# Project Status Report

**Reporting period:** **10/01/2023 - 03/31/2024**
**Project title:**
**Mid-Scale RI-1 Design Project (M1:DP):**
**Designing a Global Measurement Infrastructure to Improve Internet Security (GMI3S)**

**Principal Investigator:** kc Claffy, Bradley Huffaker (UCSD), David Clark (MIT)
**Project Manager:** Elena Yulaeva
**Lead Institution:** CAIDA, UCSD
**Other Institutions:** NSRC (U Oregon), MIT
**Cognizant PO:**

## Table of Contents

# 1. Summary of project status

A brief summary of the project's overall status on technical progress, cost and performance.

| Award Duration | Start date: 10/01/2021 | Planned close out: **09/30/2025\*** |
| --- | --- | --- |
| Project Finish Date | Planned Early Finish: | Estimated Early Finish: |
| Project %-complete | 67% (M30) | |

# 2. Near-Term Milestones

Below are milestones with the scheduled dates or actual/forecast/revised (A/F/R) dates that are in current reporting period and until the end of the project, and milestones (with past scheduled dates) that are delayed to future reporting period. (**Completed deliverables have bold font dates**. Red dates have slipped from their original schedule completion dates.) Note that we found the need to revisit some of the milestones as we learned more about other dimensions of the project and as the ecosystem evolved. The date marked with "R" indicates that the WBS element was reopened due to new requirements and was subsequently completed by that day.

| WBS | Subsystem | Milestone | Scheduled Date | Actual date (A) /Forecast Date (F) |
|---|---|---|---|---|
| **1.1** | **Design Infrastructure for Data Acquisition** | | | |
| | 1.1.3 | Monitors Requirements documented | 10/31/2024 | 10/31/2024 (F) |
| | 1.1.3.7 | DNS monitoring needs compiled | 04/30/2024 | 07/31/2024 (F) |
| | 1.1.4 | Monitor specification final report | 09/30/2024 | 03/31/2025 (F) |
| | 1.1.4.2 | Monitor specification report Draft (2) for community feedback | 05/31/2024 | 10/31/2024 (F) |
| | 1.1.4.3 | Post report for public comments | 08/31/2024 | 12/31/2024 (F) |
| | 1.1.4.4 | Incorporate public comments & publish final document | 01/31/2025 | 03/31/2025 (F) |
| | 1.1.5 | Monitor software Prototyped | 03/31/2024 | 12/31/2024 (F) |
| | 1.1.5.5 | DNS data monitoring software prototyped | 03/31/2024 | 12/31/2024 (F) |
| | **1.1.6** | **Monitor deployment prototyped** | 03/31/2024 | **03/31/2024 (A)** |
| | **1.1.6.3** | **Two-way traffic monitor deployment** | 01/31/2023 | **10/31/2023 (A)** |
| | **1.1.6.4** | **BGP data monitoring deployment** | 09/30/2022 | **01/31/2024 (R)** |
| | **1.1.6.5** | **Active probing measurements data monitoring deployment** | 03/01/2023 | **10/31/2023 (R)** |
| | **1.1.6.9** | **Additional 10 nodes of multiple measurements deployed** | 03/31/2024 | **03/31/2024 (A)** |

| | | | | |
|---|---|---|---|---|
| | 1.1.7 | Evaluation report of Data Acquisition Component | 04/30/2024 | 08/31/2025 (F) |
| | **1.1.7.1** | **Software to support active probing measurements deployed** | **09/30/2022** | **10/31/2023 (R)** |
| | **1.1.7.2** | **Third-party experiment deployed** | **01/31/2024** | **01/31/2024 (A)** |
| | 1.1.7.3 | Evaluation report of Data Acquisition Component created | 03/31/2024 | 07/31/2025(F) |
| | 1.1.7.3.1 | BGP data acquisition (GILL) evaluated and enhanced | 01/31/2025 | 03/31/2025 (F) |
| | 1.1.7.3.2 | Active measurements evaluated and enhanced | 01/31/2025 | 05/31/2025 (F) |
| | 1.1.8 | Virtualization capabilities prototyped | 03/31/3024 | 1/31/2025 (F) |
| | 1.1.8.3 | Virtualization capabilities documented and evaluated; report published | 03/31/3024 | 11/30/2024 (F) |
| | 1.1.8.5 | Scamper containerization completed | 07/31/2024 | 1/31/2025 (F) |
| | 1.1.9 | Vulnerabilities and Data Needs Report | 08/31/2024 | 07/31/2025 (F) |
| | 1.1.9.9 | Newly identified dataset added (NEW) | 04/30/2025 | 04/30/2025 (F) |
| | 1.1.9.10 | Vulnerabilities and Data Needs Report YR4 revision published (NEW) | 07/31/2025 | 07/31/2025 (F) |
| **1.2** | | **Design Infrastructure for Data Management** | | |
| | **1.2.1.2** | **Community feedback incorporated into the data storage requirements document** | 3/31/2023 | **11/30/2023 (A)** |
| | 1.2.2 | Data storage systems specifications published | 7/01/2024 | 09/30/2024 (F) |
| | **1.2.2.2** | **Data storage systems specifications documented, draft posted for stakeholders' review** | 2/28/2023 | **01/01/2024 (A)** |
| | 1.2.2.3 | Community feedback incorporated into document | 02/29/2024 | 08/31/2024 (F) |
| | 1.2.3 | Data and Metadata Standards Final Report | 08/31/2024 | 07/31/2025 (F) |
| | 1.2.3.1 | Develop metadata schema for existing CAIDA datasets | 09/30/2024 | 11/30/2024(F) |
| | 1.2.3.4 | Annual revision of data and metadata specifications | 09/30/2024 | 02/28/2025(F) |
| | 1.2.3.5 | Report on the state-of-the-art metadata generation approaches (NEW) | 01/01/2025 | 07/31/2025 (F) |

| | | | | |
|---|---|---|---|---|
| | 1.2.4 | Specification of tools for data curation and documentation | 08/31/2024 | 07/31/2025 (F) |
| | 1.2.4.3 | Specification of tools for data curation and documentation, report, annually updated | ongoing | **03/31/2024(A)** 03/31/2025 (F) |
| | 1.2.4.4 | Final report on specification of tools for data curation and documentation | 08/31/2025 | 07/31/2025 (F) |
| | 1.2.5 | Design data and metadata API for access to various datasets | 09/30/2024 | 8/31/2025 (F) |
| | 1.2.5.2 | Increase the number of supported data sources, including non-CAIDA datasets | 03/11/2024 | 07/31/2025 (F) |
| | 1.2.5.2.1 | **Incorporate in catalog datasets used by IYP** | **7/01/2023** | **03/31/2024(A)** |
| | 1.2.5.4 | Enhance and document data and metadata APIs, update annually | ongoing | 08/31/2025 (F) |
| | 1.2.5.4.3 | Year 3(&4) Update of data and  metadata APIs | 09/30/2024 | 03/31/2025 (F) |
| | 1.2.5.4.4 | Final report on data and metadata APIs published | NEW | 07/31/2025 (F) |
| | 1.2.6 | SDK libraries developed | 08/31/2024 | 07/31/2025 (F) |
| | **1.2.6.2** | **New libraries  created** | **03/31/2024** | **03/31/2024(A)** |
| | 1.2.6.3 | Libraries evaluated and enhanced, report published | 01/31/2025 | 07/31/2025 (F) |
| | 1.2.7 | Tools for additional data sources integration created | 01/31/2025 | 03/31/2025 (F) |
| | 1.2.7.1 | AS path annotations implemented | 01/31/2025 | 03/31/2025 (F) |
| | 1.2.7.2 | Report on LLM implementation | 01/31/2025 | 01/31/2025 (F) |
| **1.3** | | **Design Infrastructure for Broad Usability** | | |
| 1.3 | 1.3.1 | Data discovery tools prototyped | 07/01/2023 | 06/30/2025 (F) |
| | 1.3.1.5 | Report on integration of automated meta-data/data citation creation into catalog. | 07/31/2024 | 06/30/2025(F) |
| | 1.3.2 | Software for disclosure control developed | 01/31/2025 | 01/31/2025 (F) |

| | | | | |
|---|---|---|---|---|
| | **1.3.2.5** | **At least two practices to enable legit research access to various data types prototyped and evaluated** | **03/31/2024** | **03/31/2024 (A)** |
| | 1.3.2.7 | Resource Portal prototype deployed (NEW) | 01/31/2025 | 01/31/2025 (F) |
| | 1.3.3 | Report on Policy tools | 03/31/2024 | 03/31/2025 (F) |
| | 1.3.3.8 | New agreements designed and shared | 05/31/2025 | 01/31/2025 (F) |
| | 1.3.4 | Case studies on Extensibility | 09/30/2024 | 09/30/2024 (F) |
| | 1.3.4.2 | State of Internet report created and shared (1.3.4.2 Conduct meetings and compile data to create community-authored "State of the Internet report") | 9/30/2023 | 9/30/2024 (F) |
| | **1.3.4.3** | **"State of the DDoS attacks" report** | **3/30/2023** | **2/28/2024(A)** |
| | **1.3.4.5** | **Appropriate external datasets and tools, and organizations to include in the extensibility case studies identified, documented** | **7/31/2023** | **12/31/2023 (A)** |
| | **1.3.4.6** | **Case studies on extensibility of policy framework conducted, documented and shared** | 03/31/2024 | **03/31/2024 (A)** |
| **1.4** | **Infrastructure for Outreach** | | | |
| | 1.4.2.2 | Virtual collaboration environment evaluated and improved | 09/30/2024 | 09/30/2024(F) |
| | 1.4.3 | STEM workforce task completed | 9/30/2024 | 09/30/2024 (F) |
| | 1.4.3.3 | Video tutorials on nodes deployment and management created | 3/31/2023 | 09/30/2024(F) |
| | 1.4.4 | Quarterly calls conducted, minutes shared | ongoing | 07/31/2025(F) |
| | 1.4.5 | Project Presentations | ongoing | 08/31/2025 (F) |

# 3. Executive Summary

This progress report introduces the latest advancements in the CAIDA NSF MSRI design project, which aims to develop the next generation of Internet measurement infrastructure to enhance the security and utility of Internet measurements. Our project focuses on creating innovative platforms and tools for data collection, curation and utilization, particularly targeting data related to the security vulnerabilities within the packet carriage layer of the Internet, which often lead to significant harm.

The core components on which we focus our new capabilities for measurement include the Internet's

addressing architecture, the Border Gateway Protocol (BGP), the Domain Name System (DNS), the Certificate Authority system, and Distributed Denial of Service attacks (DDoS). These elements are crucial for the overall functionality of the Internet, requiring a unified approach to mitigate risks and address the often misaligned incentives to take preventive actions. Distributed Denial of Service (DDoS) attacks exploit vulnerabilities in network nodes and the basic packet-forwarding functions to overwhelm network segments. Addressing DDoS attacks is critical as effective mitigation often depends on broader operational practices within the Internet ecosystem, beyond just measures taken by the direct victims.

This report summarizes our efforts and outlines our design of the infrastructure tailored to various types of data, including traffic data, routing (BGP) data, active measurements, and DNS data. We discuss the current capabilities, their limitations, community requirements, and strategies for managing security and privacy concerns. Further, we have developed preliminary specifications for the required infrastructure components and developed methodologies designed to optimize data collection and storage. We have also explored current and potential approaches to data analysis and visualization, addressing the needs for standardization, interoperability, AI readiness, and compliance as we advance in our design and prototyping phase.

In this executive summary we first summarize the state-of-the-art and our progress by each data component (traffic data, routing data, active measurement, DNS data), followed by an update on tasks that span all data types.

**Traffic: One-Way: Telescope Data Monitor**

> **Current state of the art and its limitations**: Over the last two decades, CAIDA (at UCSD) has operated the world's largest Internet traffic observatory (UCSD-NT) to capture Internet background radiation (IBR) from a darknet. CAIDA's UCSD-NT platform enables researchers to access the captured IBR traffic data for security studies, e.g., characterizing distributed denial of service attacks (DDoS), network censorship, and spread of malware. This network telescope is a passive traffic monitoring system, capturing unsolicited traffic directed toward a large segment of mostly unutilized IPv4 address space, although it does encompass a few blocks of utilized IP addresses. The infrastructure captures (unprecedented for the network research community) O(1TB) per day. The data collection pipeline includes capturing raw packets and processing them into a more compressed flow record format for archiving. In parallel, we also extract thousands of time-series statistics directly from the packet headers. Over the years we have invested considerable effort in evaluating existing data management and curation capability to determine how to lower the barrier to research use of a large volume traffic data set. Also, given the scarcity of IPv4 address space needed to create such instrumentation, the only realistic way to sustain such data collection is to extract unsolicited traffic from active Internet address space, and during the course of this project we wrote a CICI proposal (STARNOVA, see next paragraph) to develop new campus infrastructure to explore design and deployment of such *greynets* in the specific UCSD campus infrastructure scenario. Another barrier is inconsistent data formats; each instrumentation in the community uses different formats for their traffic, flow, or attack inference data, making it challenging, expensive, or impossible to compare them.

> **Progress during reporting period:** During this reporting period, we completed our plan for deployment of a new hardware prototype for telescope deployment, in cooperation with a new NSF CICI project funded in 2023. However, in the meantime, the existing infrastructure against which we are measuring our design encountered several problems with data integrity that we investigated in collaboration with security researchers who brought them to our attention. These

6

efforts informed our new hardware and software design, but they were also unforeseen risks that took time away from our design efforts. A related challenge had to do with a multistakeholder working group we led to use the data for a "state of the Internet" report related to DDoS attacks, which required overhauling our own data pipeline to infer DDoS attacks from this data source, in a way that was comparable to data sources shared by other members of the working group.

**Traffic: Two-way: Internet backbone 100Gbs monitor.**

**<u>Current state of the art and its limitations</u>**: For decades it has been virtually impossible for researchers to get access to passively collected traffic data from Internet backbone links due to privacy concerns. Based on trust relationships that have been maintained for over two decades, CAIDA has been able to measure strategic links in the backbone so long as CAIDA could provide funding for the monitor. Since April 2008, CAIDA's passive traces dataset contains traces collected from high-speed monitors on a commercial backbone link, and anonymized for sharing with the research community. Six times in the last 20 years this backbone infrastructure has been upgraded beyond the scope of the budget CAIDA has for monitoring (OC3, OC12, OC48, and OC192, 40GB, and now 100GB ). CAIDA's (and, to our knowledge, the Internet's) last remaining single last point of public insight into the commercial Internet backbone was lost in January 2019.

**<u>Progress during reporting period:</u>** Although this phase took us 5 years, it was this reporting period when we finally got the two-way traffic monitor installed, recorded two one-hour captures (February and March), analyzed the data and produced statistical metadata. (This work was reported in our 2024 ILANDS CNS-2120399 report)

**<u>Work left to complete</u>**: In the next year we will be working with a small group (of the hundreds of users of the previous such data sets) of bet users to explore the best options for sharing this O(1TB) data set, and integrate the lessons learned into this data acquisition and management component.

**Interdomain Routing (BGP) data.**

**<u>Current state of the art</u>**: Our evaluation of the current state of BGP data collection revealed that scaling up data collection to keep pace with the growth of the Internet routing system would require an enormous increase in data volume and number of peers. Collecting global BGP data faces a fundamental cost-benefit trade-off. The information-hiding character of BGP requires collecting routes from as many BGP routers, (vantage points or VPs) as possible. But in practice the BGP protocol extensively propagates connectivity messages, leading to highly redundant (along with significant unique) information coming from each peer. The result is a data set with enormous redundancy and yet dangerous visibility gaps. The platforms' policies to store a snapshot of the aggregated data every few hours, as well as every BGP update received in between these snapshots, exacerbates the storage of redundant data. Continued growth of the Internet (~75k ASes and ~1M globally announced prefixes) and increasing connectivity between networks further burden data collection and use. Researchers often resort to sampling the data, using only a sample of the VPs, neglecting the connectivity uniquely visible to other VPs. Finally, the manual vetting of new peers also strains platform scalability. The platforms collectively peer with only~1% of the observably active ASes on the global Internet. Despite continued addition of peers, RouteViews and RIPE RIS' coverage in terms of fraction of ASes they peer with has remained flat for two decades.

**Progress during reporting period:** Our biggest achievement to date is a new design for a global BGP monitoring platform that will scale to orders of magnitude more vantage points (peers), to overcome many limitations of the current systems. The design was led by a PhD student who visited CAIDA/UCSD for summer 2023. We presented the peer-reviewed initial design to ACM SIGCOMM Hotnets and submitted a fleshed out design to ACM SIGCOMM in January 2024. We deployed a prototype system at https://bgproutes.quest where we have invited R&E networks to peer. To support interpretation of the massive amount of collected BGP data, we also designed and prototyped an infrastructure platform that would automatically infer the geolocation semantics of BGP communities. In the meantime, as with the telescope, continued pressure on existing infrastructure required streamlining of pipelines and updating of software, which further informed our redesign and in particular the need for automated ingesting of new vantage points.

**Work left to complete**: We will experiment with an expanded deployment of the new data acquisition components described above, and evaluate the data-compression tool chain (development of which is funded by CNS-2120399).

## Active Measurement platform

**Current state of the art and its limitations**: Network operators and researchers often require the ability to conduct active measurements of networks from a specific location in order to understand some property of the network. However, obtaining access to a vantage point at a given location is challenging, as significant trust barriers may prevent access: a platform operator has to provide vantage point hosts with guarantees about the activity that their vantage points will exhibit, and a platform operator has to trust that users will stay within those guarantees. Current access control to active measurement infrastructure has two extremes: access that allows trusted users to run arbitrary code, and API access that allows arbitrary users to schedule a (limited) set of measurements and obtain their results. Prior efforts in active measurement design (Scriptroute, Packetlab) have focused on interfaces that allow a user to request a host to send arbitrary packets, leaving the implementation of the measurement to the user. However, this design pattern limits the guarantees that a platform operator can provide a vantage point host regarding what their vantage point will actually do. A domain-specific language for conducting active measurements can alleviate these concerns because it (1) allows a platform operator to specify the measurements that a user can run, and communicate to the host what their vantage point will do, (2) provides users reference implementations of measurement applications that act as building blocks to more complex measurements. We present our recent accomplishments in this context.

**Progress during reporting period:** We made tremendous progress on this component. Our biggest achievement to date is a new design for a global community-oriented active measurement infrastructure that will scale to orders of magnitude more vantage points (peers). We developed new software libraries that provide measurement primitives available on each of the active measurement nodes. We implemented the domain-specific language (DSL) approach and connected 120 Ark nodes to a server running software we implemented to centralize execution of domain-specific measurement. We also designed and tested and evaluated a new Kafka infrastructure to support coordination among nodes, a certificate authority system to allow authentication of remote nodes, and a monitoring management platform. We wrote a series of blog entries to describe how to use the new software tooling. We created 17 accounts from 13 different organizations for researchers in the community to test and evaluate the prototype system.

Due to an unforeseen loss of senior personnel, we accelerated some of our implementation work, including a complete redesign and reimplementation of the software development packaging for creating several of our flagship data sets, modernizing and automating all components. We

leveraged this as an opportunity to evaluate our new data acquisition components. We also undertook third-party experiments by researchers in the U.S. and global R&E community. The unexpected retirement meant that we had to reallocate some resources to complete an advanced topology query system developed to enable discovery of the full potential value of extensive raw Internet end-to-end path measurement datasets.

Finally, we completed the design for a new IXP-based active measurement infrastructure to allow us to gain potentially hundreds of vantage points with one hardware deployment. This unforeseen new data acquisition component resulted from learning that we could not execute our originally proposed design of integrating the active measurements with existing RouteViews BGP route collectors, due to emerging hesitation from that subcontractor.

**Work left to complete**: We will deploy and evaluate a prototype new IXP-based Ark node, including routing, management, and containerization software to support this new architecture.

### Domain Name System (DNS) data platforms

**Current state of the art and its limitations:** We have undertaken extensive study of DNS research measurement needs, and recognize that this category of data could merit its own entire MSRI Design project. We have expanded the categories of DNS data we considered in our last report to include: (1) active probing of DNS infrastructure; (2) passive DNS measurements (traffic capture of queries and responses); (3) zone files (4) domain blacklist data; (5) logs from DNS servers; (6) registration metadata (e.g., owner, hosting registrar); (7) domain pricing information; (8) evidence of role of DNS in various attack chains, e.g., misdirection; (9) estimates of actual harms due to DNS-related attacks. The best known and most comprehensive platform for this type of data is the U. Twente's OpenIntel project (https://openintel.nl), although this project only probes from one vantage point and the data is not widely shared. In terms of passive data, the best data source today for studying macroscopic attack surfaces in the DNS is ICANN's Centralized Zone Data Service (CZDS) which is limited in coverage (the subset of TLDs under the governance of ICANN contracts) and granularity (24-hour daily snapshots). ICANN supports no public archive of historical snapshots, so CAIDA has supported an indexed database of Top Level Domain (TLD) zone files stretching back over a decade via our DNS Zone Database (DZDB.caida.org).

**Progress during reporting period:** Leveraging our advisory committees and collaborators, we gained a thorough understanding of how researchers consider future needs for active measurement of DNS infrastructure, and to what extent the redesigned Ark infrastructure can support it. We also demonstrated a new design of a platform that relies on active measurement and DNS data to extensively label Internet topology (with geolocation and infrastructure ownership). Our other thrust for the last year has been understanding the strength and limitations of the existing Top Level Domain (TLD) Zone files that are collected and shared on a daily basis by ICANN. We put effort into preserving an open source version of this database to serve as the basis for an implementation phase. Finally, we initiated the design of a new "rapid zone updates" platform to remove the visibility gaps of this data source, creating an audit trail for changes made to DNS and the origin of those changes.

**Work left to complete**: This dimension of the project fell the furthest behind while we focused on BGP and active measurement. We need an additional six months to prototype the rapid zone updates platform described above, and also to investigate the other types of DNS data and how to accommodate them in an implementation phase.

**Supporting infrastructure for data management.**

We made significant progress on our other tasks, including specification of storage systems, integration of infrastructure (vantage points), and test and evaluation of new dissemination approaches for traffic data. We also completed the design of a new Resource Role Portal (Figure 2) for centrally managing data requests. Feedback from the community prompted us to start two new working groups related to data management design. The first one is exploring the integration of existing CAIDA data sets with UCSD's Data Science and Machine Learning Platform (DSMLP), The second one is exploring the potential of Large Language Models for use in automatic curation and annotation of data sets, as well as extraction of metadata from scientific publications on what data sets they used. We would like to take another six months to complete these conversations and integrate the results of these efforts into our design.

**Prototyping repeatable practices to enable legitimate research access to sensitive data.**

We experimented with three repeatable practices to advance data sharing from industry. In the first, we led a multi-stakeholder working group that explored and applied a novel approach to sharing sensitive statistics about DDoS attacks, resulting in a submission to the Usenix Security conference that compares data across many different industry and academic sources. In the second, we expanded our mutual data sharing agreement with industry partner Domain Tools, where we increased the amount and scope of data we shared with each other under the existing MoU. Third, we investigated the limitations of differential privacy as a privacy-enhancing technology and submitted a paper discussing its constraints and proposing potential enhancements for managing cybersecurity data.

**Design for new data-driven Internet routing security framework.**

We completed our design of a new routing security auditing framework that relies on public BGP data to enable measurable advances against global BGP security threats, notably BGP origin and path hijacking. We also demonstrated the need for additional vantage points (peers) to improve visibility into deployment of routing security practices.

**Outreach and engagement.**

We conducted another intense week-long workshop (Oct 30 - Nov 3) with participants from academia and industry to converge on various aspects of the design project. We presented a keynote at the annual FABRIC/NRP workshop, where we summarized progress on this project and how it relates to the research infrastructure community. These meetings focused on AI/ML activities, making it clear that the next phase of NSF's AI initiatives needed to be focused on the need for data to support model training and validation. We believe our project is in an excellent position to support the community's needs for labeled data about critical cyberinfrastructure.
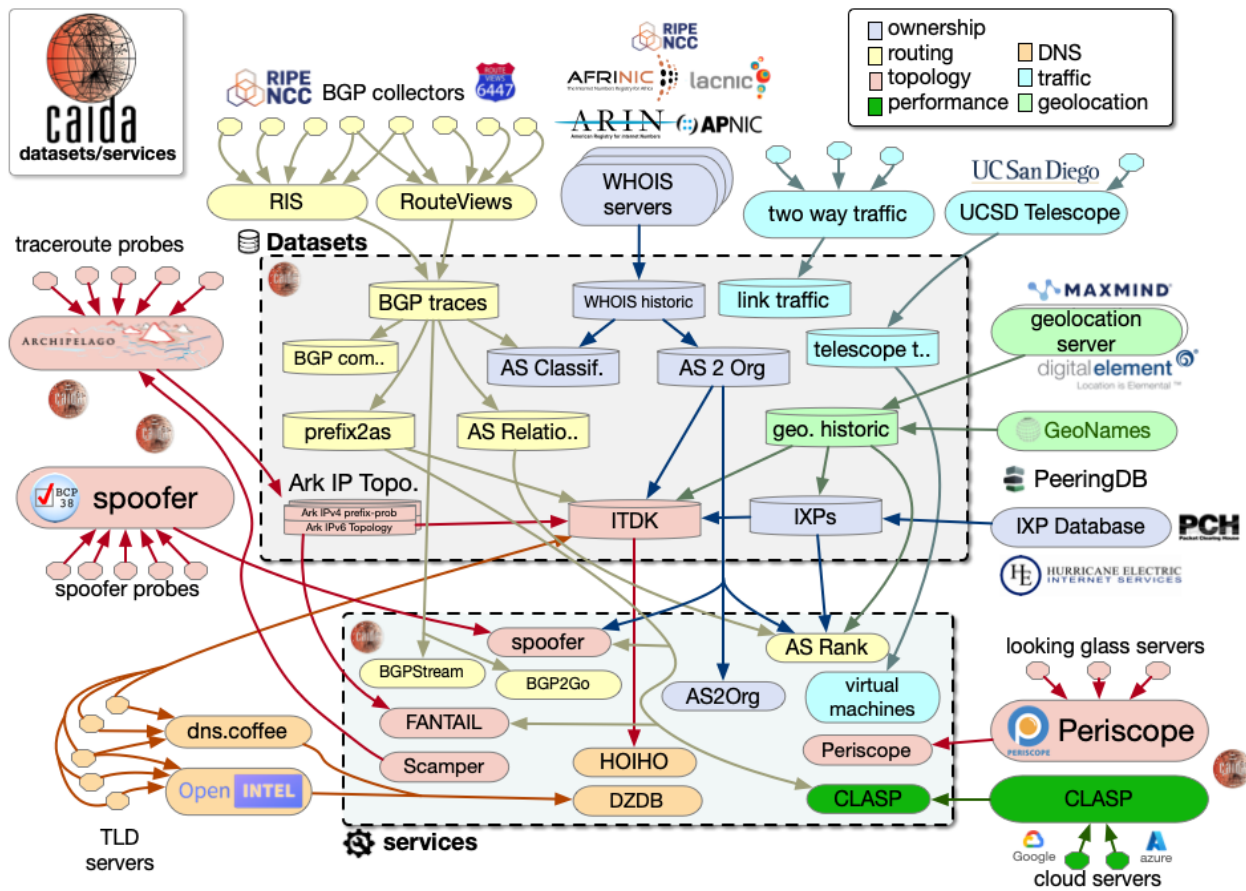
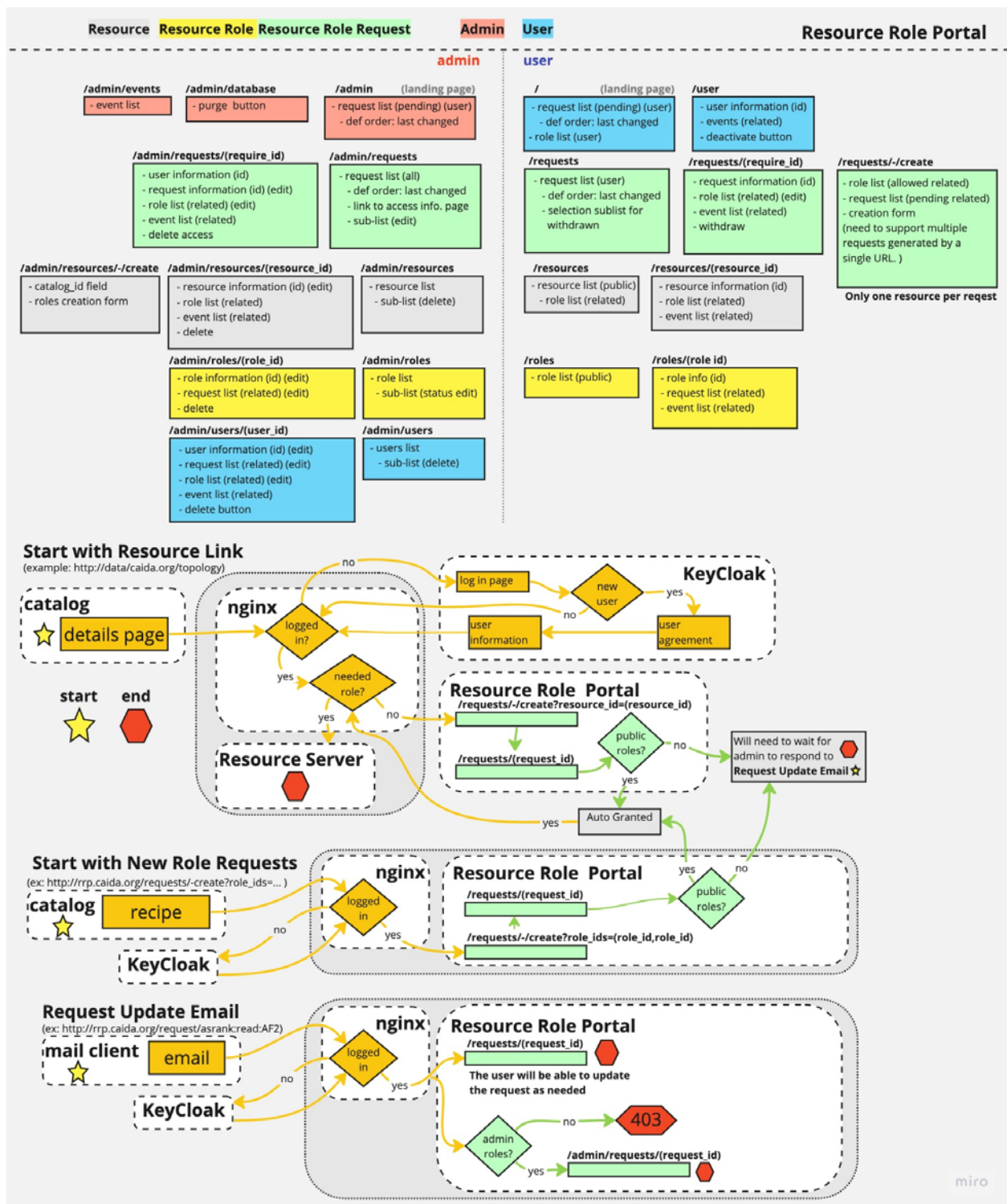Figure 1: *GMI Data acquisition components (Task 1)*

admin                     user

**/admin/events**
- event list

**/admin/database**
- purge  button

**/admin**        (landing page)
- request list (pending) (user)
- def order: last changed

**/**        (landing page)
- request list (pending) (user)
- def order: last changed
- role list (user)

**/user**
- user information (id)
- events (related)
- deactivate button

**/admin/requests/(require_id)**
- user information (id)
- request information (id) (edit)
- role list (related) (edit)
- event list (related)
- delete access

**/admin/requests**
- request list (all)
- def order: last changed
- link to access info. page
- sub-list (edit)

**/requests**
- request list (user)
- def order: last changed
- selection sublist for withdrawn

**/requests/(require_id)**
- request information (id)
- role list (related) (edit)
- event list (related)
- withdraw

**/requests/-/create**
- role list (allowed related)
- request list (pending related)
- creation form
(need to support multiple requests generated by a single URL. )

**/admin/resources/-/create**
- catalog_id field
- roles creation form

**/admin/resources/(resource_id)**
- resource information (id) (edit)
- role list (related)
- event list (related)
- delete

**/admin/resources**
- resource list
- sub-list (delete)

**/resources**
- resource list (public)
- role list (related)

**/resources/(resource_id)**
- resource information (id)
- role list (related)
- event list (related)

Only one resource per reqest

**/admin/roles/(role_id)**
- role information (id) (edit)
- request list (related) (edit)
- delete

**/admin/roles**
- role list
- sub-list (status edit)

**/roles**
- role list (public)

**/roles/(role id)**
- role info (id)
- request list (related)
- event list (related)

**/admin/users/(user_id)**
- user information (id) (edit)
- request list (related) (edit)
- role list (related) (edit)
- event list (related)
- delete button

**/admin/users**
- users list
- sub-list (delete)

**Start with Resource Link**
(example: http://data/caida.org/topology)

catalog
★ details page

KeyCloak

nginx — logged in?

log in page — new user — yes

no

user information — user agreement — user

needed role?

**Resource Role  Portal**
/requests/-/create?resource_id=(resource_id)

/requests/(request_id)

**Resource Server**

start  end

public roles? — no → Will need to wait for admin to respond to **Request Update Email ★**

yes → Auto Granted

**Start with New Role Requests**
(ex: http://rrp.caida.org/requests/-/create?role_ids=... )

catalog ★ recipe

nginx — logged in

KeyCloak

**Resource Role  Portal**
/requests/(request_id)

/requests/-/create?role_ids=(role_id,role_id)

public roles? — yes / no

**Request Update Email**
(ex: http://rrp.caida.org/request/asrank:read:AF2)

mail client ★ email

nginx — logged in

KeyCloak

**Resource Role Portal**
/requests/(request_id)

The user will be able to update the request as needed

admin roles? — no → 403

yes → /admin/requests/(request_id)

miro

Figure 2: *Diagram for the workflow in the Resource Access Management Portal*