# NSDF-Catalog: Toward a Lightweight Indexing Service for the National Science Data Fabric

Jakob Luettgau, Giorgio Scorzelli, Naweiluo Zhou, Glenn Tarcea, Jay Lofstead, Christine Kirkpatrick, Valerio Pascucci, Michela Taufer

PI/Co-PIs:
Valerio Pascucci (U Utah)
**Michela  Taufer (UTK)**
Alexander  Szalay (JHU)
John  Allison (U. Michigan)
Frank  Wuerthwein (UCSD / SDSC)

**http://nationalsciencedatafabric.org/**

We are building a trans-disciplinary **testbed** that will **democratize data-driven scientific discovery** by **connecting an open network of institutions**, including minority serving institutions and with **a shared, modular, containerized data delivery environment**.
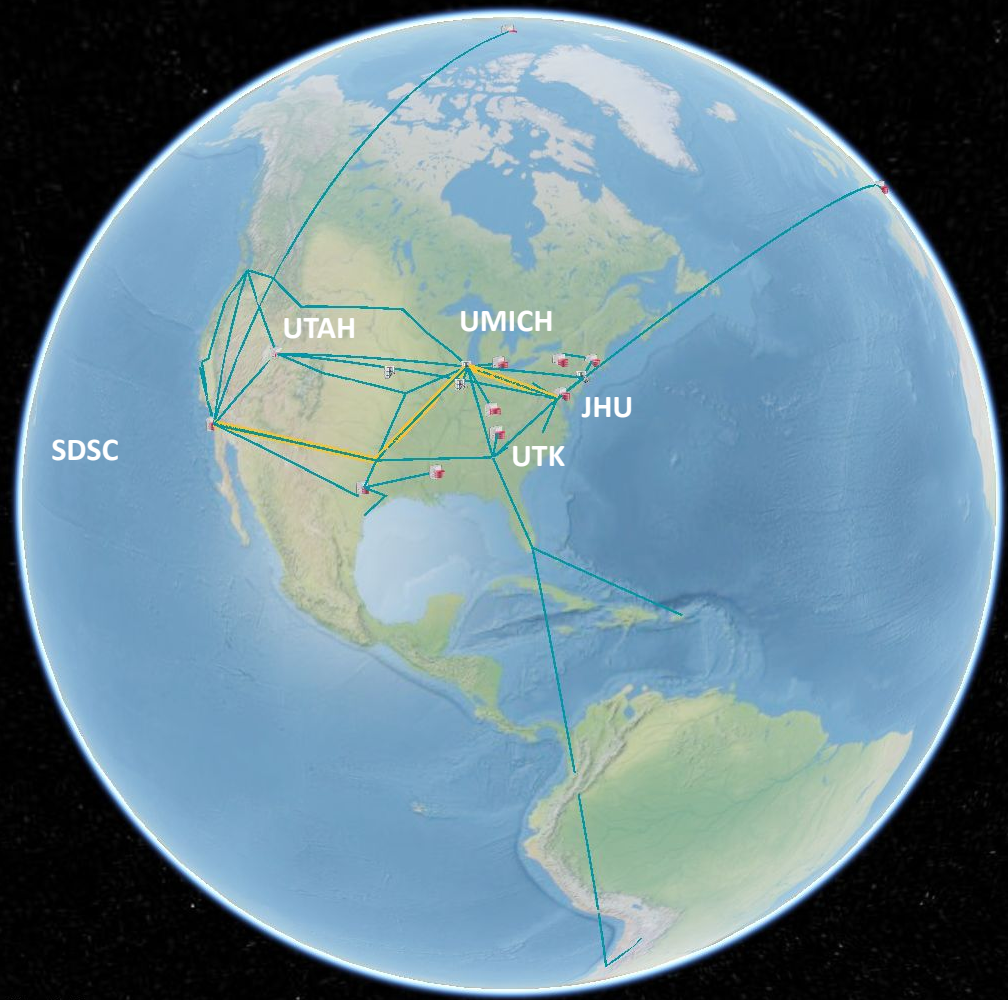
We are building a trans-disciplinary **testbed** that will **democratize data-driven scientific discovery** by **connecting an open network of institutions**, including minority serving institutions and with **a shared, modular, containerized data delivery environment**.

→ **Need for global index of data available within National Science Data Fabric**

100G Core
Terabit Core

NSDF EntryPoints

OSG StashCaches

Both

UTAH

UMICH

JHU

SDSC

UTK

Today
Aug 8 2022
19:51:04 UTC

Aug 9 2022 00:00:00 UTC    Aug 9 2022 04:00:00 UTC    Aug 9 2022 08:00:00 UTC    Aug 9 2022 12:00:00 UTC    Aug 9 2022 16:00:00 UTC    Aug 9 202

2 20:00:00 UTC

Full screen

IceCube

XenonNT

PRISMS
Materials
Commons
UMICH

CHESS/Cornell
MGHPCC + OSN

UTAH/SCI

JHU

SDSC + OSG

UTK
MS-CC
Digital Rocks

100G Core
Terabit Core

NSDF EntryPoints

OSG StashCaches

Both

Today
Aug 8 2022
19:51:04 UTC

CESIUM ion    Upgrade for commercial use.  Data attribution

Full screen

Aug 9 2022 00:00:00 UTC    Aug 9 2022 04:00:00 UTC    Aug 9 2022 08:00:00 UTC    Aug 9 2022 12:00:00 UTC    Aug 9 2022 16:00:00 UTC

**Entry Point**

**Entry Point**

NSDF Entry Point

User Layer
Command Line Tools (CLI)
Interactive Notebooks

**Scientists Educators Students**

**Application Developers**

Data Layer
Tier 3 — Workflows and Automation
Tier 2 — Data Discovery | Data Curation | Data Processing | Data Analytics | Data Mapping | Visualization
Tier 1 — Data Management | Co...

**Middleware/Service Developers**

Extensible Content Delivery Network
SDK
APIs
Microservices
CDN Kernel and Plug-Ins

Support Services
Security
Data Meta-Catalog
Logging and monitoring
Provenance
Containers and Orchestration

Commodity Appliance

Local, Additional Existing Storage and Computing Resources
Domain Agnostic Software Stack

**Infrastructure Providers**

Regional Ceph Storage
PRP DTN Node (FIONA)
OSN STORAGE POD
VDC HUB
OSG StashCache Node
OSG Data Federation
National CI
Large Facilities and Repositories
Edge Data Streams

HPC and Other Remote File System Storage
Public Cloud and Other Object Storage
OSG Origin Node

**Entry Point**

NSDF Entry Point

User Layer
- Command Line Tools (CLI)
- Interactive Notebooks

**Application Developers**

**Scientists Educators Students**

Support Services
- Security
- Data Meta-Catalog
- Logging and monitoring
- Provenance
- Containers and Orchestration

Data Layer
- Tier 3: Workflows and Automation
- Tier 2: Data Discovery | Data Curation | Data Processing | Data Analytics | Data Mapping | Visualization
- Tier 1: Data Management

**Middleware/Service Developers**

Domain Agnostic Software Stack

Local, Additional Existing Storage and Computing Resources

Extensible Content Delivery Network
- SDK
- APIs
- Microservices
- CDN Kernel and Plug-Ins

Commodity Appliance

**Infrastructure Providers**

- Regional Ceph Storage
- PRP DTN Node (FIONA)
- OSN STORAGE POD
- VDC HUB
- OSG StashCache Node
- OSG Data Federation
- National CI
- Large Facilities and Repositories
- Edge Data Streams
- OSG Origin Node

**HPC**

**Commercial Clouds**

**Science Clouds/Grids**

HPC and Other Remote File System Storage

13

# Services & Building Blocks



Data transfer services. Like Globus but object-oriented. Set-the-transfer and forget it using the NSDF infrastructure
**CHESS Thinkmate rack 190TiB**
**UMIch Supermicro Rack 144TiB**

**NSDF-dts**

Mount object-storage as a file system. Simplify data access to computer scientists with a friendly almost-POSIX file system
**short paper - GITHUB**

**NSDF-fuse**

Create mini cluster (globally distributed and hybryd) on-the-fly. Supporting either educational clouds (xsede,chamaleon,cyverse) and commercial clouds (AWS,Google, Azure etc)
**short paper - GitHub**

**NSDF-cloud**

**NSDF-plugin**

Globally distributed key-value store for metadata storage. Fast (milliseconds) access. Supports unique global namespace with multi-regional UID. **paper?**

Add repositories to our federation (Material Commons, Digital Rocks, Chess, NIST etc)

**NSDF-catalog**

**NSDF-monitor**

Track data movements and data accesses. View them. First prototype on UCSD infrastructure

Services for parallel and fast conversion from scientific file formats (netCDF, HD5 etc) to Analysis ReadtyCloud Optimized (ARCO) File Format

**NSDF-workflow**

**NSDF platform**

**NSDF-website**

and socials, and webinar

Streaming and visualization services to access the data remotely. Coarse-to-fine

**NSDF-stream**

**NSDF-edu**

Educational materials (e.g. Jupyter Notebooks, Playground etc)

material commons to catalog all openvisus DB with jupyter notebook?

**NSDF-openvisus-commons**

Support for the Inter Planetary file system for peer-to-peer distributed storage. More oriented to never-loose-your-data / cloud-storage- Glacier like accesses

**NSDF-ipfs**

User Communities

Commercial Partners

**NSDF-cdn**

Content delivery network with multi-layered caching and object distribution. Monitor data flows, push-pull policies, cache evictions, dashboards etc.

14

# Services & Building Blocks



**Transfer Abstraction** with, e.g,, Globus or xrootd

**Mapping object storage into POSIX namespaces for legacy support**

**Data Inventory and Discovery**

**Common APIs across cloud providers**

**System Health & Optimization**

NSDF platform

**NSDF-dts**

NSDF-fuse

**NSDF-cloud**

*clouds (AWS,Google, Azure etc)* **short paper - GitHub**

**NSDF-plugin**

*Globally distributed key-value store for metadata storage. Fast (milliseconds) access. Supports unique global namespace with multi-regional UID.* **paper?**

**NSDF-catalog**

*(Material Commons, Digital Rocks, Chess, NIST etc)*

**NSDF-workflow**

*Services for parallel and fast conversion from scientific file formats (netCDF, HD5 etc) to Analysis ReadtyCloud Optimized (ARCO) File Format*

**NSDF-stream**

*Streaming and visualization services to access the data remotely. Coarse-to-fine*

**NSDF-openvisus-commons**

*material commons to catalog all openvisus DB with jupyter notebook?*

**NSDF-monitor**

*accesses. View them. First prototype on UCSD infrastructure*

**NSDF-website**

*and socials, and webinar*

**NSDF-edu**

*Educational materials (e.g. Jupyter Notebooks, Playground etc)*

**NSDF-ipfs**

*Support for the Inter Planetary file system for peer-to-peer distributed storage. More oriented to never-loose-your-data / cloud-storage- Glacier like accesses*

**NSDF-cdn**

*Content delivery network with multi-layered caching and object distribution. Monitor data flows, push-pull policies, cache evictions, dashboards etc.*

User Communities

Commercial Partners

*Data transfer services. Like Globus but object oriented. Set... ...mputer scien... ...st-PO... ...rt pap...*

*Mount object ... ...system. Simp...*

15
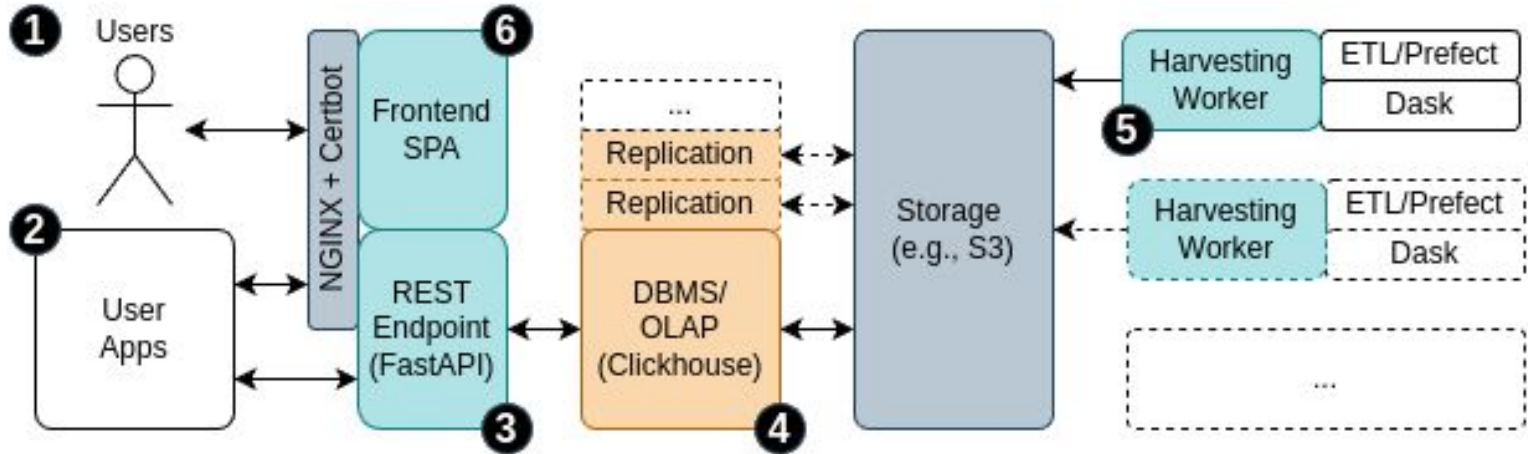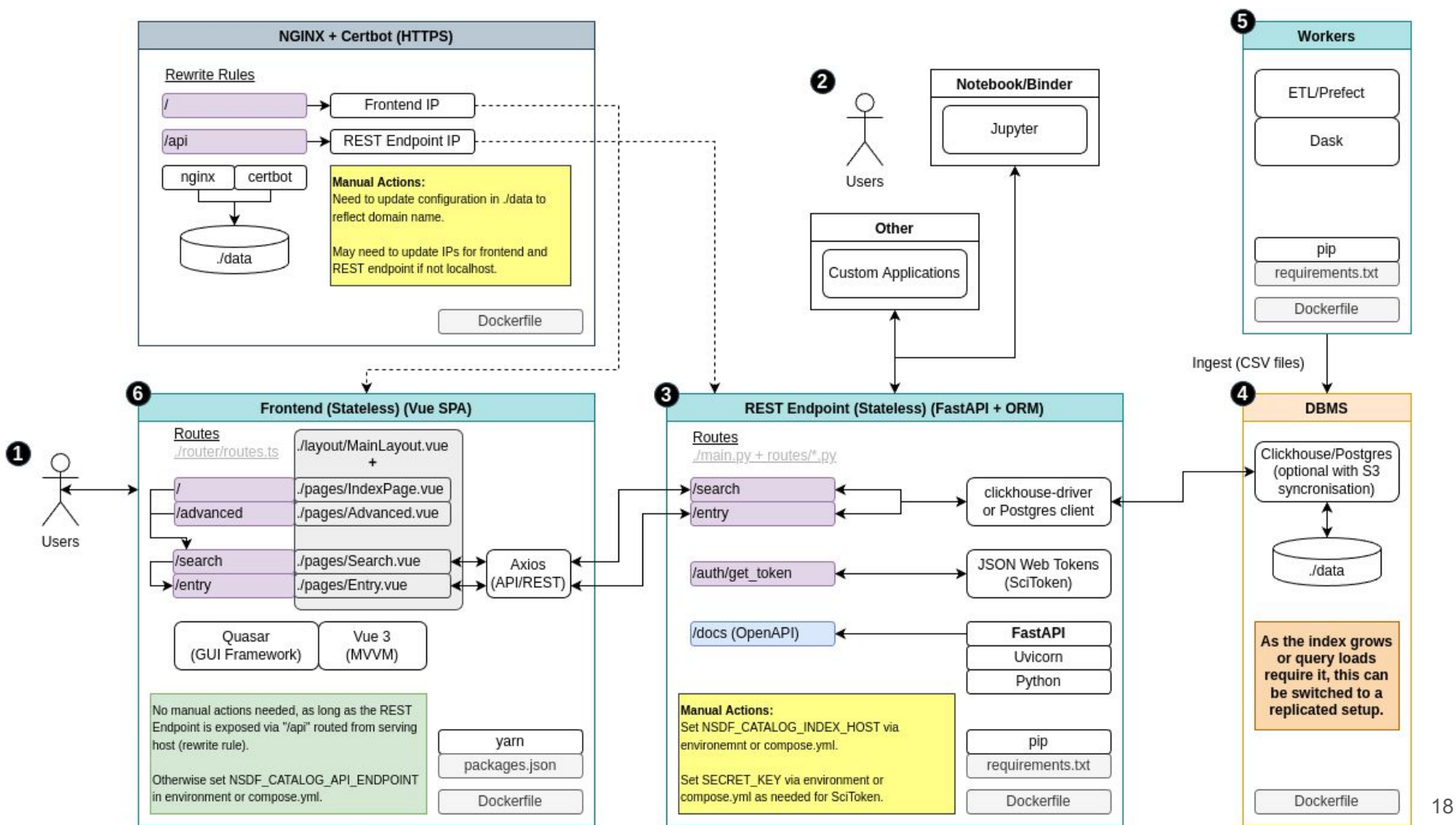
# NSDF-Catalog: Challenges

- Support wide range of existing repositories

- Catalog must be scalable to many entries with a pathway to trillions of datasets, files, or objects

- Containerized to easy maintenance and scaling

- Catalog must be federated for efficient indexing and to offer user/provider control over their data
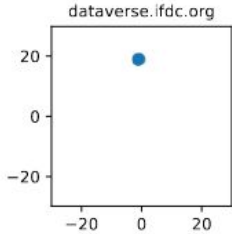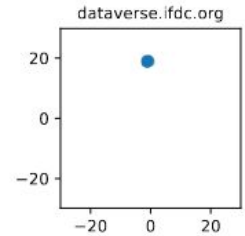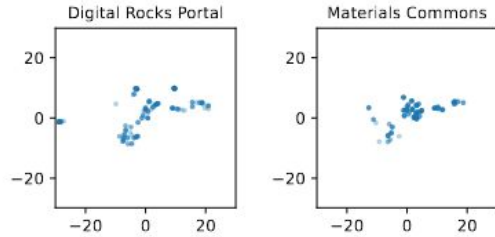
# NSDF-Catalog: Architecture Overview

**NGINX + Certbot (HTTPS)**

Rewrite Rules

/ → Frontend IP

/api → REST Endpoint IP

nginx | certbot → ./data

**Manual Actions:**
Need to update configuration in ./data to reflect domain name.

May need to update IPs for frontend and REST endpoint if not localhost.

Dockerfile

❷ Users

**Notebook/Binder**

Jupyter

**Other**

Custom Applications

❺ **Workers**

ETL/Prefect

Dask

pip
requirements.txt

Dockerfile

Ingest (CSV files)

❻ **Frontend (Stateless) (Vue SPA)**

Routes
./router/routes.ts

./layout/MainLayout.vue
+

/ → ./pages/IndexPage.vue

/advanced → ./pages/Advanced.vue

/search → ./pages/Search.vue

/entry → ./pages/Entry.vue

Axios
(API/REST)

Quasar
(GUI Framework)

Vue 3
(MVVM)

No manual actions needed, as long as the REST Endpoint is exposed via "/api" routed from serving host (rewrite rule).

Otherwise set NSDF_CATALOG_API_ENDPOINT in environment or compose.yml.

yarn
packages.json

Dockerfile

❶ Users

❸ **REST Endpoint (Stateless) (FastAPI + ORM)**

Routes
./main.py + routes/*.py

/search

/entry

/auth/get_token

/docs (OpenAPI)

clickhouse-driver
or Postgres client

JSON Web Tokens
(SciToken)

**FastAPI**
Uvicorn
Python

**Manual Actions:**
Set NSDF_CATALOG_INDEX_HOST via environemnt or compose.yml.

Set SECRET_KEY via environment or compose.yml as needed for SciToken.

pip
requirements.txt

Dockerfile

❹ **DBMS**

Clickhouse/Postgres
(optional with S3 syncronisation)

./data

**As the index grows or query loads require it, this can be switched to a replicated setup.**

Dockerfile

18

# NSDF-Catalog: Challenges

| Repository | # Collections | # Entries | Size (Bytes) |
|---|---|---|---|
| Digital Rocks Portal | 148 | 17,285 | 6.1 TiB |
| Materials Commons | 70 | 258,576 | 10.2 TiB |
| Materials Data Facility | 178 | 1,075,706 | 4.8 TiB |
| Arecibo Observatory | 221 | 2,045,049 | 447.4 TiB |
| AWS Open Data | 397 | 1,617,966,022 | 50,400.0 TiB |
| TACC/Ranch | 184 | 1,091,321 | 20,500.0 TiB |
| zenodo.org | 1,001,459 | 3,461,517 | 339.5 TiB |
| Dataverse | 154,472 | 2,306,495 | 104.9 TiB |

# NSDF-Catalog: Finding Similarities Across Research Repositories



dataverse.ifdc.org

# NSDF-Catalog: Finding Similarities Across Research Repositories



Digital Rocks Portal

Materials Commons

dataverse.ifdc.org

# NSDF-Catalog: Finding Similarities Across Research Repositories

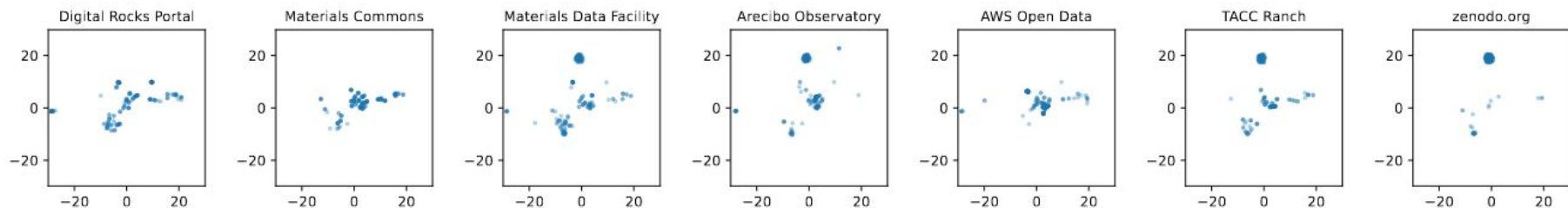# NSDF-Catalog: Finding Similarities Across Research Repositories
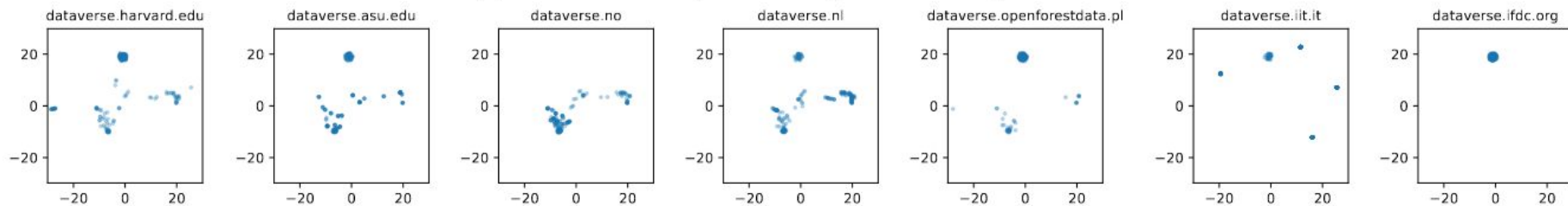


(a) Set of domain-specific and general data repositories.

(b) Selected Dataverse repositories.

# NSDF-Catalog: Finding Similarities Across Research Repositories



(a) Set of domain-specific and general data repositories.

(b) Selected Dataverse repositories.

# Summary

- Building a lightweight index for large amounts of scientific data is feasible spread across a variety of different existing repositories

- There is structure across catalogs that can be leveraged to improve search and also to optimize performance for the National Science Data Fabric

Outlook:
- We looking to collaborate with science teams that are sharing their or depend on other data to better understand needs for NSDF-Catalog

- We are looking for scientists performing cross-disciplinary research that would leverage our search and are willing to discuss their use case