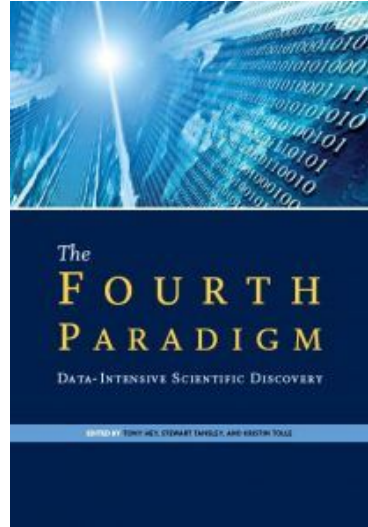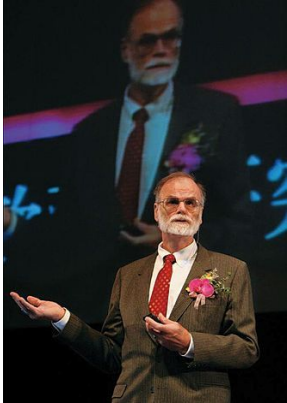SciServer

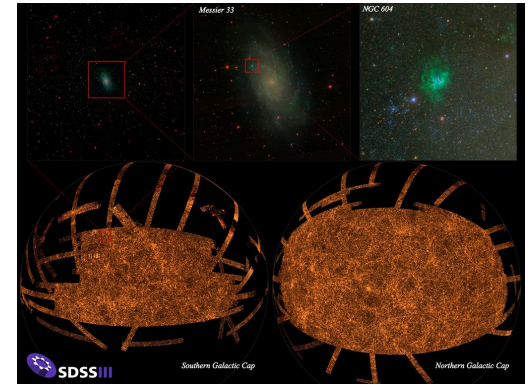A Collaborative Science Platform developed by IDIES

# Outline

- Background
- SciServer overview
- Use case highlights
- Architecture & Deployments
- New/hidden features/accelerated computing
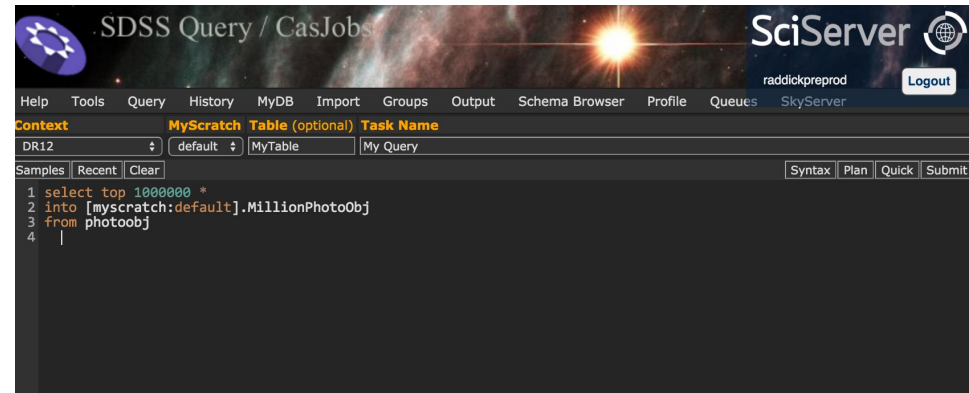- Future directions
- Questions

# Beginnings







The fourth paradigm: **Data exploration**.

More data being produced than can be analyzed, increasing dependence on tools and machine learning to help make sense.

# SDSS / CasJobs

- SQL interface to large catalog of object data
- Server-side execution, download small result-set of interest
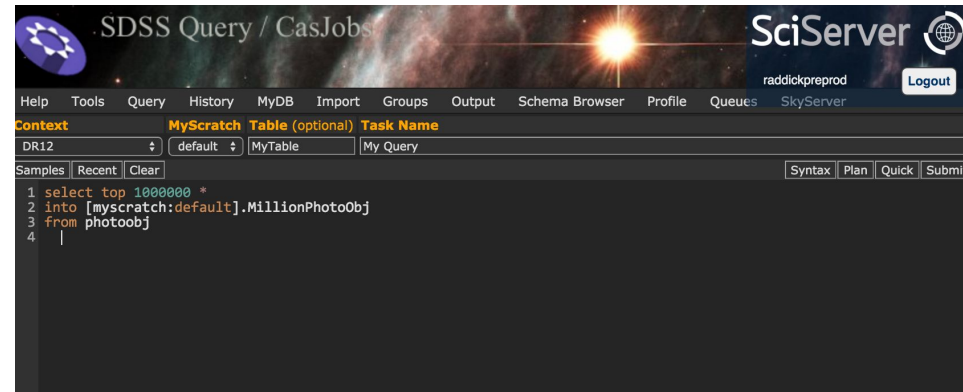- Personal DB space
- Sharing features

# SDSS / CasJobs

- SQL interface to large catalog of object data
- Server-side execution, download small result-set of interest
- Personal DB space
- Sharing features

- Still need to download
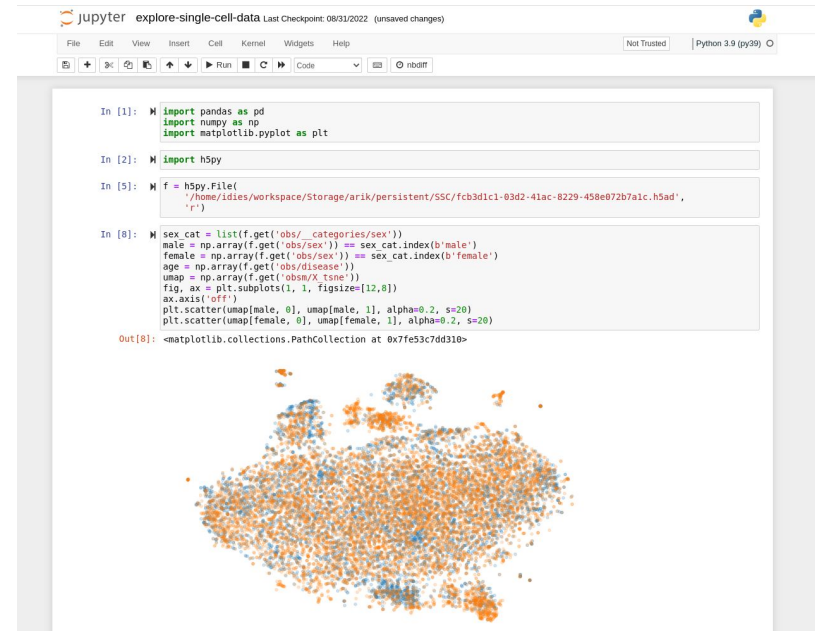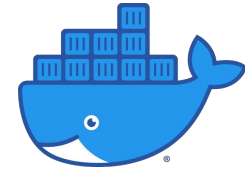- No solution for file-based data
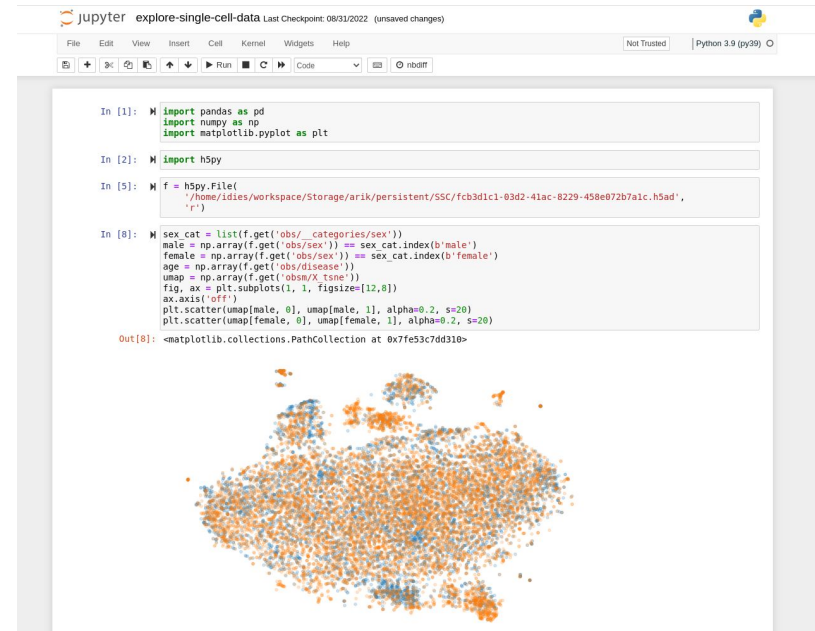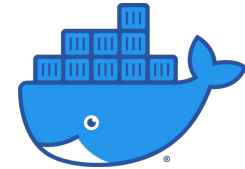- limited/no visualization support

# Beyond SQL

- Add Jupyter
- Docker makes scheduling simple
- Jupyter can run arbitrary codes
- Inline visualization
- Common dependencies included
- SDK for CasJobs (get a pandas DF)
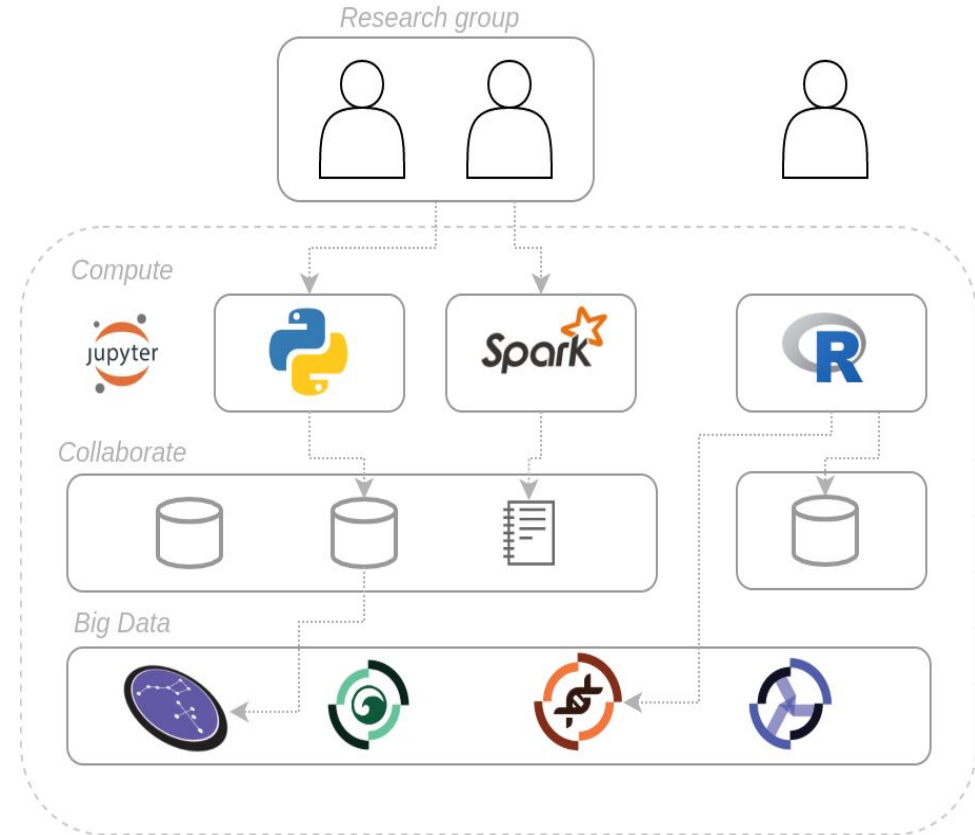
# Beyond SQL

- Add Jupyter
- Docker makes scheduling simple
- Jupyter can run arbitrary codes
- Inline visualization
- Common dependencies included
- SDK for CasJobs (get a pandas DF)


- Server-side collaboration?
- Expanding scope?

# SciServer

- Add in an ecosystem
- Sharing/resource access control
- Web based portals
- Curated data
  - Astronomy
  - Genomics
  - Earth Science
  - Simulation data,...

*A Science Platform*

# Basic Usage

- Go to [https://apps.sciserver.org](https://apps.sciserver.org)
- Create an account
  - You can sign in with your institution via globus, your google account or orcid by using "sign in with Globus"

# Basic Usage

- Arrive at dashboard
- A number of links to services
  - Files
  - Groups
  - Compute
  - Casjobs

# Basic Usage - Groups

- Create groups
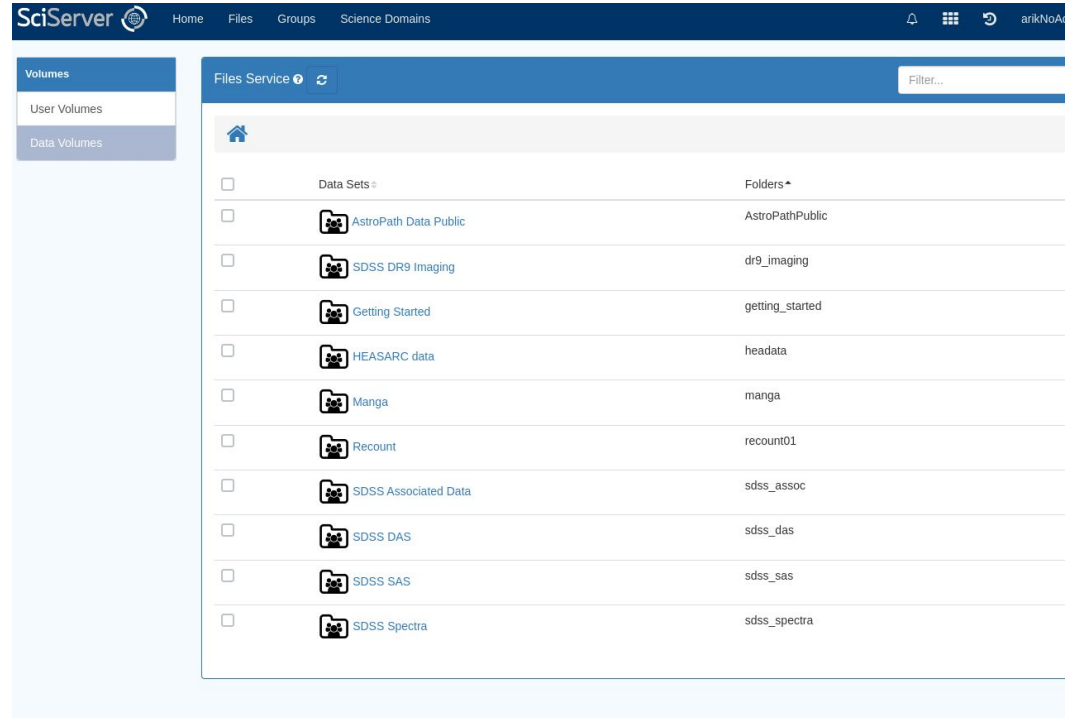- Groups are a unit of sharing
- Invite other users

# Basic Usage - User Volumes

- Create "User Volumes"
- Volumes can be shared
- Store your own data
- Default limited space
- Temporary space:
  - For intermediate results
  - "Unlimited"
  - time-to-live

# Basic Usage - Data Volumes

- "Data Volumes"
- Read-only
- Curated data (typically)
- Shared by "admins"

# Basic Usage - Compute

- Create a "container"
- Backed by docker
- Web service listening at specific port (Jupyter)
- Mount volumes needed

# Basic Usage - Compute

- Do analysis
- Volumes are mounted
  - Posix access to regular files
- Lots of software already included
- Token "injected" for API access to SciServer and CasJobs

# Basic Usage - Jobs

- Asynchronous Jobs
- Time-limited
- Large/specialized resources
- Queued, fair scheduling
- Clearer provenance

# Use Case Highlights

# Use Case: Turbulence

- Get cutouts of a large-scale turbulence simulation (~½ PB)
- Web interface, actions backed by SciServer jobs
- Use notebooks to analyze results

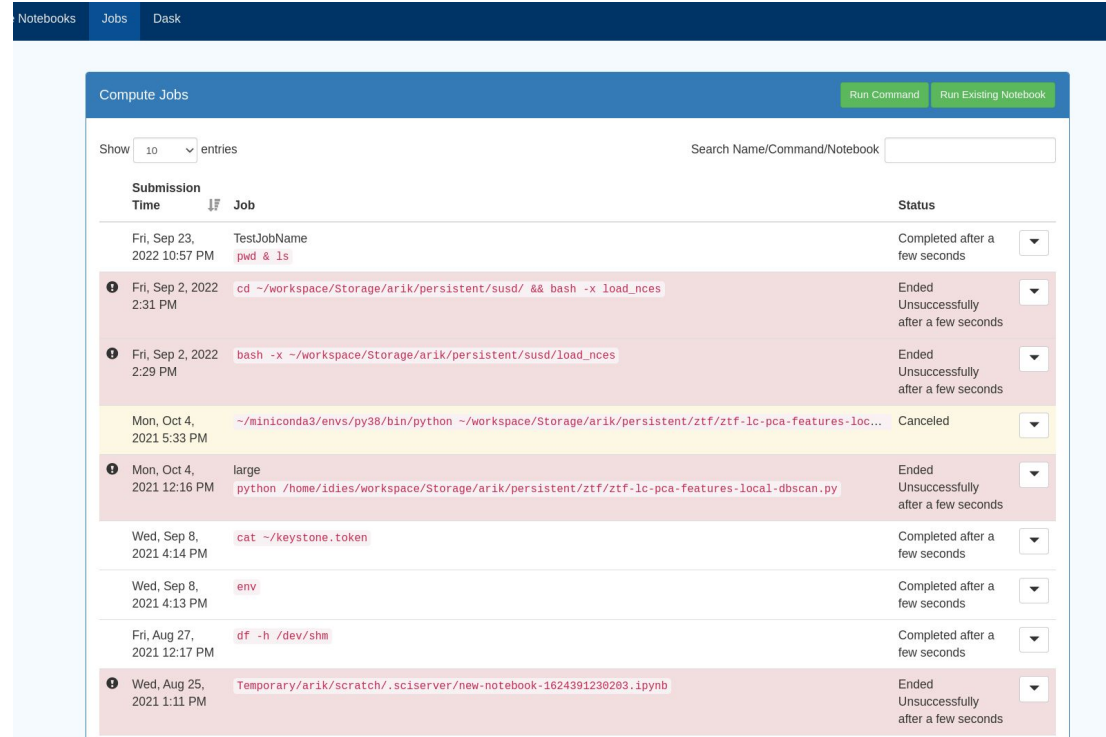http://turbulence.pha.jhu.edu/



SciUI - JHU Turbulence DB - Cutout Service

Log in



```
In [2]:    import pandas as pd
           import numpy as np
           import matplotlib.pyplot as plt
           import h5py
```

```
n [21]:    f = h5py.File('/home/idies/workspace/Temporary/arik/scratch/jobs/__turbcutout__/20221013/20221013120636-195226/channel.h5', 'r')
```

```
n [27]:    fig, a = plt.subplots(1, 10, figsize=[16,4])
           for i in range(10):
               plt.sca(a[i])
               plt.imshow(f['Velocity_0001'][:,:,i,0])
```
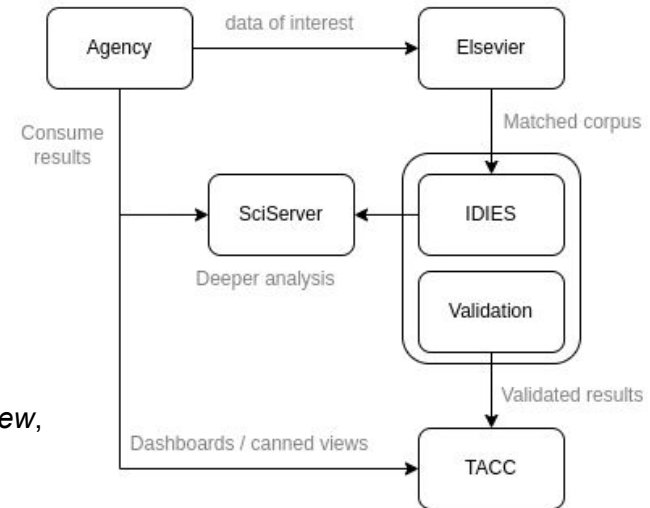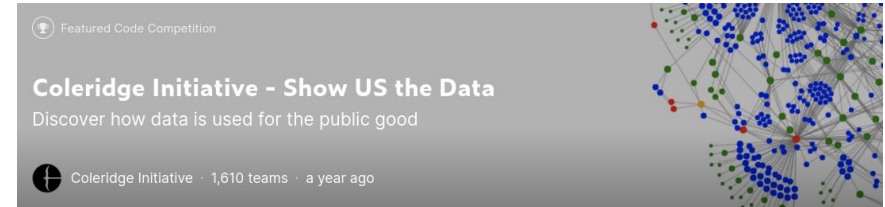
In [ ]:

# Use Case: Democratizing Data (Show Us The Data)

- Context: track dataset usage for gov agencies
- Rich data model, SQL backed
- Simple APIs for Tableu (hosted at TACC), but…
- SciServer enables deeper dive for those interested



See aso: Potok, N. (2022). Show US the Data. *Harvard Data Science Review*, *4*(2). https://doi.org/10.1162/99608f92.9d13ba15

# SciServer as a service

# Architecture

# Deployments

- All the same features:
  - In another datacenter on own hardware
  - Behind a VPN
  - For specific insiders
  - In a HIPAA-safe environment
- A number of partners
- Active development

@JHU (sciserver.org)

*On premise*

@NIST

*AWS VPC, k8s*

@NAOJ

*On premise, k8s*

@MPE

*On premise, k8s*

@Precision-Medicine

*Azure, secure-desktop*

# Deployments

- Branding
  - Theme customizations
  - links
  - etc

# Deployments

- Leverage Kubernetes and helm
- Bundle SciServer + 3rd party dependencies
- Configurable with a single file

```
$ helm install sciserver -f config.yaml sciserver-k8s/charts/sciserver
```

# Extra Features

# Extra features

- Accelerated computing
  - GPU
  - Spark
  - Dask
- Interactive and/or Jobs

# Extra Features

- Not Just Jupyter
- Compute exposes service that listens to single port
- VNC, for example

# Extra Features
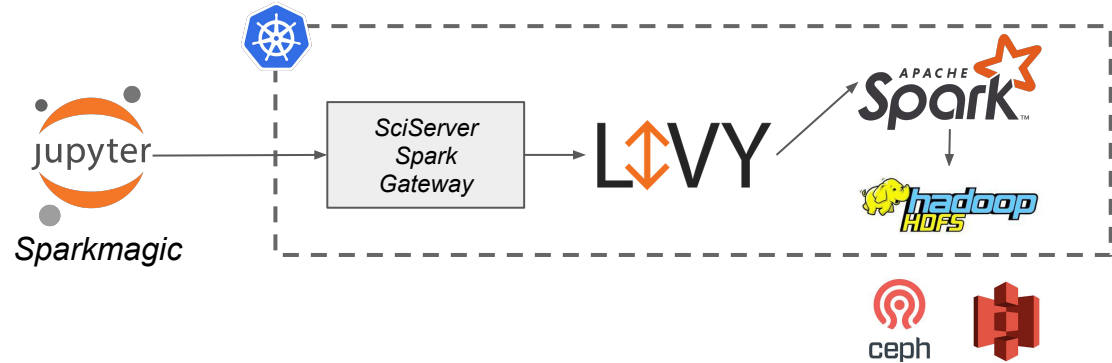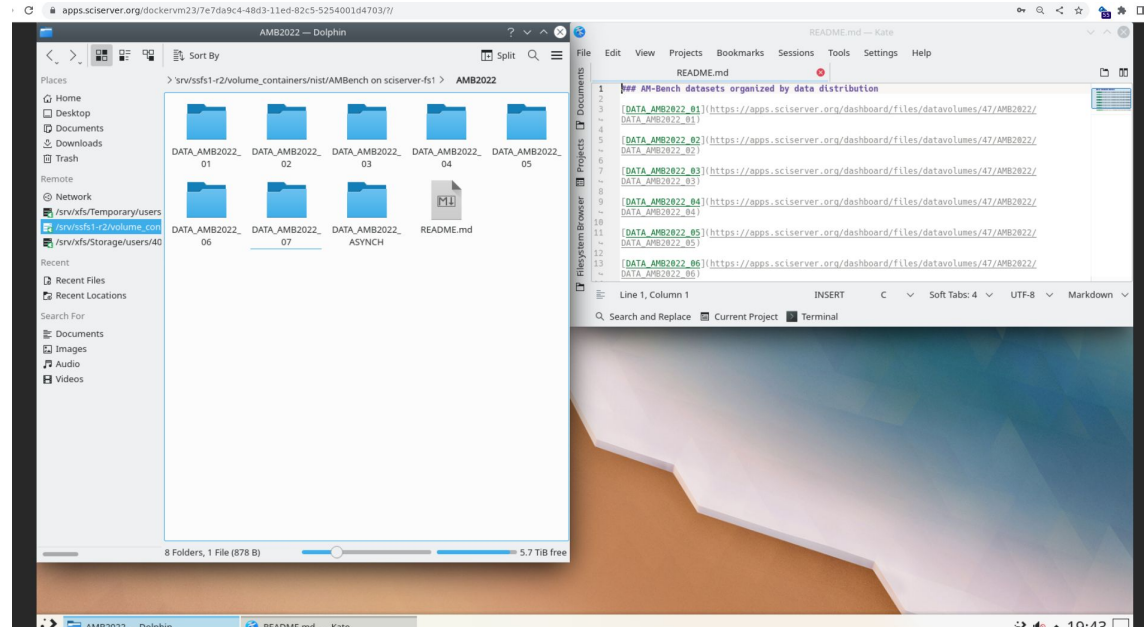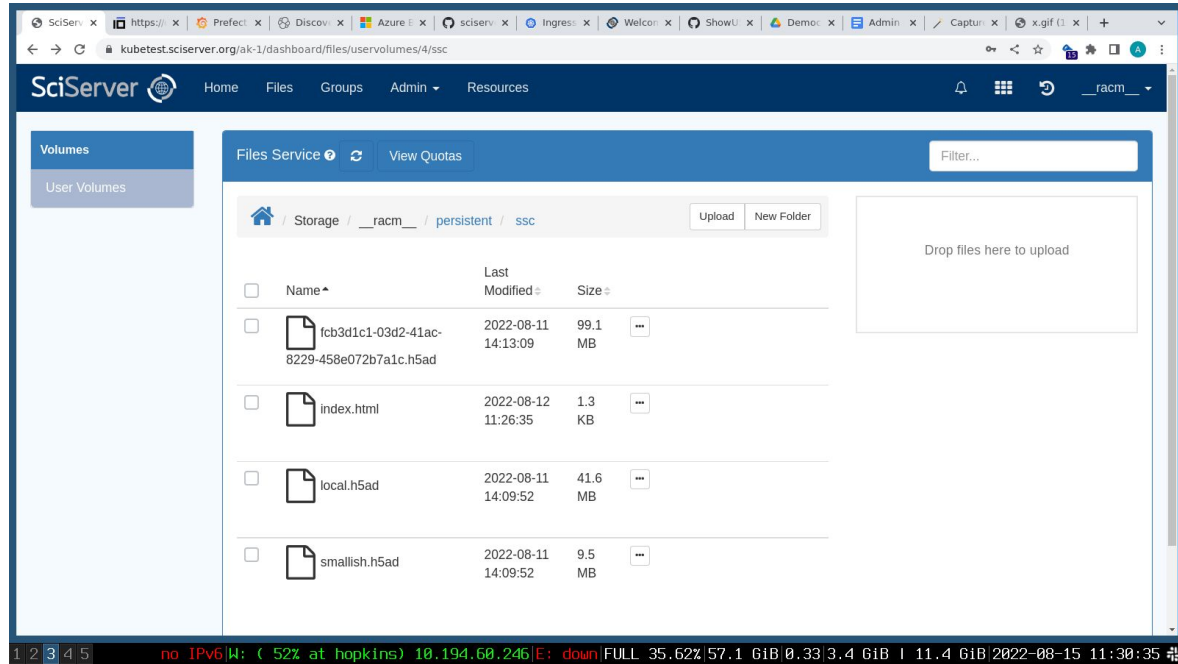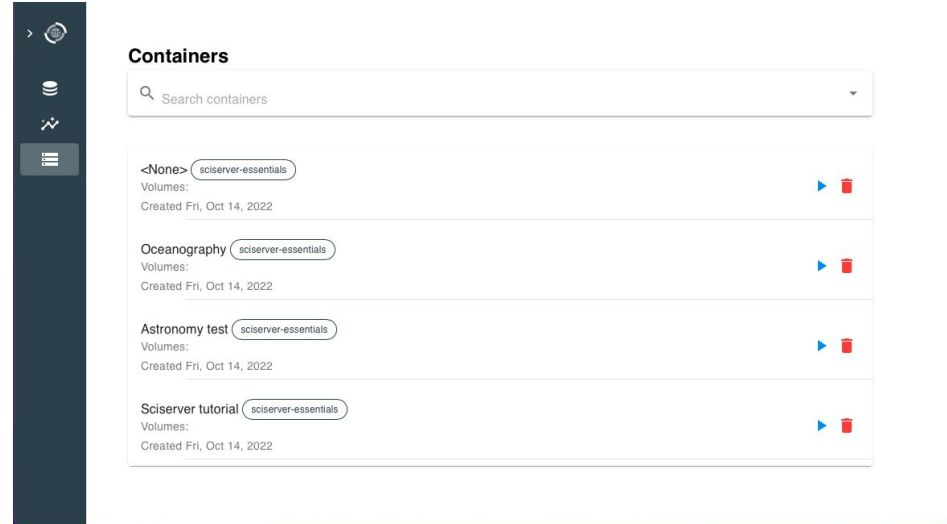
- "One-click apps"
- Shareable links
- Argument passing
  - Can do in one-click what would take steps using generic compute containers
- Operate on files
- Embedded feel

# Future Directions

# Future Directions

- UI consolidation / improvements
  - Simplify development
  - Simplify branding and customizations
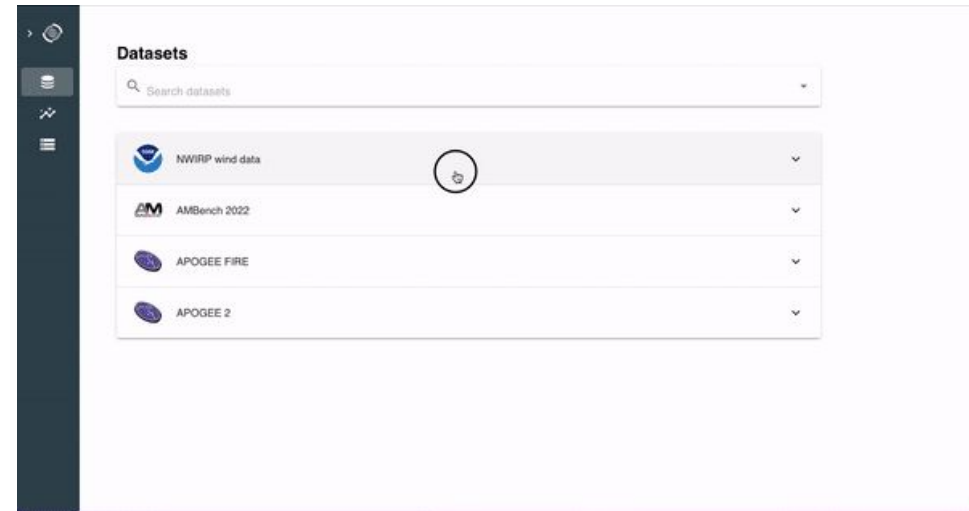  - Improved layout and features

# Future Directions

- UI consolidation / improvements
- Data discoverability
  - "Documents close to the data"
  - Leverage built-in access control
  - Human and machine writable
  - Just documentation, not full ontology

# Future Directions

- UI consolidation / improvements
- Data discoverability
- One-click simplicity / link-share
  - Combine with data discoverability
  - Notebooks, quicker start
  - Jupyter widgets based apps
  - Sends links, not wordy instructions

# Future Directions

- UI consolidation / improvements
- Data discoverability
- One-click simplicity / link-share
- Cloud integrations
- Others on the radar
  - Extended login sessions / api keys
  - Scheduled / cron-like jobs
  - Anonymous file sharing

# Thanks!