
| | Elastiflow | Y | N | N | Y | N | N | Y |
| | GÉANT Portal | Y | Y | N | N | N | N | N |
| | Gloriad (historical) | Y | Y | N | Y | Y | N | N |
| | InMon | Y | N | N | Y | N | Y | N |
| | Kentik | Y | N | Y | Y | N | Y | N |
| | my.es.net | Y | Y | N | Y | N | N | Y |
| | Nagios | Y | N | N | N | N | Y | Y |
| | NFSEN | Y | N | N | Y | N | Y | Y |
| | ntop | Y | Y | N | Y | N | Y | Y |
| | SolarWinds | Y | Y | Y | Y | N | Y | N |
| | WLCG Dashboard | Y | Y | Y | Y | N | Y | Y |
| | WorldView | Y | Y | N | N | N | N | N |

response to researcher needs. Longer term work includes adding in alarms and alerts and exploring adaptations needed to use a NetSage Deployment in a campus environment.

## References

1. ANA: Advanced North Atlantic Network. https://internet2.edu/network/initiatives-partnerships/global-networks-and-partnerships/advanced-north-atlantic-ana/
2. Arecibo Observatory. https://www.naic.edu/
3. Argus. https://openargus.org/
4. ASN: Autonomous System Numbers. https://www.arin.net/resources/guide/asn/
5. Babik, M., McKee, S., Bockelman, B., Hernandez, E., Martelli, E., Vukotic, I., Weitzel, D., Zvada, M.: Improving WLCG networks through monitoring and analytics. EPJ Web Conf. **214**, 1 (2019). https://doi.org/10.1051/epjconf/201921408006
6. BGP: Boarder Gateway Protocol. https://en.wikipedia.org/wiki/Border_Gateway_Protocol
7. Cacti. https://www.cacti.net/
8. CAIDA AS to Org: Center for Applied Internet Data Analysis (CAIDA) UCSD AS to Organization Mapping Dataset. https://www.caida.org/data/as_organizations.xml
9. Carroll, M., Britos, L., Goldman, S.: Becoming a Design Thinker. Berg Publishers Open University, London (2012)

10. Case, J., Fedor, M., Schoffstall, M., Davin, C.: RFC 1157: Simple Network Management Protocol. Tech. rep., IETF (1990). https://doi.org/10.17487/rfc1157

11. Chen, Y.C., Lee, S., Hur, H., Leigh, J., Johnson, A., Renambot, L.: Design an interactive visualization system for core drilling expeditions using immersive empathic method. In: CHI '09 Extended Abstracts on Human Factors in Computing Systems, pp. 2671–2674 (2009). https://doi.org/10.1145/1520340.1520382

12. Claise, B., Trammell, B., Aitken, P.: Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of Flow Information. Tech. Rep. RFC 7011, IETF (2013). https://tools.ietf.org/html/rfc7011

13. Cole, G., Kim, D., Sobieski, J., Li, J., Riley, D.: GLORIAD: Global Ring Network for Advanced Applications Development (GLORIAD),IRNC:ProNet, NSF #0963058, 2010-2016. https://www.nsf.gov/awardsearch/showAward?AWD_ID=0963058

14. Continuing Arecibo's Legacy. https://www.tacc.utexas.edu/-/continuing-arecibo-s-legacy

15. D3 Data Driven Documents. http://d3js.org

16. Dart, E., Wehner, M.F., Prabhat: An assessment of data transfer performance for large-scale climate data analysis and recommendations for the data infrastructure for CMIP6. CoRR (2017). arXiv:1709.09575

17. Deepfield. https://www.nokia.com/networks/solutions/deepfield/

18. Defense in Depth. https://en.wikipedia.org/wiki/Defense_in_depth_(computing)

19. Event Horizon Telescope: Downloading data from black holes. https://beta.nsf.gov/science-matters/downloading-data-black-holes

20. Elastic. https://www.elastic.co/products

21. Elastiflow. https://github.com/robcowart/elastiflow

22. EPOC RA: Engagement and Performance Operations Center Roadside Assistance and Consulting. https://epoc.global/wp-content/uploads/2020/02/Roadside-2-pager.pdf

23. Energy Sciences Network (ESnet). http://es.net

24. Fox, L., Johnson, R.: Pacific Wave Expansion Supporting SDX and Experimentation, IRNC: RXP, NSF #1451050, 2015-2021. https://www.nsf.gov/awardsearch/showAward?AWD_ID=1451050

25. GDPR: General Data Protection Regulation. https://gdpr-info.eu/

26. GÉANT Tools Portal. https://tools.geant.net/portal/

27. GlobalNOC at Indiana University. https://globalnoc.iu.edu/

28. GlobalNOC World View. https://docs.globalnoc.iu.edu/worldview.html

29. Gonzales, A., Leigh, J., Peisert, S., Tierney, B., Balas, E., Radulovic, P., Schopf, J.M.: Monitoring Big Data Transfers Over International Research Network Connections. In: Proceedings of IEEE Big Data Congress 2017 (2017)

30. Gonzalez, A., Leigh, J., Peisert, S., Tierney, B., Lee, A., Schopf, J.M.: NetSage: Open Privacy-Aware Network Measurement, Analysis, and Visualization Service. In: Proceedings of the TNC Conference 2016 (2016)

31. Grafana. https://grafana.com

32. Ibarra, J., Clark, R., Morgan, H.L.: AtlanticWave-Software Defined Exchange: A Distributed Intercontinental Experimental Software Defined Exchange (SDX), IRNC: RXP, NSF #1451024, 2015-2021. https://www.nsf.gov/awardsearch/showAward?AWD_ID=1451024

33. Ibarra, J., Morgan, H., Cox, D.: AmLight Express and Protect (ExP), IRNC:Backbone, NSF #1451018, 2015-2021. http://nsf.gov/awardsearch/showAward?AWD_ID=1451018

34. Ifa: Institute for astronomy at the university of hawaii. http://www.ifa.hawaii.edu/

35. IN@IU: International Networks at Indiana University. https://in.iu.edu

36. InMon. https://inmon.com/

37. IRNC: International Research and education Network Connections Program . https://www.nsf.gov/funding/pgm_summ.jsp

38. Kentik. https://www.kentik.com/

39. Lassner, D., Jacobs, G., Johnson, R., Fox, L.: PIREN: SXTransPORT Pacific Islands Research and Education Network, IRNC: Backbone, NSF #1451058, 2015-2021. https://www.nsf.gov/awardsearch/showAward?AWD_ID=1451058

40. Logstash. https://www.elastic.co/logstash

41. Looking Glass. https://en.wikipedia.org/wiki/Looking_Glass_server

42. Losoff, B.: Shibboleth: A project of the Internet2 middleware initiative. Collaborative Librarianship 1, 27–27 (2009). https://doi.org/10.29087/2009.1.1.04

43. Mambretti, J., Brown, M., DeFanti, T., Chen, J.H.: StarLight SDX: A Software Defined Networking Exchange for Global Science Research and Education, IRNC: RXP, NSF #1450871, 2015-2022. http://www.nsf.gov/awardsearch/showAward?AWD_ID=1450871

44. MaxMind GeoIP Service. https://www.maxmind.com/

45. Mellia, M., Lo Cigno, R., Neri, F.: Measuring IP and TCP behavior on edge nodes with Tstat. Comput. Netw. 47, 1–21 (2005). https://doi.org/10.1016/j.comnet.2004.06.026

46. CERN MONIT Dashboards. https://monit.web.cern.ch/monit/

47. MyESnet Portal. http://my.es.net

48. Nagios. https://www.nagios.org/

49. NCAR: National Center for Atmospheric Research. https://ncar.ucar.edu/

50. NERSC: National Energy Research Scientific Computing Center. https://www.nersc.gov/

51. NetFlow - A Technical Overview (2012). https://tinyurl.com/2cyca26s

52. NetSage Source Code on GitHub. http://https://github.com/netsage-project

53. NetSage International Portal Bandwidth Dashboard. https://international.netsage.global

54. NetSage Privacy Policy. http://www.netsage.global/home/netsage-privacy-policy

55. NetSage Science Registry. https://scienceregistry.netsage.global

56. NetSage Project Website. http://www.netsage.global

57. NetSage All Data Portal. https://all.netsage.global

58. NetSage All Data Portal Project Dashboard showing Pan-STARRS Project Details February 2021. https://tinyurl.com/pcxycns9

59. NetSage ANA Portal Bandwidth Dashboard. https://ana.netsage.global

60. NetSage FRGP Portal Flow Data Dashboard Slope Graph. https://tinyurl.com/2x4rtknt

61. NetSage Front Range GigaPop Portal Top Talkers. https://tinyurl.com/fz2smnsp

62. NetSage Great Plains Network Portal Top Talkers over Time Dashboard. https://tinyurl.com/2kyk76px

63. NetSage iLight Portal Science Discipline Dashboard. https://tinyurl.com/ur8k5645

64. NetSage International Portal Flow Data Dashboard for TransPAC Guam-Hong Kong circuits for September 1-December 31, 2018. https://tinyurl.com/2rfmem7e

65. NetSage International Portal Latency Pattern Dashboard. https://tinyurl.com/39zdc6f2

66. NetSage International Portal Flow Analysis Dashboard for NEAAR circuit, January 25-February 11, 2018. https://tinyurl.com/yrj8wxm5

67. NetSage International Portal Flow Statistics Dashboard for the first week of March 2021. https://tinyurl.com/achc9yh8

68. NetSage KINBER Portal Individual Flow Data Dashboard for National Library of Medicine, January-May 2021 . https://tinyurl.com/4meev2c2

69. NetSage PacificWave Portal Individual Flows Dashboard for Zoom Video Communications for February 2020 - April 2020. https://tinyurl.com/y42ve8kc

70. NetSage SOX Portal Flow Data per Organization Dashboard for Emory University, January-March 2021 . https://tinyurl.com/fw5ackud

71. NetSage TACC Portal Individual Flows Dashboard showing data transfers between Arecibo, Puerto Rico, and the TACC science data archive. https://tinyurl.com/2x8538ja

72. NetSage International Portal Analysis of VLBI Traffic. https://tinyurl.com/aa4s4ska

73. nfdump and NfSen Overview. https://www.first.org/conference/2006/papers/haag-peter-papers.pdf

74. ntop. https://www.ntop.org/

75. OmniSOC. https://omnisoc.iu.edu/

76. PerfSONAR. http://www.perfsonar.net/

77. PerSONAR Nodes Worldwide. http://stats.es.net/ServicesDirectory/

78. Phaal, Panchen, Mckee: InMon Corporations sFlow: A Method for Monitoring Traffic in Switched and Routed Networks. Tech. Rep. RFC 376, IETF (2001). https://tools.ietf.org/html/rfc3176

79. Sankey Diagram. https://en.wikipedia.org/wiki/Sankey_diagram

80. Schopf, J., Addleman, H.: TransPAC4 - Pragmatic Application-Driven International Networking, IRNC:BackBone, NSF #1450904, 2015-2021. http://www.nsf.gov/awardsearch/showAward?AWD_ID=1450904

81. Schopf, J., Leigh, J., Lake, A.: NetSage - An Open, Privacy-Aware, Network Measurement, Analysis, and Visualization Service, IRNC:AMI, NSF # 1540933, 05/2015-04/2022. https://www.nsf.gov/awardsearch/showAward?AWD_ID=1540933

82. Schopf, J., Moynihan, E., Fryer, T.: NEAAR: Networks for European, American, and African Research, IRNC: Backbone, NSF #1638863, 2016-2021. http://www.nsf.gov/awardsearch/showAward?AWD_ID=1638863

83. SolarWinds. https://www.solarwinds.com

84. Shared Who Is Project (SWIP). https://en.wikipedia.org/wiki/Shared_Whois_Project

85. TACC: Texas Advanced Computing Center. https://www.tacc.utexas.edu/

86. Tcptrace. https://linux.die.net/man/1/tcptrace

87. Tierney, B., Boote, J., Boyd, E., Brown, A., Grigoriev, M., Metzger, J., Swany, M., Zekauskas, M., Zurawski, J.: perfSONAR: Instantiating a global network measurement framework". In: Proceedings of the SOSP Workshop on Real Overlays and Distributed Systems (2009)

88. Trammell, B., Casas, P., Rossi, D., Bar, A., Houidi, Z., Leontiadis, I., Szemethy, T., Mellia, M.: mPlane: an intelligent measurement plane for the internet. IEEE Commun. Mag. **52**, 148–156 (2014). https://doi.org/10.1109/MCOM.2014.6815906

89. TSDS: Time Series Data System. https://docs.globalnoc.iu.edu/software/measurement/tsds.html

90. Tstat: TCP Statistic and Analysis Tool. http://tstat.polito.it/

91. Tufte, E.: The Visual Display of Quantitative Information. Graphics Press, Cheshire (1983)

92. Tufte, E.: Envisioning Information. Graphics Press, Cheshire (1995)

93. Turner, K., Balas, E., Baveja, S., Doyle, D., Ensman, L., Faci, S., Gonzalez, A., Khanal, M., Lake, A., Leigh, J., Seto-Mook, T., Southworth, D., Tierney, B., Schopf, J.M.: The NetSage Measurement Framework: Design, Development, and Discoveries. In: Innovating the Network for Data-Intensive Science (INDIS) 2020 (2020)

94. Zeek. https://zeek.org/

**Jennifer M. Schopf** is the Director of International Networks at Indiana University, where she oversees several NSF-funded networking awards, TransPac, and Networks for European, American, and African Research (NEAAR), all of which support network infrastructure for international collaborations. She also leads the NetSage collaboration, which provides measurement and monitoring for NSF-funded international circuits and the Engagement and Performance Operations Center(EPOC) which supports both domestic and international science engagement to improve research outcomes. As part of all these projects, there is an emphasis on the need to understand the end-user needs to enable collaborations and science engagement. Prior to IU, she was an NSF program officer and helped to develop pragmatic networking solicitations, as well as supporting several data and cyberinfrastructure solicitations.



**Katrina Turner** has a BEd in Education and Mathematics. She is currently a graduate student in Computer Science at the University of Hawai'i at Mānoa and a research assistant at the Laboratory for Advanced Visualizations and Applications (LAVA).



**Dan Doyle** is a senior developer and manager at GlobalNOC in Bloomington, Indiana. He and his team are responsible for data collection and analysis, as well as day to day operations of all the systems.

**Andrew Lake** is an ESnet software engineer and project leader at Lawrence Berkeley National Laboratory (LBNL). He is currently a lead on perfSONAR, an open source distributed network measurement and monitoring framework with thousands of deployments around the globe. He was also one of the original developers of ESnet's On-Demand Secure Circuits and Advance Reservation System (OSCARS), an award-winning production bandwidth reservation system. His research interests include measurement, monitoring, advanced troubleshooting techniques, and automation of distributed systems.



**Jason Leigh** is the Director of LAVA: the Laboratory for Advanced Visualization & Applications, and Professor of Information & Computer Sciences at the University of Hawai'i at Mānoa. He is also Director Emeritus of the Electronic Visualization Lab and the Software Technologies Research Center at the University of Illinois at Chicago, where he maintains appointments in the Computer Science and Communications departments. His research expertise includes: Big data visualization; virtual reality; high performance networking; and video game design. He is co-inventor of the CAVE2 Hybrid Reality Environment, and SAGE: Scalable Adaptive Graphics Environment software, which has been licensed to Mechdyne Corporation & Vadiza Corporation, respectively.



**Brian L. Tierney** is a project consultant. He was a Staff Scientist and groupleader of the ESnet Advanced Network Technologies Group at LawrenceBerkeley National Laboratory (LBNL) before retiring in 2017. His research interests include high-performance networking and network protocols;distributed system performance monitoring and analysis; network tuningissues; and the application of distributed computing to problems inscience and engineering.