

# SciServer



Collaborative data-driven science

## SciServer Compute

Mike Rippin



JOHNS HOPKINS  
UNIVERSITY



idies

# Agenda

- ▶ Background and History of SciServer
- ▶ Major Objectives
- ▶ Current System
- ▶ SciServer Compute – Now
- ▶ SciServer Compute – Future
- ▶ Q&A

# SciServer



Collaborative data-driven science

## Background



“The Project aims to create a sustainable collaborative ecosystem built around several large scientific data sets for the broader science community, based upon the expertise developed for the Sloan Digital Sky Survey (SDSS) SkyServer and associated projects.”

# Project Management

- ▶ NSF Cooperative Agreement
- ▶ 5 years duration, just completed first 3
- ▶ Development of Cyberinfrastructure
- ▶ Science Driven

# Motivation and History

- ▶ Started with the SDSS SkyServer
- ▶ Goal: instant access to rich content
- ▶ Idea: bring the analysis to the data
- ▶ Interactive access at the core

# Where Are We Going?

- ▶ Interactive science on petascale data
- ▶ Create scalable open numerical laboratories
- ▶ Large footprint across many disciplines
- ▶ Use commonly shared building blocks
- ▶ Major national and international impact

# SciServer



Collaborative data-driven science

## Current System





# SciServer : Core Functions

Cyber  
Infrastructure

Science  
Collaboration

SDSS  
Integration

Outreach &  
Education

# SciServer : Core Functions

Cyber  
Infrastructure

Science  
Collaboration

SDSS  
Integration

Outreach &  
Education

## Main components:

Database storage & Query:

**CasJobs**



Data analysis:

**Compute**



File storage:

**SciDrive**



Data exploration:

**SkyServer**

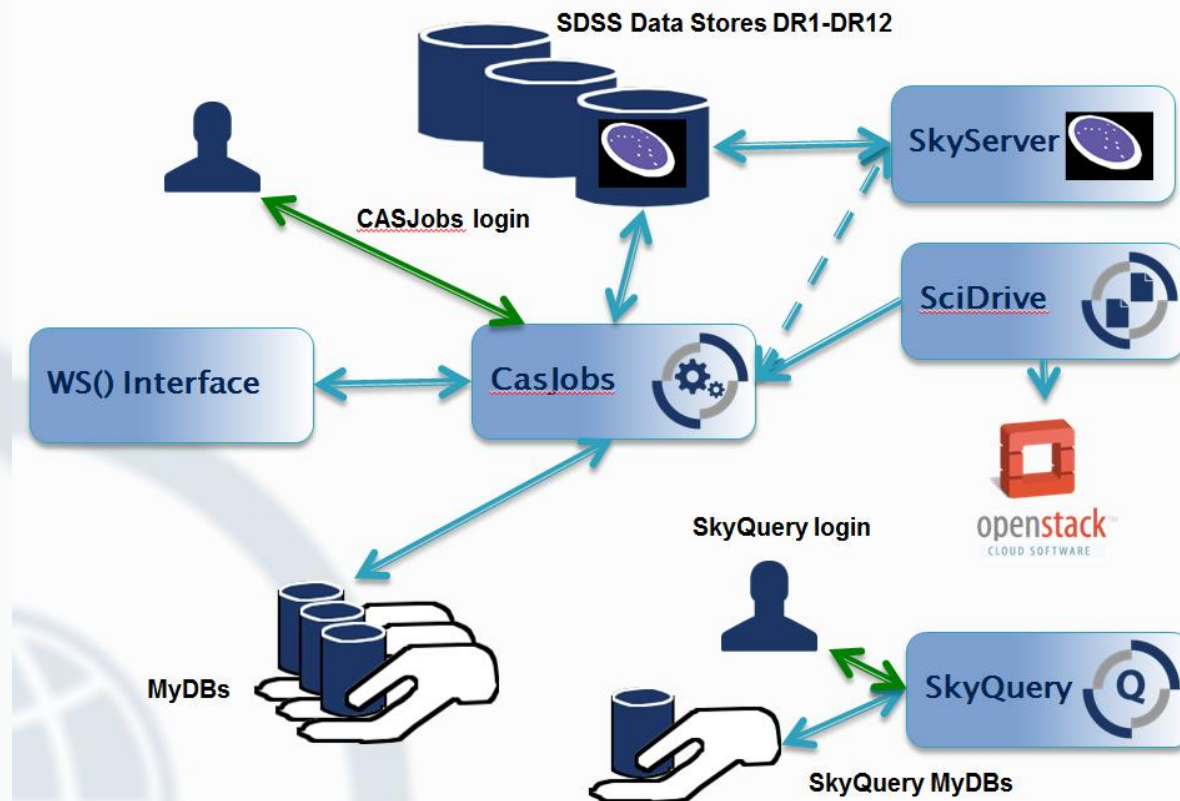


User sign-on:

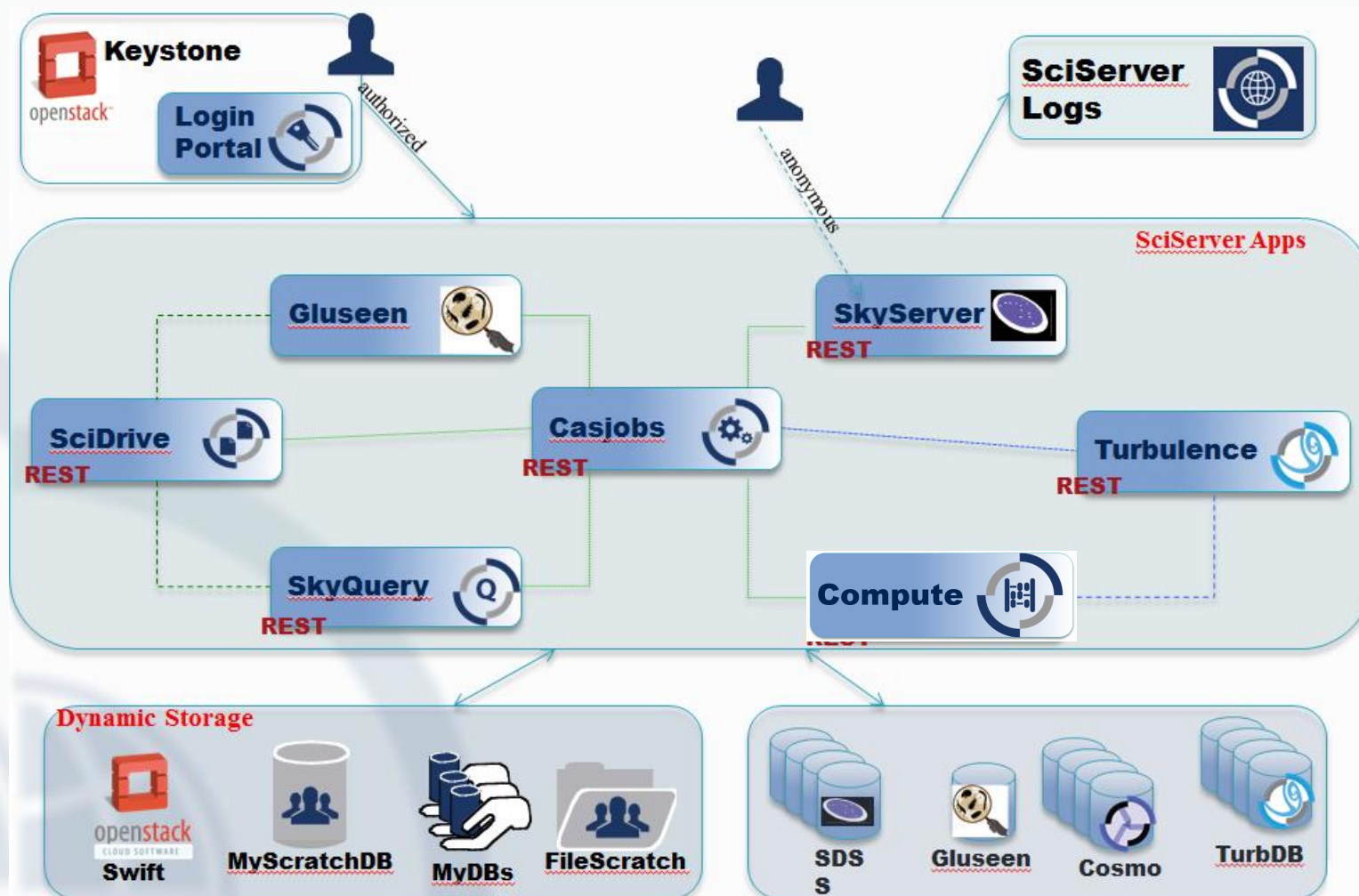
**Portal**



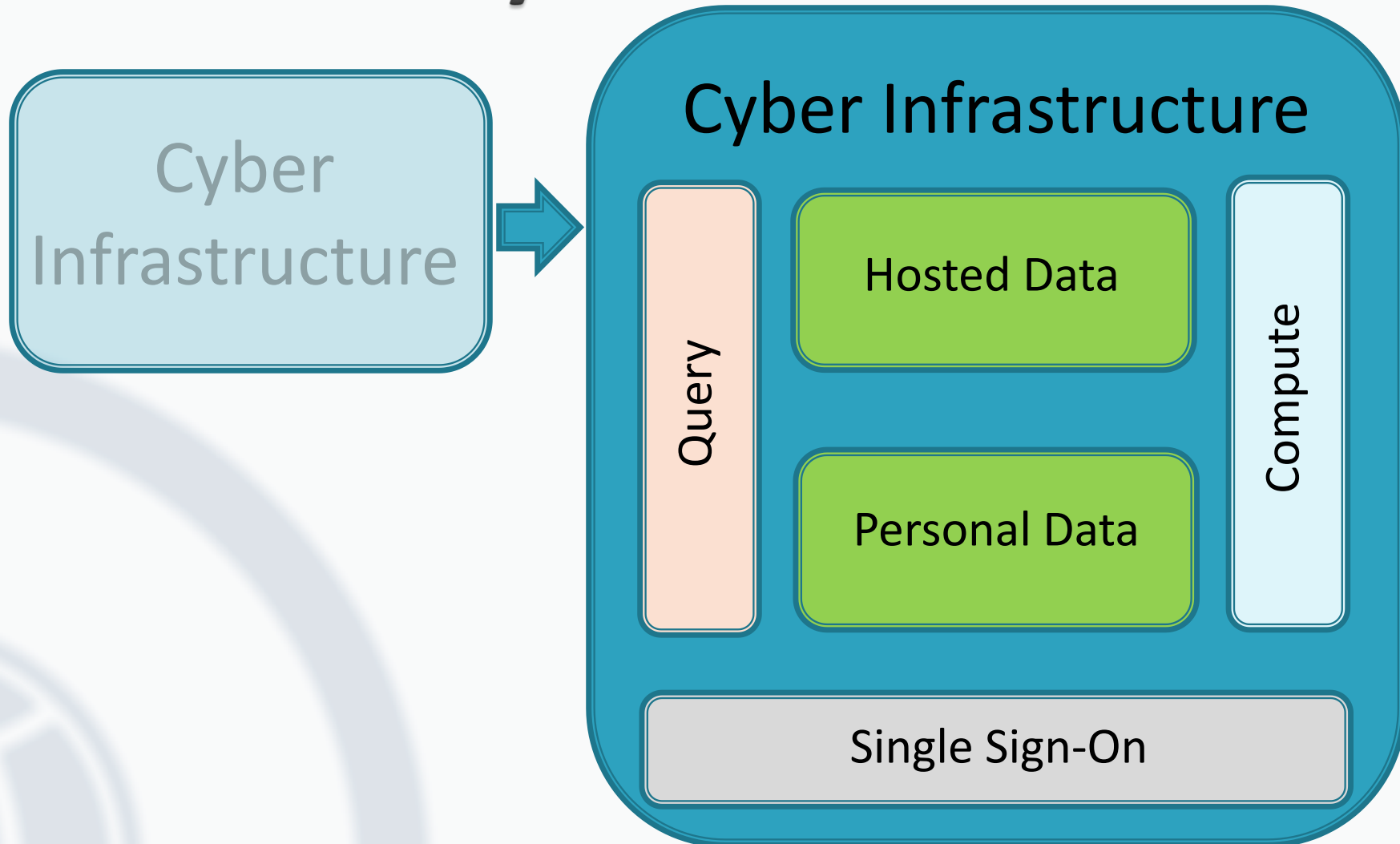
# How they fit together : Original



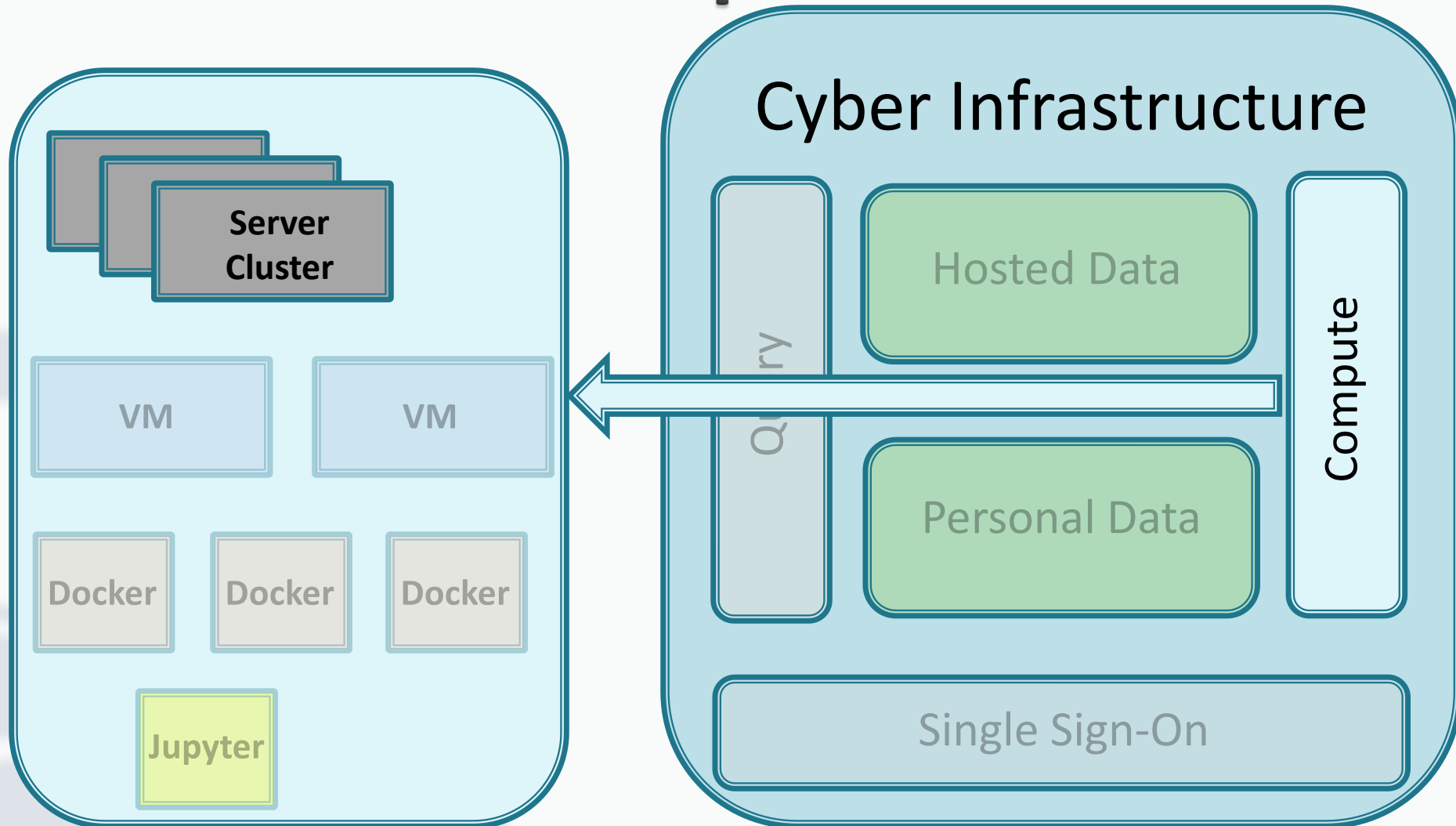
# Intermediate Refactoring



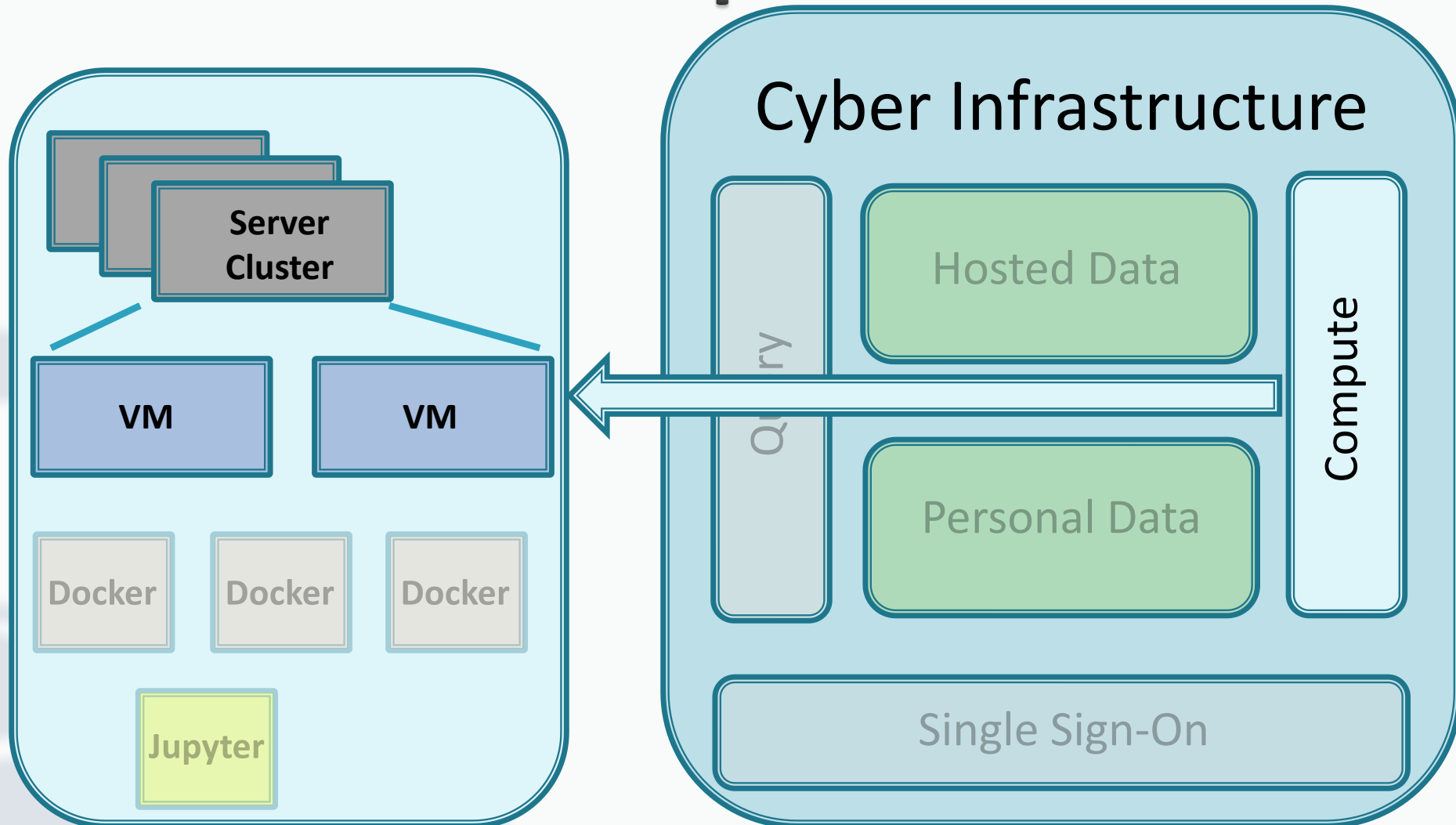
# SciServer : CyberInfrastructure



# SciServer : Compute

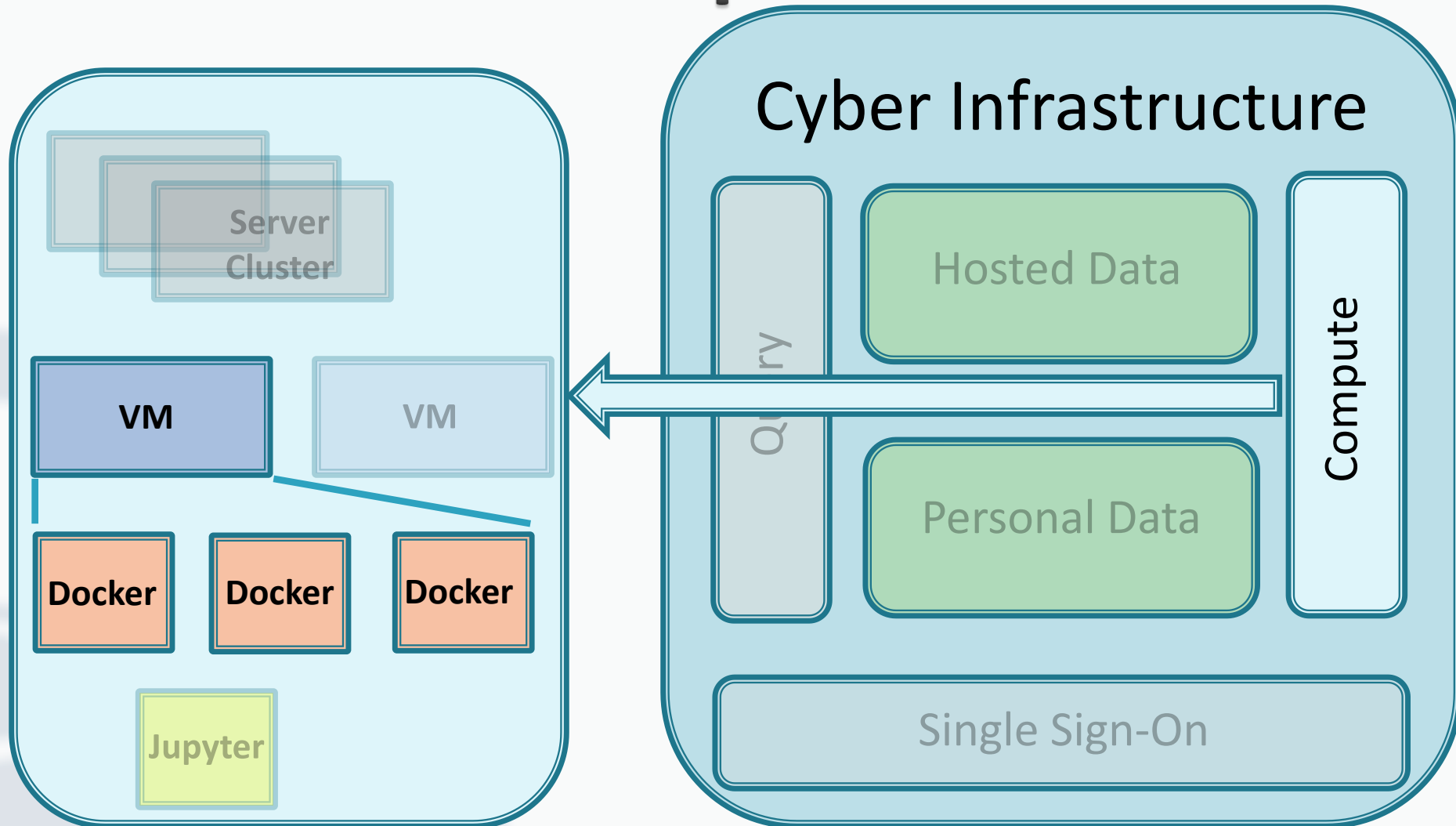


# SciServer : Compute

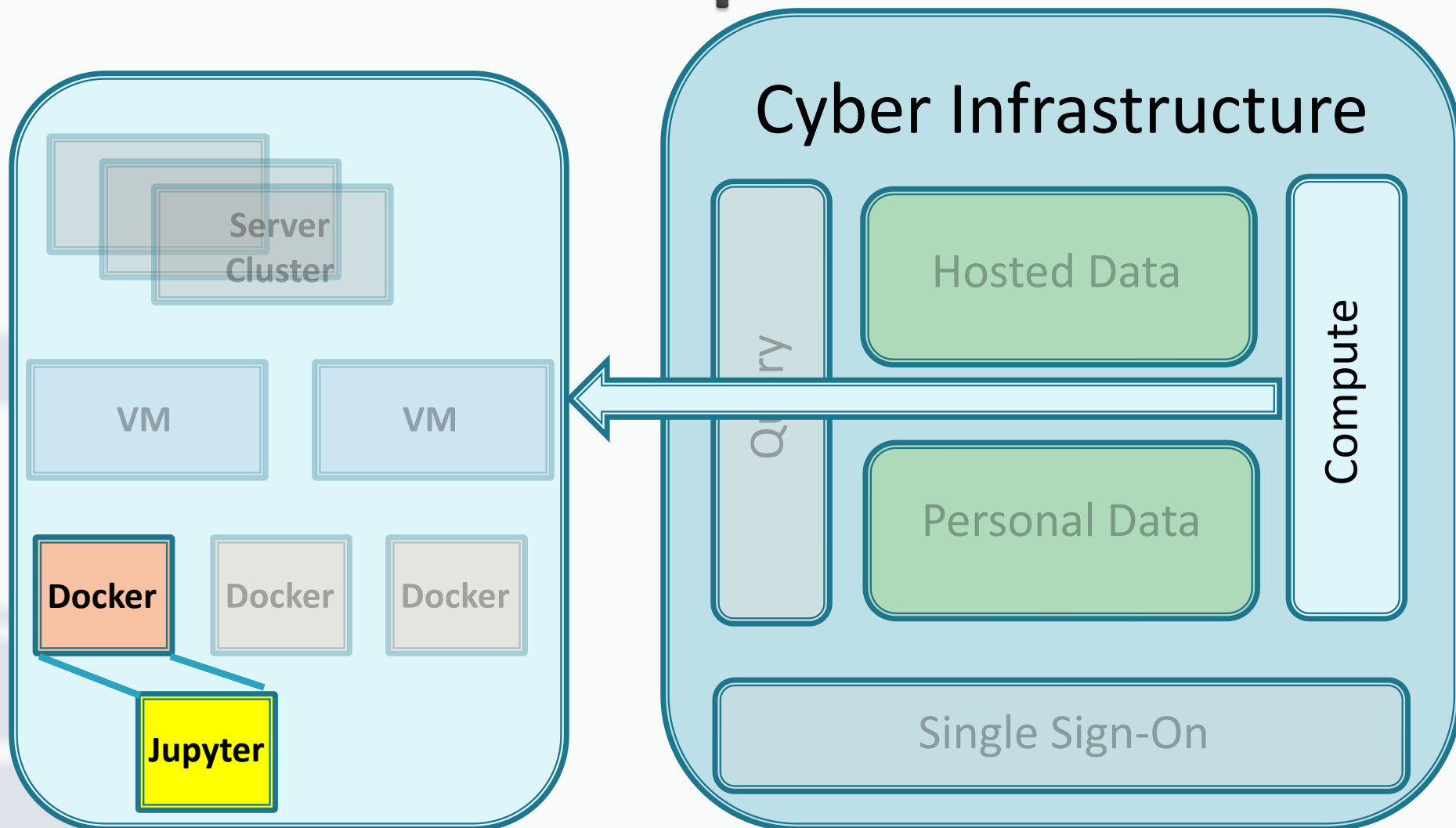




# SciServer : Compute



# SciServer : Compute



# SciServer : Compute

**INTERACTIVE &  
SYNCHRONOUS**

Cyber Infrastructure

Compute

Personal Data

Single Sign-On

Docker

Docker

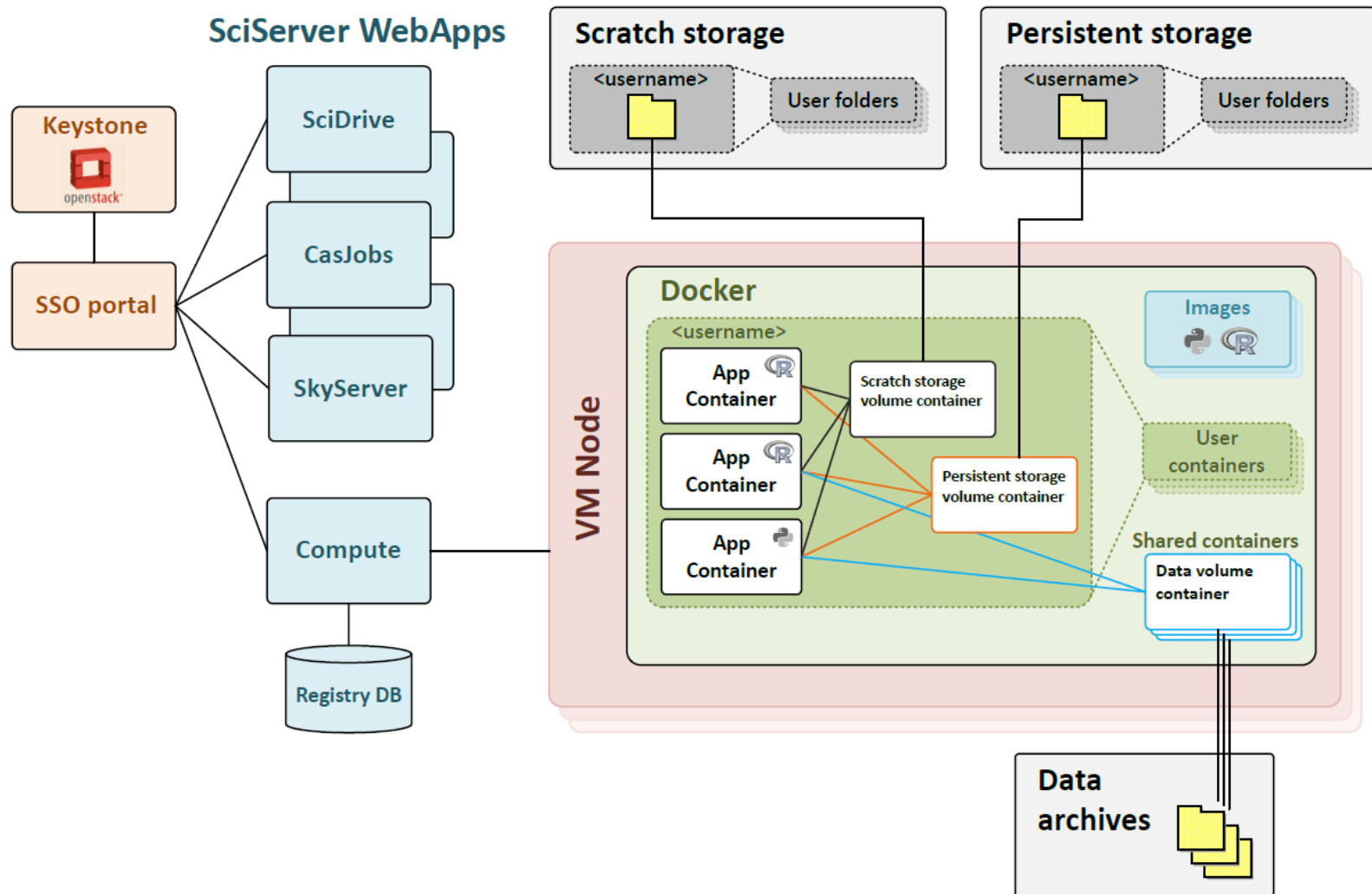
Docker

Jupyter

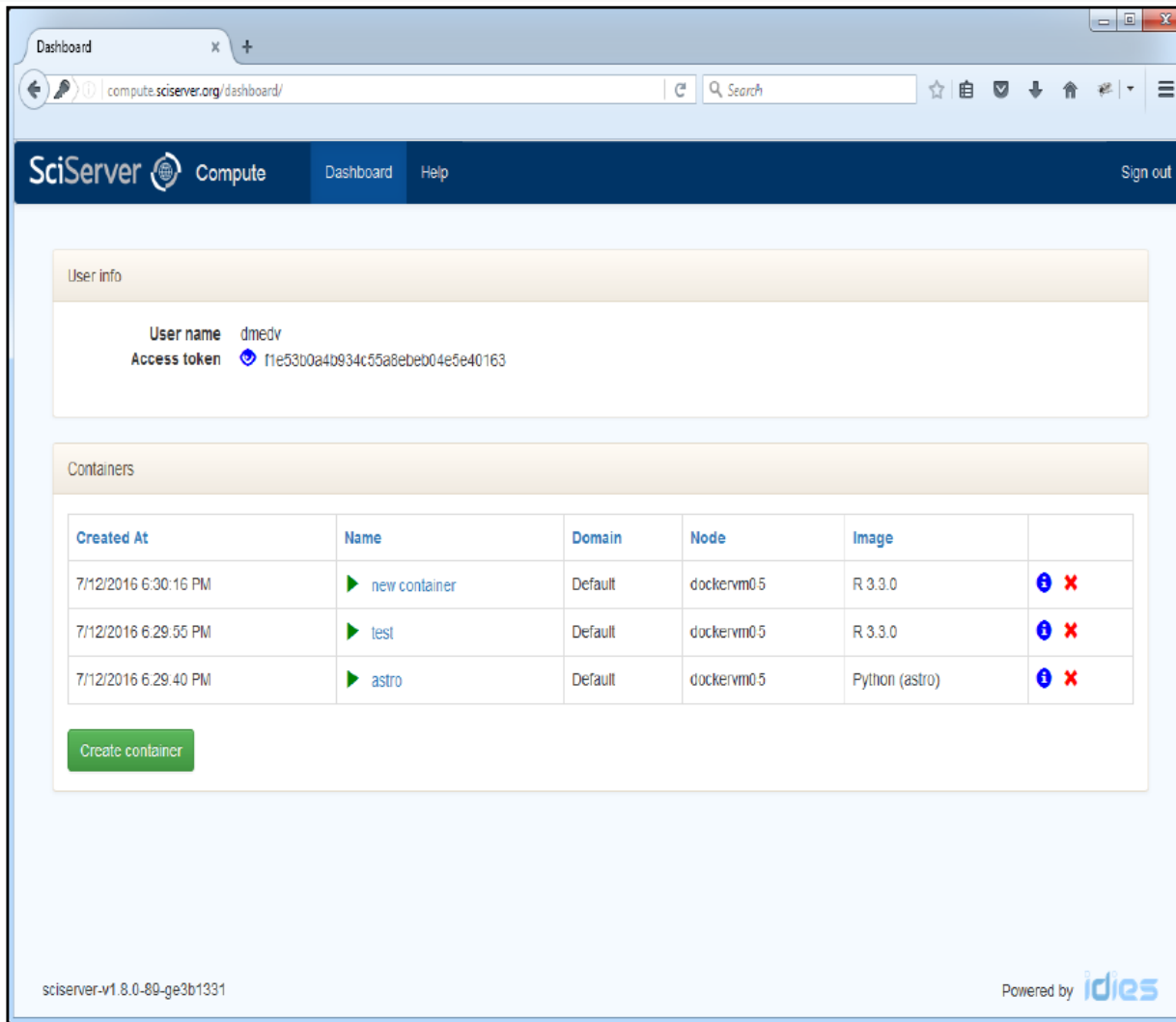
# SciServer : Compute

- ▶ “Engine” for executing analysis on data sets
- ▶ Environment for executing Python Notebooks **Interactively**
- ▶ Utility API Libraries in Python and R
- ▶ Interacts with ALL other SciServer components that have a WS API:
  - Login Portal for authentication
  - CASJobs for Queries
  - SkyServer and SkyQuery for Astronomy data
  - SciDrive for Storage

# Compute Architecture




# Compute Dashboard













The screenshot shows the SciServer Compute Dashboard in a web browser. The browser's address bar shows the URL `compute.sciserver.org/dashboard/`. The dashboard has a dark blue header with the SciServer logo, 'Compute' tab, 'Dashboard' and 'Help' links, and a 'Sign out' button. Below the header, there are two main sections: 'User info' and 'Containers'.


**User info**

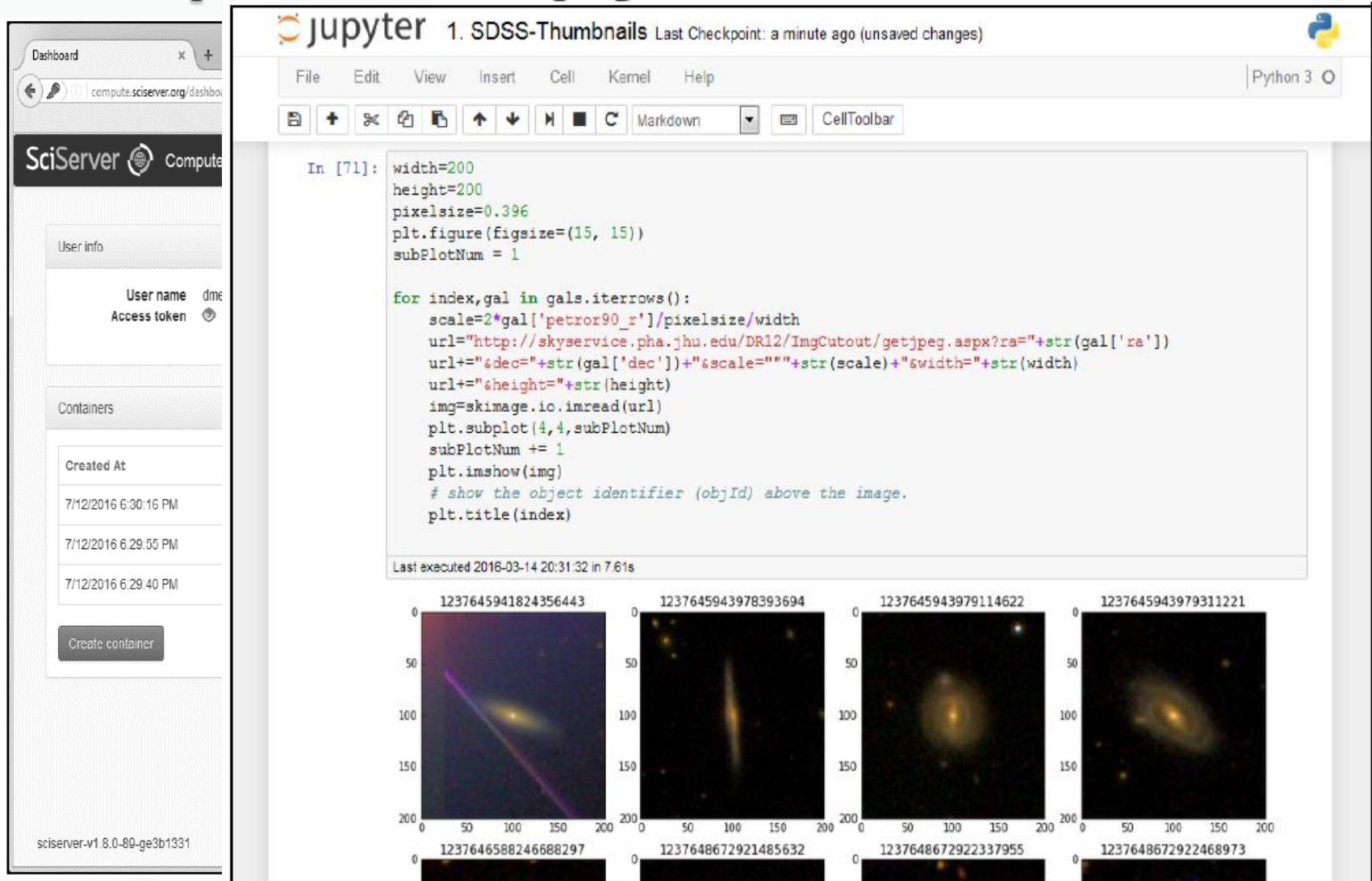
User name: dmedv  
Access token:  f1e53b00a4b934c55a8eb04e5e40163

**Containers**

Created At	Name	Domain	Node	Image	
7/12/2016 6:30:16 PM	 new container	Default	dockervm05	R 3.3.0	 
7/12/2016 6:29:55 PM	 test	Default	dockervm05	R 3.3.0	 
7/12/2016 6:29:40 PM	 astro	Default	dockervm05	Python (astro)	 



sciserver-v1.8.0-80-ge3b1331 Powered by 





# Example : Astronomy (SDSS)

jupyter 2. Stripe82-coadd-Copy1 Last Checkpoint: 3 minutes ago (unsaved changes)

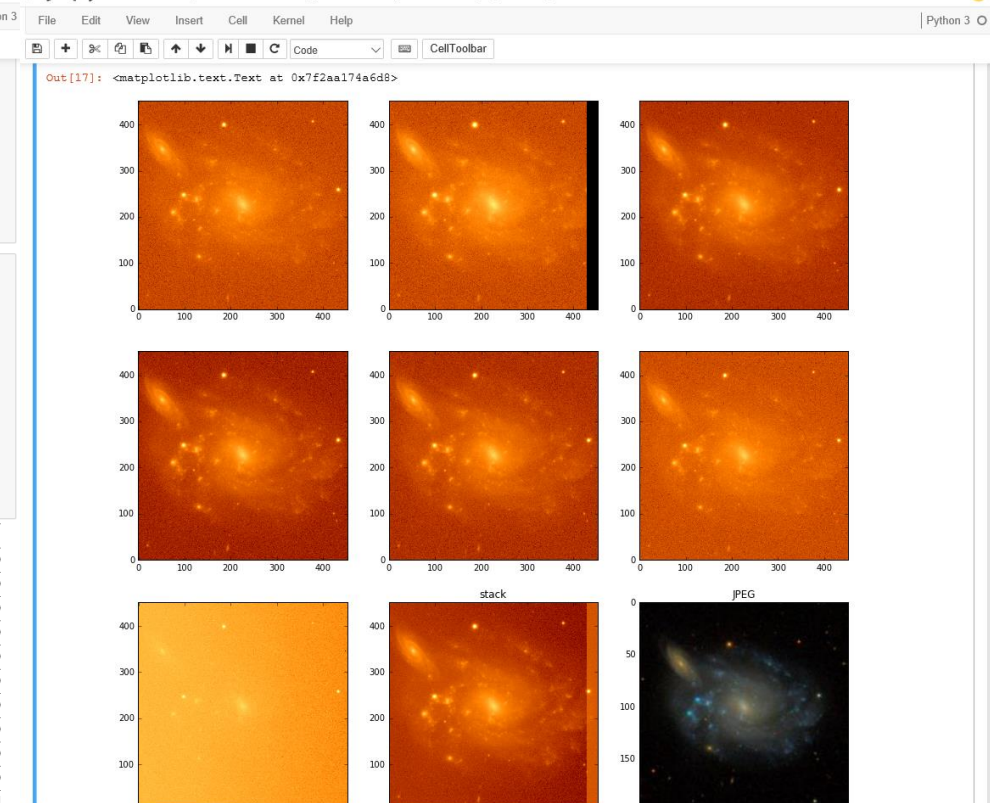
```
File Edit View Insert Cell Kernel Help Python 3
import SciServer.Login as login
token=login.getToken()
import pandas
import tables
import numpy as np
import astropy
from astropy.io import fits
from astropy import wcs
import skimage.io
import urllib
import os
import matplotlib
import matplotlib.pyplot as plt

In [18]: sql="""
SELECT a.objid as head, c.objid2 as match, b.matchcount,
       p.fieldid as head_field, d.fieldid as match_field,
       dbo.fGetUrlFitsCFrame(d.fieldid, 'g') as fits_g,
       dbo.fGetUrlFitsCFrame(d.fieldid, 'r') as fits_r,
       dbo.fGetUrlFitsCFrame(d.fieldid, 'z') as fits_z,
       p.ra, d.ra as match_ra, p.dec, d.dec as match_dec
       , p.petroz90_r
from (select top 1 * from galaxy where objid=8658194378960928809) a
join matchhead b on a.objid=b.objid -- join with matchhead
join photoobj p on a.objid=p.objid -- get matchhead photoobj
join match c on c.objid1=b.objid -- join with all the matches
join photoobjall d on c.objid2=d.objid -- get match photoobj
order by d.fieldid
"""
queryResponse = SciServer.CasJobs.executeQuery(sql, "Stripe82", token=token)
obs = pandas.read_csv(queryResponse, index_col=None)
obs[:10]

Out [18]:
```

	head	match	matchcount	head_field	match_field	fits_g
0	8658194378960928809	865819443049955320	57	8658194378960928768	8658194430499553280	<a href="http://das.sdss.org/imaging/5622/40/corr/">http://das.sdss.org/imaging/5622/40/corr/</a>
1	8658194378960928809	8658194477742948377	57	8658194378960928768	8658194477742948352	<a href="http://das.sdss.org/imaging/5633/40/corr/">http://das.sdss.org/imaging/5633/40/corr/</a>
2	8658194378960928809	8658194516375371821	57	8658194378960928768	8658194516375371776	<a href="http://das.sdss.org/imaging/5642/40/corr/">http://das.sdss.org/imaging/5642/40/corr/</a>
3	8658194378960928809	8658194585083510793	57	8658194378960928768	8658194585083510784	<a href="http://das.sdss.org/imaging/5658/40/corr/">http://das.sdss.org/imaging/5658/40/corr/</a>
4	8658194378960928809	8658194804163018771	57	8658194378960928768	8658194804163018752	<a href="http://das.sdss.org/imaging/5709/40/corr/">http://das.sdss.org/imaging/5709/40/corr/</a>
5	8658194378960928809	8658194954470752297	57	8658194378960928768	8658194954470752256	<a href="http://das.sdss.org/imaging/5744/40/corr/">http://das.sdss.org/imaging/5744/40/corr/</a>
6	8658194378960928809	8658195018907910161	57	8658194378960928768	8658195018907910144	<a href="http://das.sdss.org/imaging/5759/40/corr/">http://das.sdss.org/imaging/5759/40/corr/</a>
7	8658194378960928809	8658195044651040803	57	8658194378960928768	8658195044651040768	<a href="http://das.sdss.org/imaging/5765/40/corr/">http://das.sdss.org/imaging/5765/40/corr/</a>
8	8658194378960928809	8658195066151239724	57	8658194378960928768	8658195066151239680	<a href="http://das.sdss.org/imaging/5770/40/corr/">http://das.sdss.org/imaging/5770/40/corr/</a>
9	8658194378960928809	8658195113395748890	57	8658194378960928768	8658195113395748864	<a href="http://das.sdss.org/imaging/5781/40/corr/">http://das.sdss.org/imaging/5781/40/corr/</a>

jupyter 2. Stripe82-coadd-Copy1 Last Checkpoint: a minute ago (autosaved)





**Jupyter** | **FragData - Analysis: Demo** Last checkpoint 4 minutes ago (autosaved)

File Edit View Insert Cell Kernel Help Python 3

---

**In [1]:**

```
import SciServer.LoginPortal as Login
login=Login.getLogin()
import SciServer.CasJobs as CJ
import pandas
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.cm as cm
from matplotlib.collections import PolyCollection
from matplotlib.collections import LineCollection
```

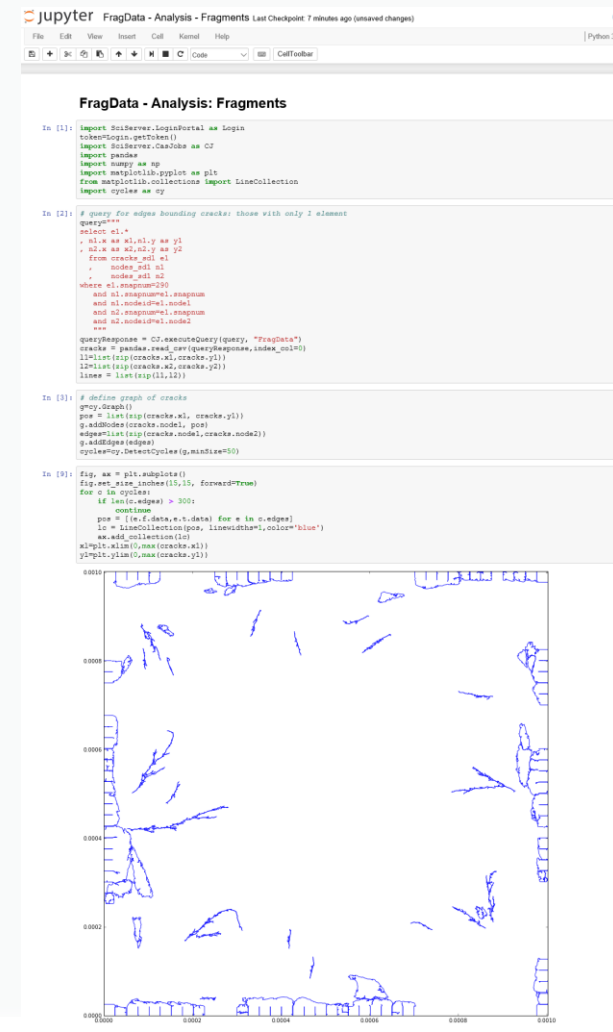
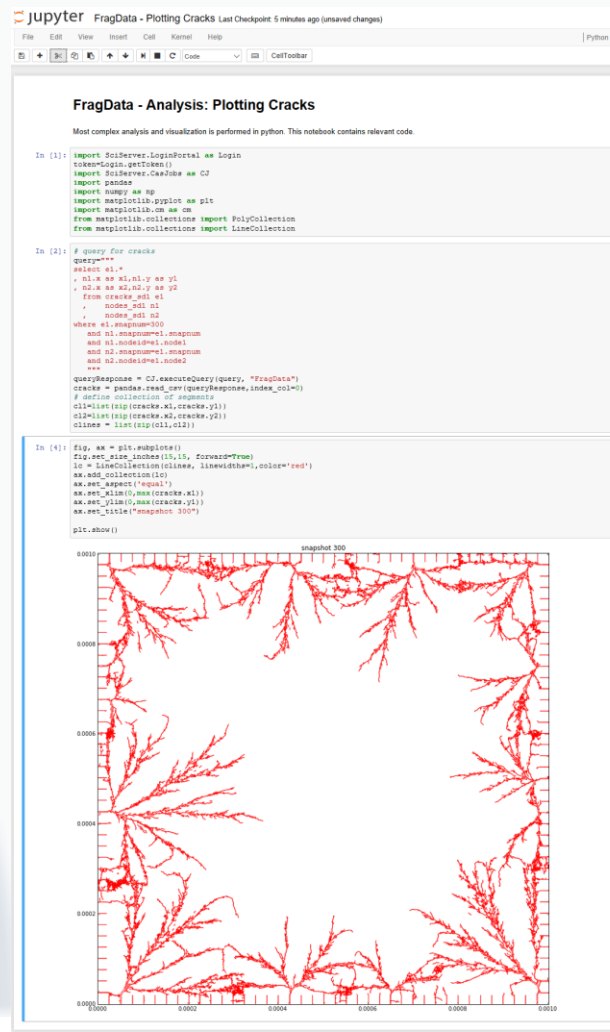
**In [2]:**

```
query="""
select e.elementid
      , e.x1,y1,e.x2,y2,e.k1,e.y0
      , e.k0,e.s1,e.s2,e.d
      , 0.5*(e.k0+e.s1) as t, e.k0*e.k0+e.s1*e.s2 as d
from elements_e_sdl e
where e.elementid < 10
and e.k1 between 0.00018 and 0.0002
and e.y1 between 0.00018 and 0.0002
"""

# query CasJobs table. Using FragData as context
sr = CJ.executeQuery(query,"FragData")
elements_pandas.read_sql(sr,(login.portalName))
l1=[x[0] for x in sr.elements_1]
l2=[x[1] for x in sr.elements_2]
l3=[x[2] for x in sr.elements_3]
l4=[x[3] for x in sr.elements_4]
verts=np.vstack((l1,l2,l3))
l1_elements,t=np.array(elements_1).elements_1-vertices,d
colorMap_array(elements_1)
```

**In [4]:**

```
fig, ax = plt.subplots()
fig.set_size_inches(10,10, forward=True)
coll = PolyCollection(verts, array=colors, cmap=cm.jet, edgecolors='none')
ax.add_collection(coll)
ax.autoscale_tight()
ax.set_aspect('equal')
ax.set_xlimmin(elements_1),max(elements_1))
ax.set_ylimmin(elements_2),max(elements_2))
plt.show()
```



# Example : Cosmology

jupyter CosmoUC\_5\_PlotHalo Last Checkpoint: 11 minutes ago (autosaved)

File Edit View Insert Cell Kernel Help  
Notebook saved Python 3

## Dark-matter Halos in a Cosmological Simulation

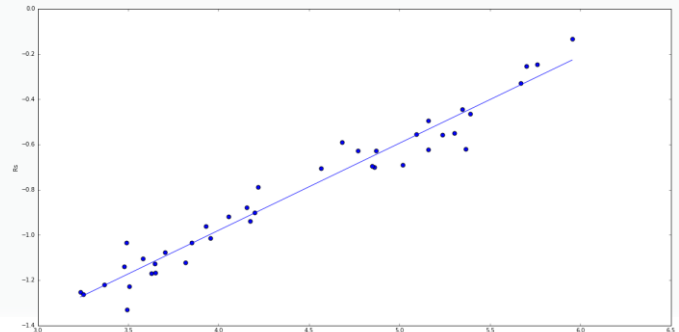
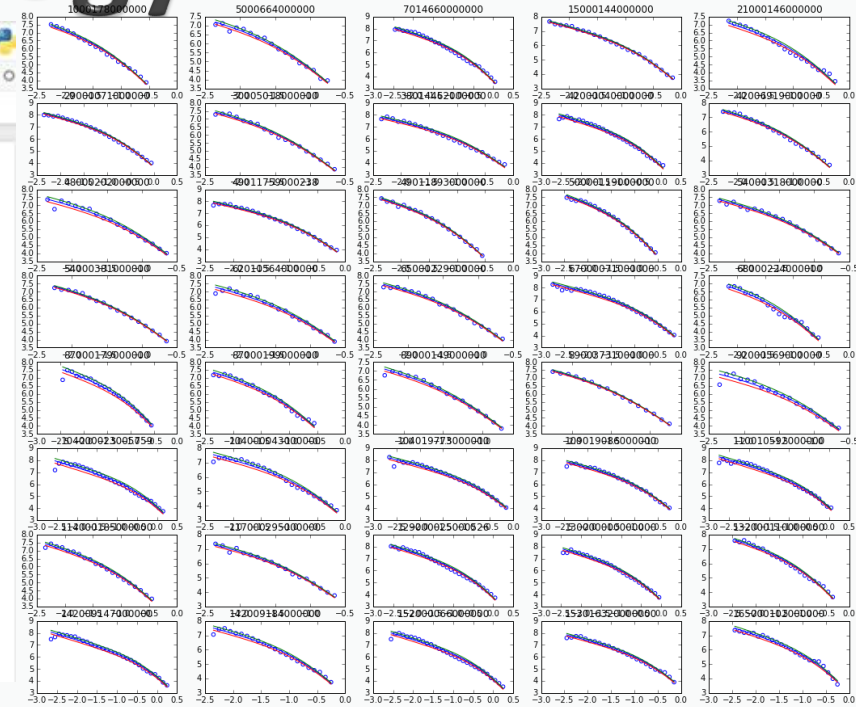
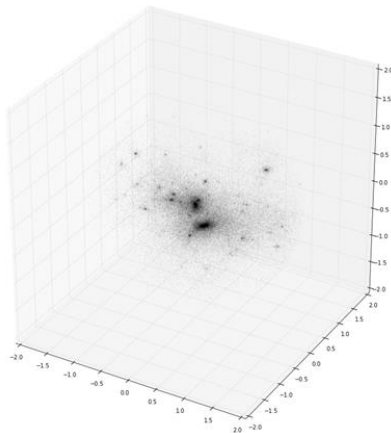
```
In [3]: import SciServer.LoginPortal as Login
token = Login.getToken()
import SciServer.CasJobs
import pandas
import tables
import numpy as np
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
```

```
In [10]: %%time
queryString = """
select top 100000 p.x-hh.x as x,p.y-hh.y as y,p.z-hh.z as z
from mpahaltrees.mr hh
cross apply dbo.MillenniumParticles(hh.snapnum,
dbo.Sphere::New(hh.x,hh.y,hh.z,3*hh.halfmassradius).ToString()) p
where hh.haloId=8400007000000 order by newid()
"""
responseStream = SciServer.CasJobs.executeQuery(queryString, token=token, context="SimulationDB")
df = pandas.read_csv(responseStream, index_col="None")

CPU times: user 331 ms, sys: 184 ms, total: 535 ms
Wall time: 5.27 s
```

```
In [13]: fig = plt.figure(figsize=(15, 15))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(df.x, df.y, df.z, s=0.001)
```

```
Out[13]: <mpl_toolkits.mplot3d.art3d.Path3DCollection at 0x7fe86428b9e8>
```



# SciServer

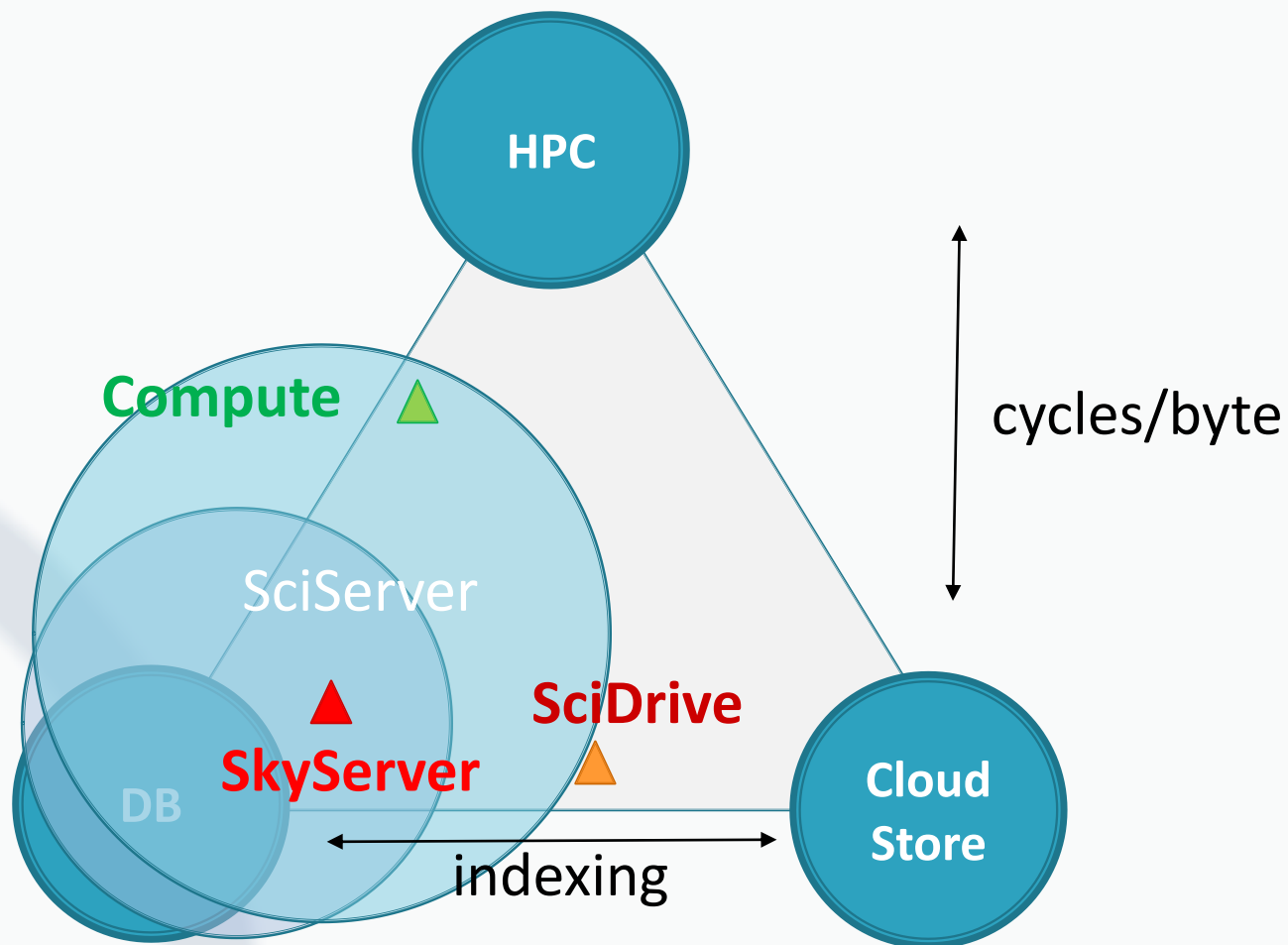


Collaborative data-driven science

## Compute Next Stage

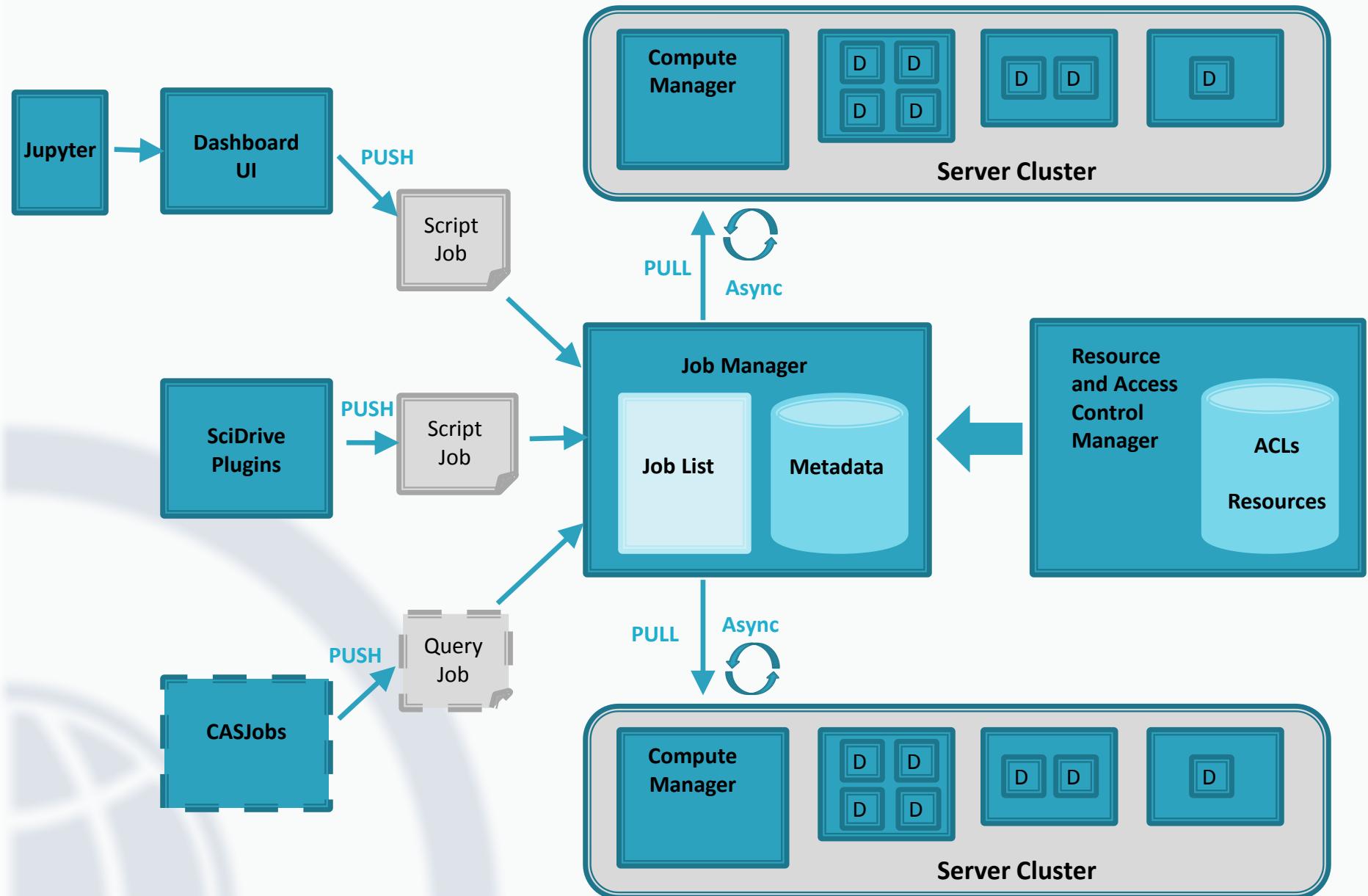


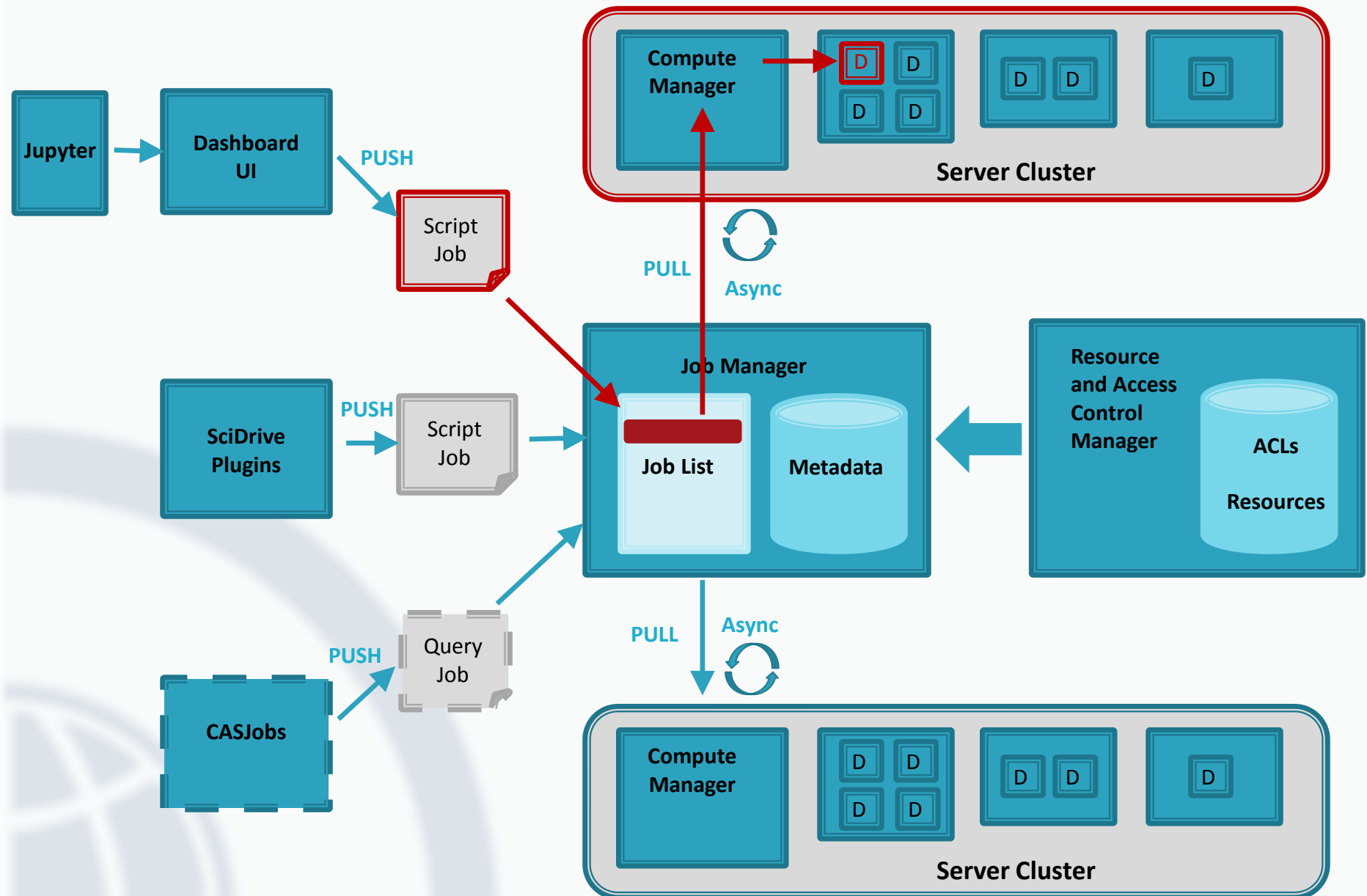
# System Balance

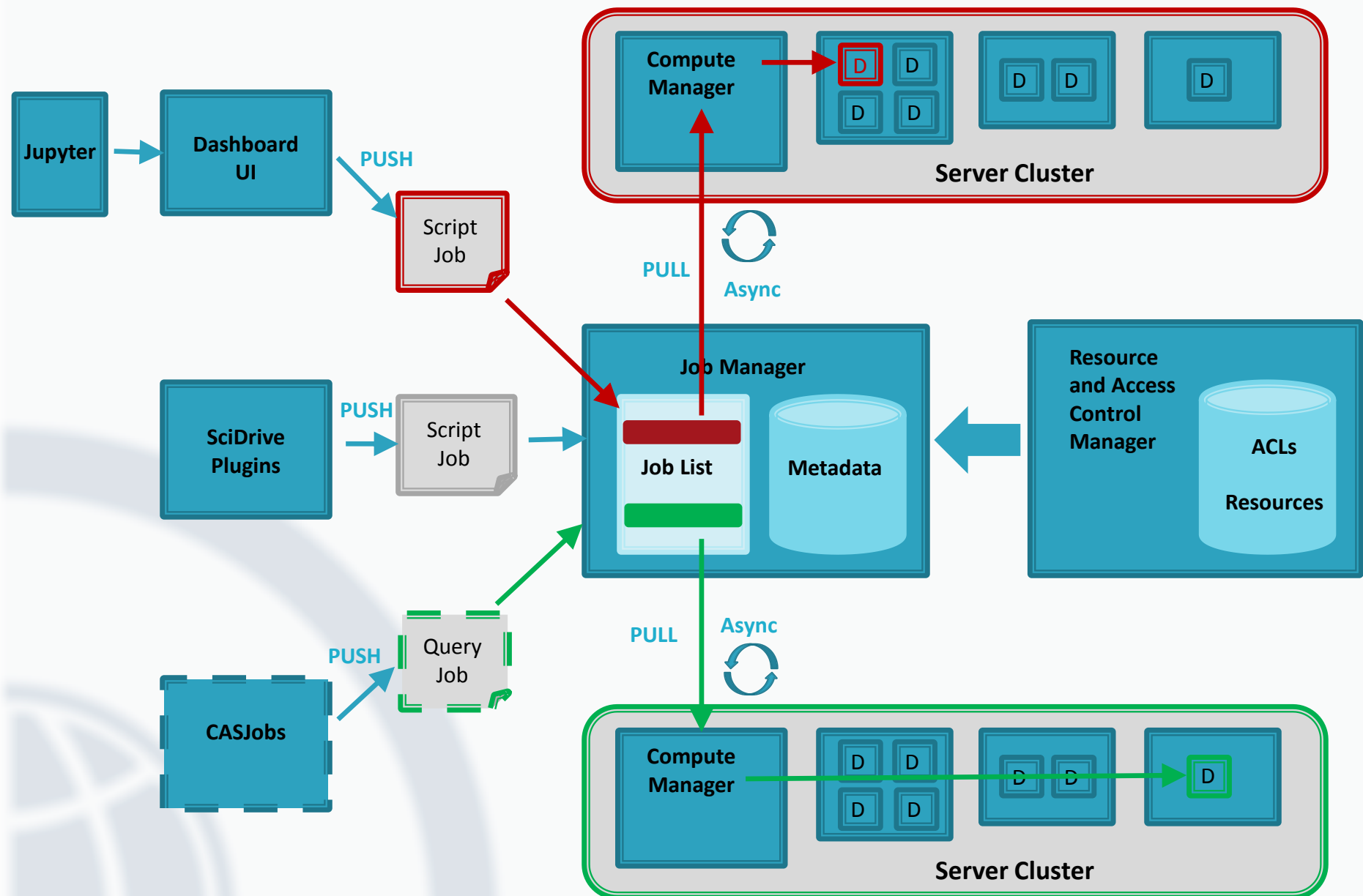


# Compute Development

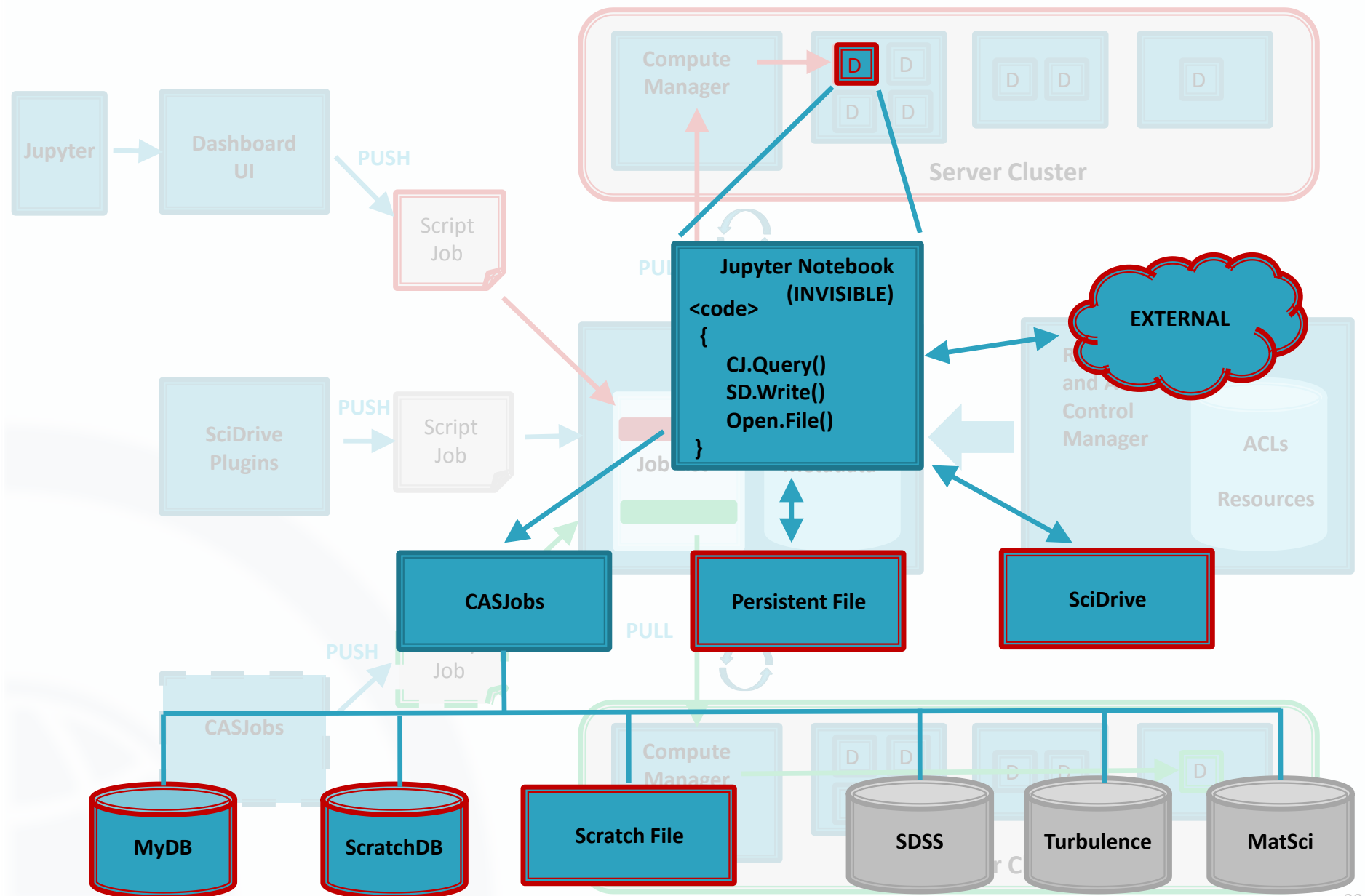
- ▶ Build on VM/Docker Architecture
- ▶ Scalable **non-interactive, asynchronous** Job management (JOBM)
- ▶ Rich Access Controls (RACM)
- ▶ Distributed compute execution (COMPM)
- ▶ Support Python, R, Matlab











# Summary

- ▶ **SciServer Compute Interactive is live now**
- ▶ **Supports Python, R, Jupyter**
- ▶ **Runs on a 4 node cluster**
- ▶ **Access to several domain databases**
- ▶ **Asynchronous Job Execution early 2017**
- ▶ **Please register with SciServer and try it out**

# SciServer



Collaborative data-driven science

## Questions?

