

SciServer, a collaborative science platform with cosmological simulations

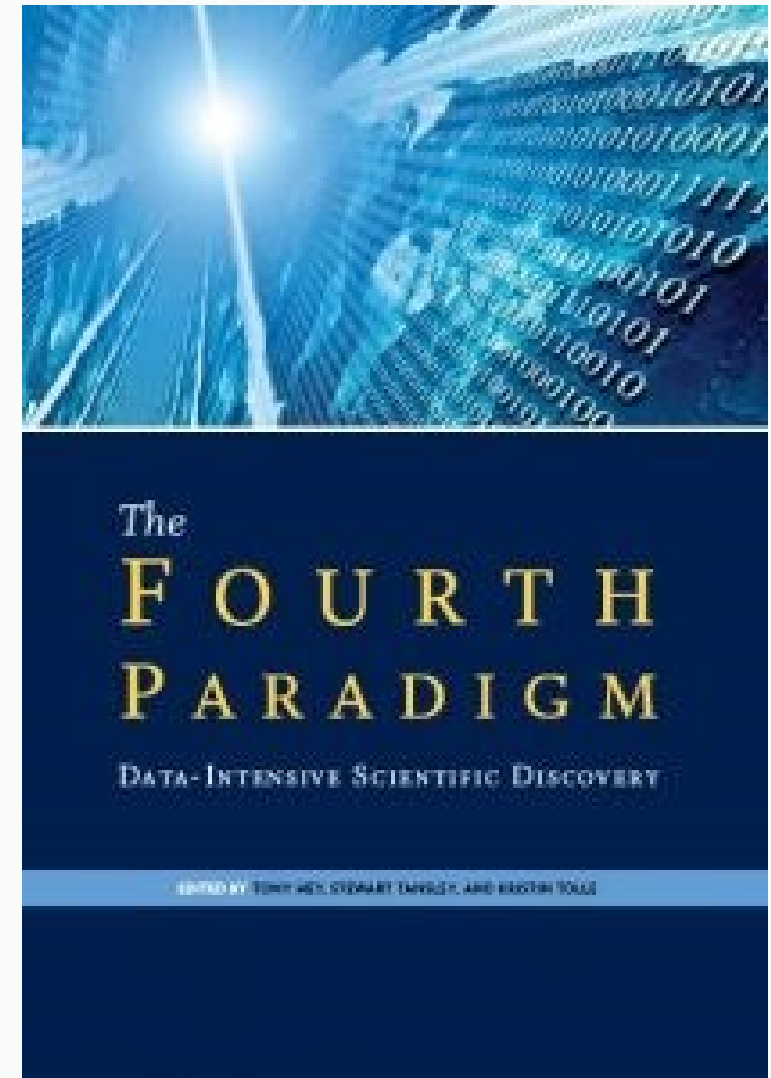
Gerard Lemson
Institute for Data Intensive Engineering and Science (IDIES)
The Johns Hopkins University

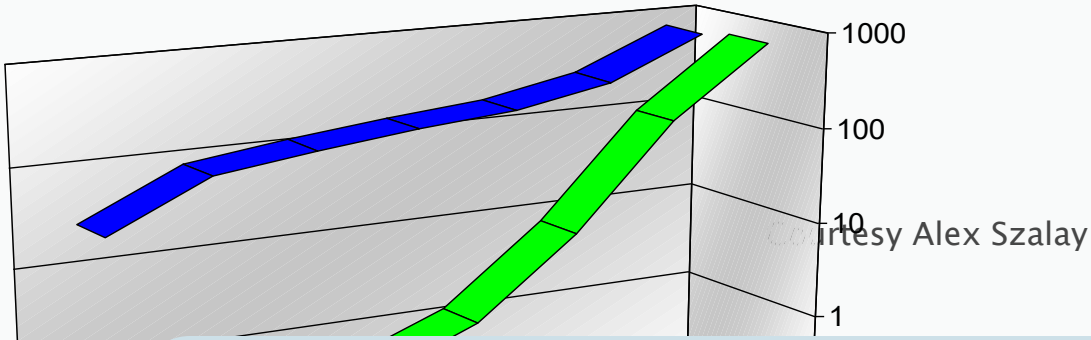


The 4th paradigm : data intensive scientific discovery

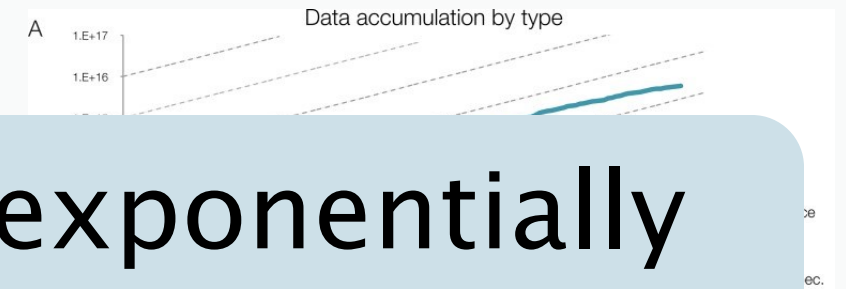
- Jim Gray

Increasingly scientific investigations require **combination of large amounts of data** from many different sources and ever more **sophisticated machine learning algorithms and tools** for their analysis.



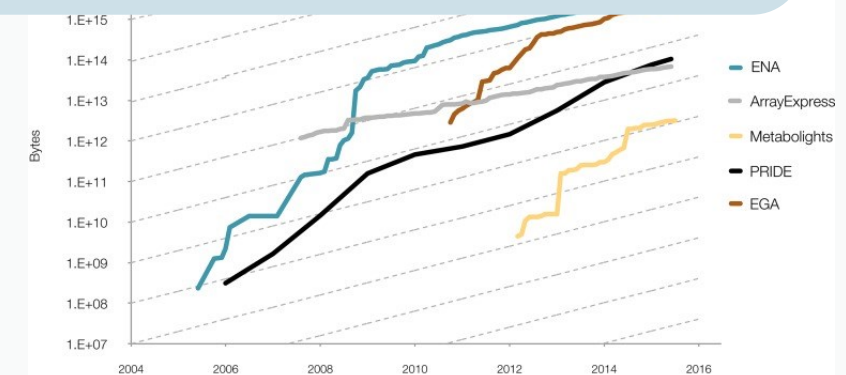
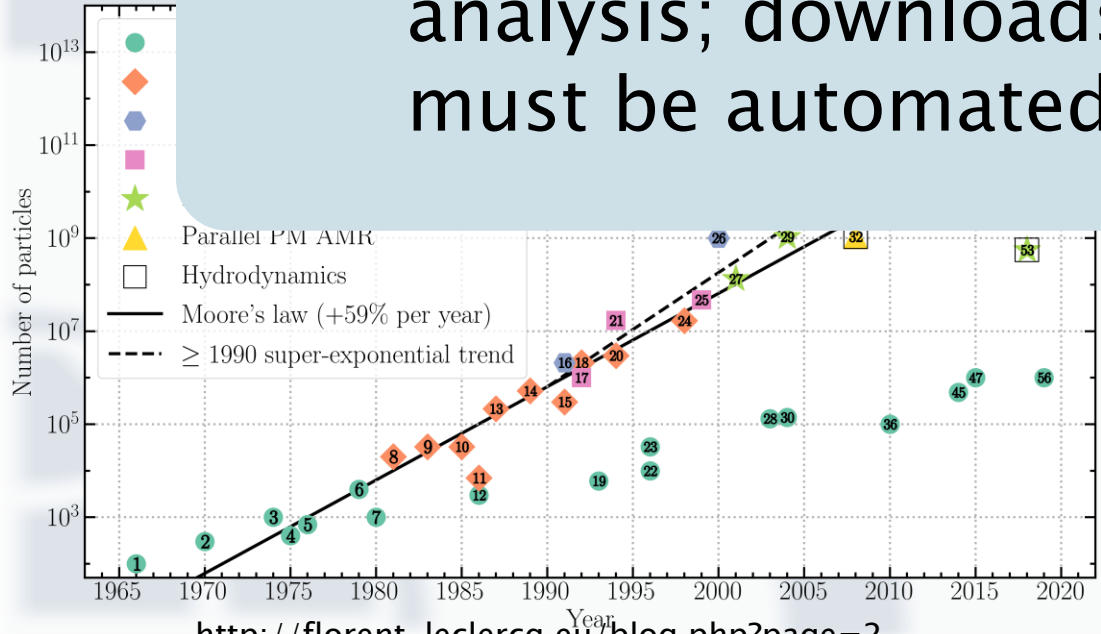


The European Bioinformatics Institute in 2016: Data growth and integration
 Charles E. Cook, Mary Todd Bergman, Robert D. Finn, Guy Cochrane, Ewan Birney, and Rolf Apweiler



Science data growing exponentially

Expertise required to optimize big data analysis; downloads unfeasible; analysis must be automated and *data-proximate*



Planck 2018 results. IX. Constraints on primordial non-Gaussianity

Show affiliations Hide authors

Planck Collaboration; Akrami, Y. ; Arroja, F.; Ashdown, M.; Aumont, J. ; Baccigalupi, C.; Ballardini, M.; Banday, A. J.; Barreiro, R. B.; Bartolo, N.; Basak, S.; Benabed, K.; Bernard, J. -P.; Bersanelli, M.; Bielewicz, P.; Bond, J. R.; Borrill, J.; Bouchet, F. R.; Bucher, M.; Burigana, C. ; Butler, R. C.; Calabrese, E.; Cardoso, J. -F.; Casaponsa, B.; Challinor, A.; Clowe, N.; Colombo, G. P. G.; Combès, M.; Cori, A.; Crill, B. P.; Cucunubá, J.; Delabrouille, J.; Douspis, M.; Enßlin, T. A. ; Fernandez-Cobate, R.; Franceschi, E.; Génova-Santo, J.; Gratton, S.; Gundersen, R.; Handley, W. ; Jaffe, A. H. ; Kiiveri, K.; Kirton, M. J.; Lamarre, J. -M.; Levrier, F.; López-Cañiego, A.; Maino, D.; Mandoùs, S.; Martin, P. G. ; McEwen, J. D.; Mennella, A.; Migliaccio, M.; Miville-Deschenes, M. -A.; Molinari, D. ; Moneti, A.; Montier, L.; Morgante, G. ; Moss, A.; Münchmeyer, M.; Natoli, P.; Oppizzi, F.; Pagano, L.; Paoletti, D. ; Partridge, B.; Patanchon, G.; Perrotta, F.; Pettorino, V.; Piacentini, F. ; Polenta, G.; Puget, J. -L.; Rachen, J. P.; Racine, B. ; Reinecke, M.; Remazeilles, M.; Renzi, A.; Rocha, G. ; Rubiño-Martín, J. A.; Ruiz-Granados, B.; Salvati, L.; Savelainen, M.; Scott, D.; Shellard, E. P. S.; Shiraishi, M.; Sirignano, C.; Sirri, G.; Smith, K.; Spencer, L. D.; Stanco, L.; Sunyaev, R.; Suur-Uski, A. -S.; Tauber, J. A.; Tavagnacco, D.; Tenti, M.; Toffolatti, L.; Tomasi, M.; Trombetti, T.; Valiviita, J. ; Van Tent, B.; Vielva, P.; Villa, F. ; Vittorio, N.; Wandelt, B. D.; Wehus, I. K.; Zacchei, A. ; Zonca, A.



The Sloan Digital Sky Survey: Technical Summary

Show affiliations Hide authors

York, Donald G.; Adelman, J.; Anderson, John E., Jr.; Anderson, Scott F.; Annis, James; Bahcall, Neta A.; Bakken, J. A.; Barkhouser, Robert; Bastian, Steven; Berman, Eileen; Boroski, William N.; Bracker, Steve; Briegel, Charlie; Briggs, John W.; Brinkmann, J.; Brunner, Robert; Burles, Scott; Carey, Larry; Carr, Michael A.; Castander, Francisco J.; Chen, Bing; Colestock, Patrick L.; Connolly, A. J.; Crocker, J. H.; Okamura, Sadanori; Ostriker, Jeremiah P.; Owen, Russell; Pauls, A. George; Peoples, John; Peterson, R. L.; Petravick, Donald; Pier, Jeffrey R.; Pope, Adrian; Pordes, Ruth; Prosapio, Angela; Rechenmacher, Ron; Quinn, Thomas R. ; Richards, Gordon T. ; Richmond, Michael W.; Rivetta, Claudio H.; Rockosi, Constance M.; Ruthmansdorfer, Kurt; Sandford, Dale; Schlegel, David J.; Schneider, Donald P.; Sekiguchi, Maki; Sergey, Gary; Shimasaku, Kazuhiro; Siegmund, Walter A.; Smee, Stephen; Smith, J. Allyn ; Snedden, S.; Stone, R.; Stoughton, Chris; Strauss, Michael A.; Stubbs, Christopher; SubbaRao, Mark; Szalay, Alexander S. ; Szapudi, Istvan ; Szokoly, Gyula P.; Thakar, Anirudda R.; Tremonti, Christy; Tucker, Douglas L. ; Uomoto, Alan; Vanden Berk, Dan; Vogeley, Michael S.; Waddell, Patrick; Wang, Shu-i.; Watanabe, Masaru; Weinberg, David H.; Yanny, Brian; Yasuda, Naoki; SDSS Collaboration

Full length article

The illustris simulation: Public data release ☆

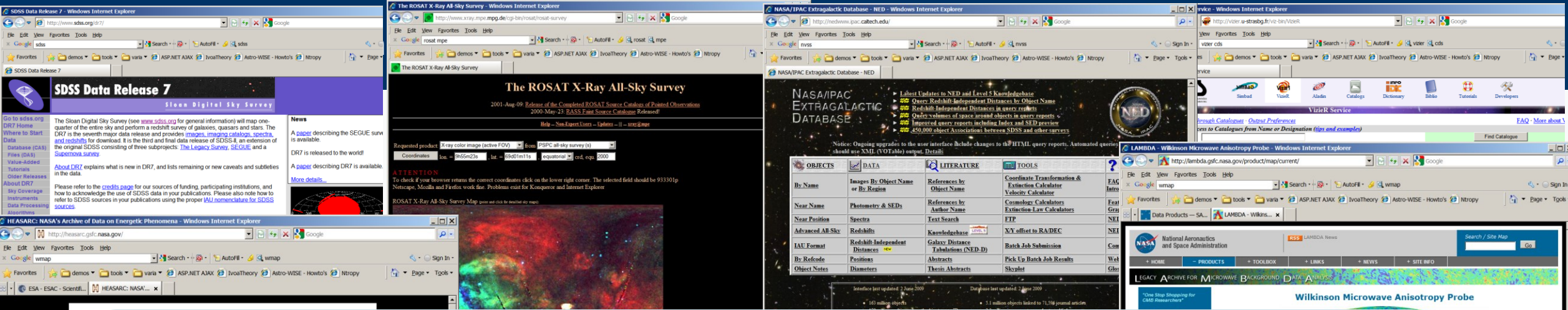
D. Nelson ^a  , A. Pillepich ^a, S. Genel ^{b, a, 1}, M. Vogelsberger ^c, V. Springel ^{d, e}, P. Torrey ^{c, g}, V. Rodriguez-Gomez ^a, D. Sijacki ^f, G.F. Snyder ^h, B. Griffen ^c, F. Marinacci ^c, L. Blecha ^j, L. Sales ⁱ, D. Xu ^d, L. Hernquist ^a



Science increasingly collaborative

Data and analysis sharing often ad hoc

Cambridge, MA,
street, New
f Physics, MIT,
nweg 35, 69118
ofstr. 12-14,
ty of
^a TAPIR, Mailcode 350-17, California Institute of Technology, Pasadena, CA 91125, USA
^b Space Telescope Science Institute, 3700 San Martin Dr, Baltimore, MD 21218, USA
ⁱ Department of Physics and Astronomy, University of California, Riverside, 900 University Avenue, Riverside, CA 92521, USA
^j University of Maryland, College Park, Department of Astronomy and Joint Space Science Institute, USA



Data increasingly open/public

Improve science by combining data sets from different sources



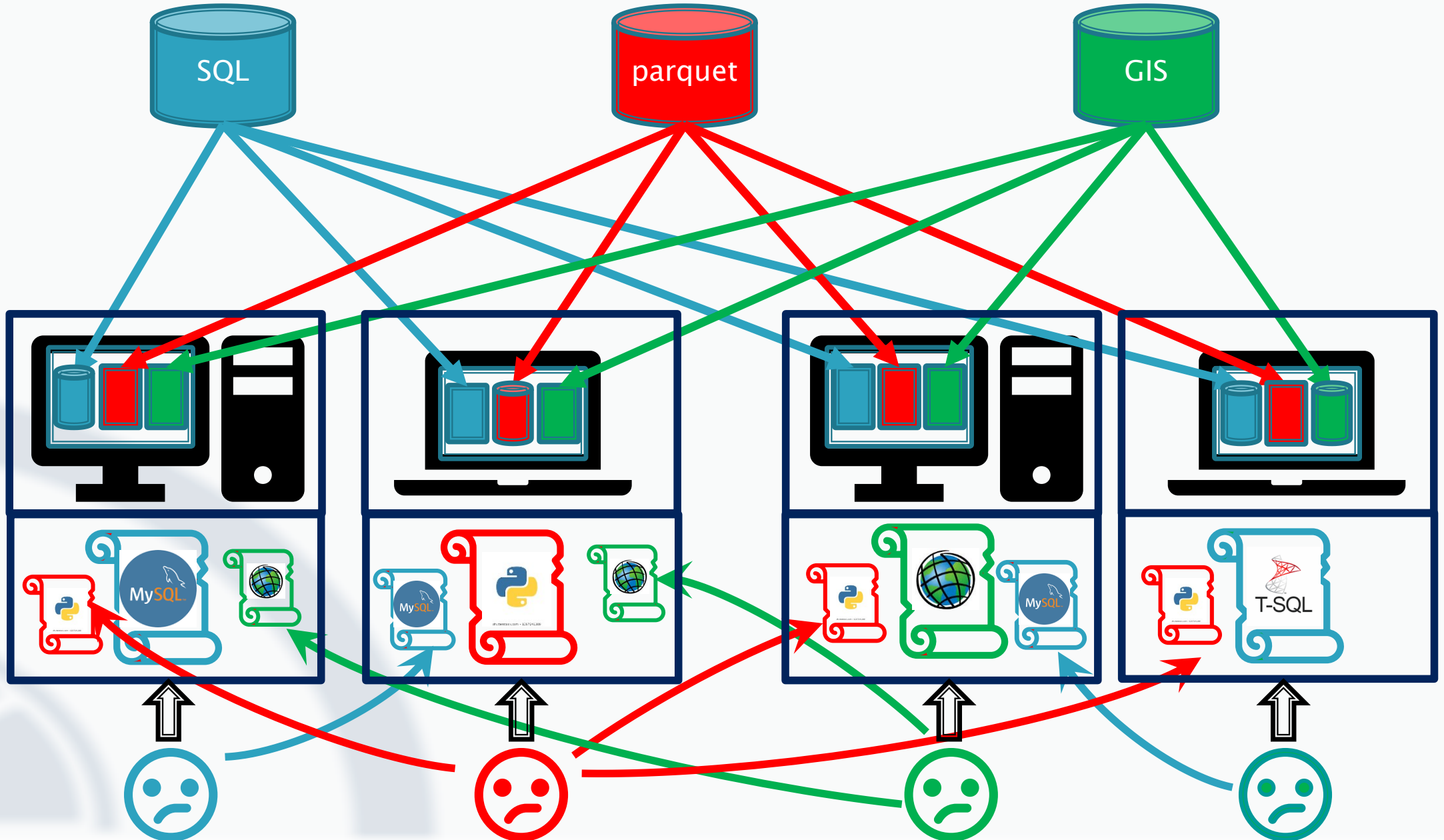
Data sources/formats heterogeneous

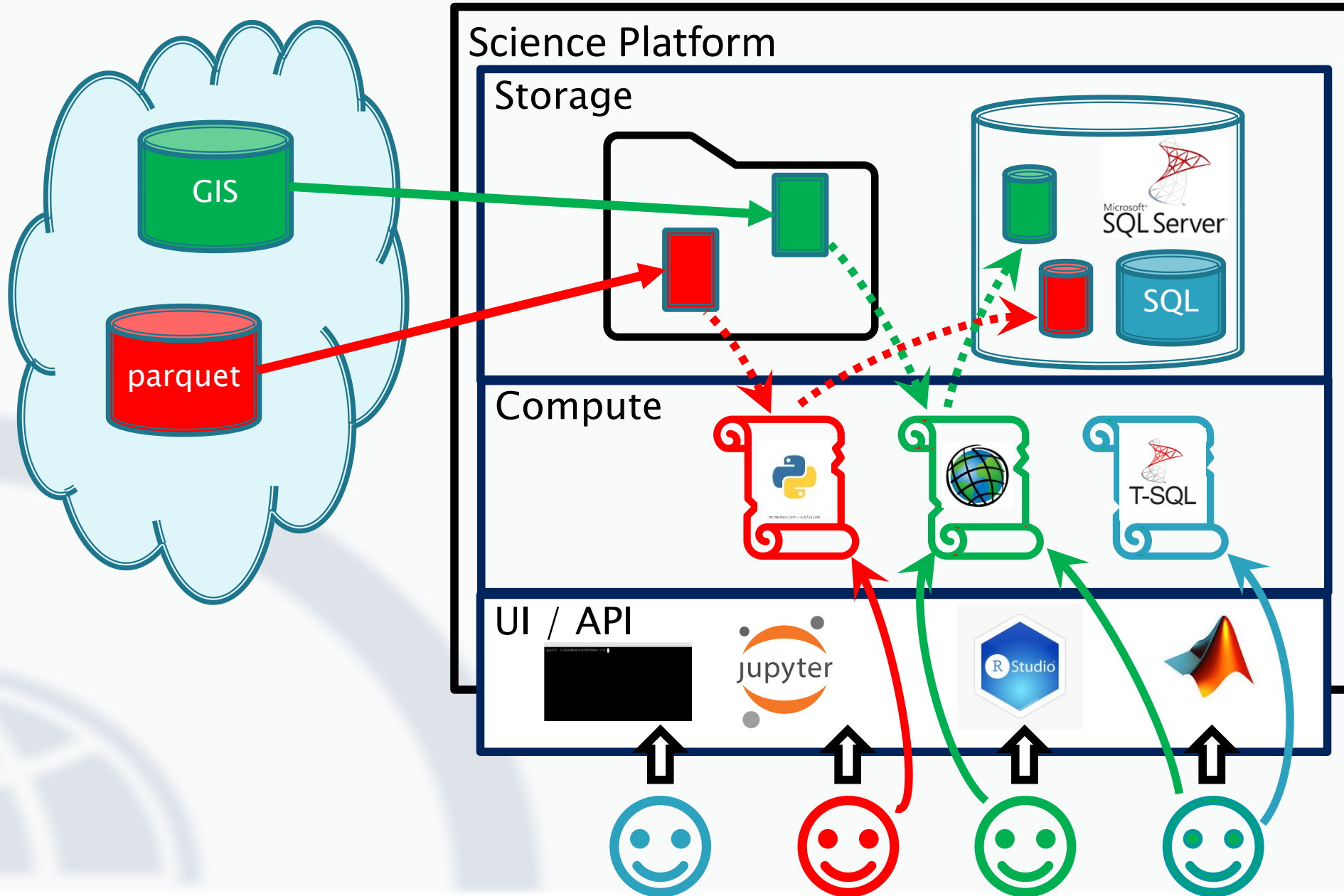
Diverse set of skills and knowledge required, both technical and domain knowledge



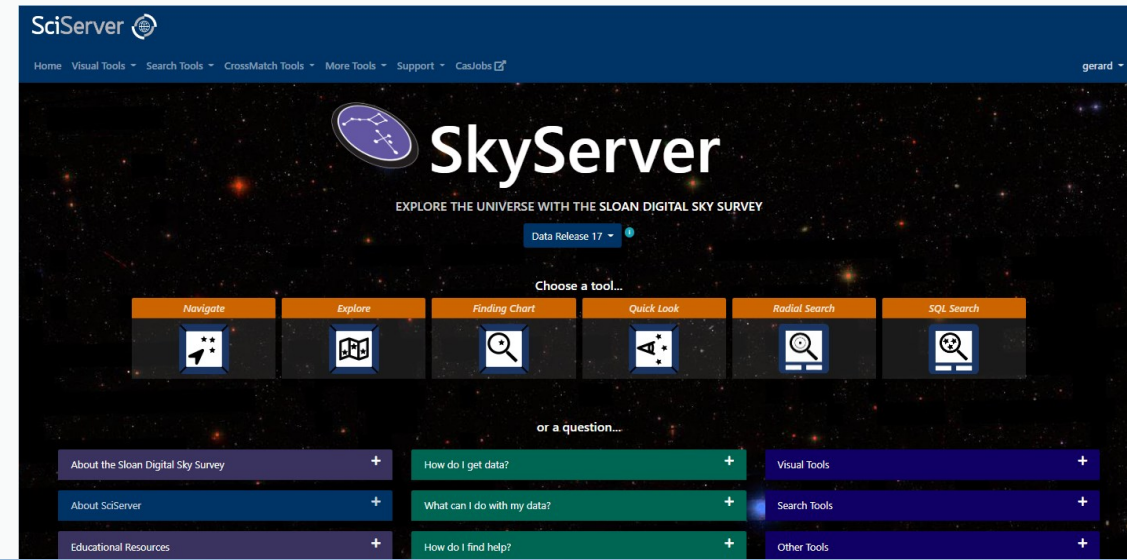
Question:

- ▶ How can we improve the support for **collaborative science** projects producing and requiring **large, heterogeneous data sets** with **geographically distributed** partners with **varying expertise**?





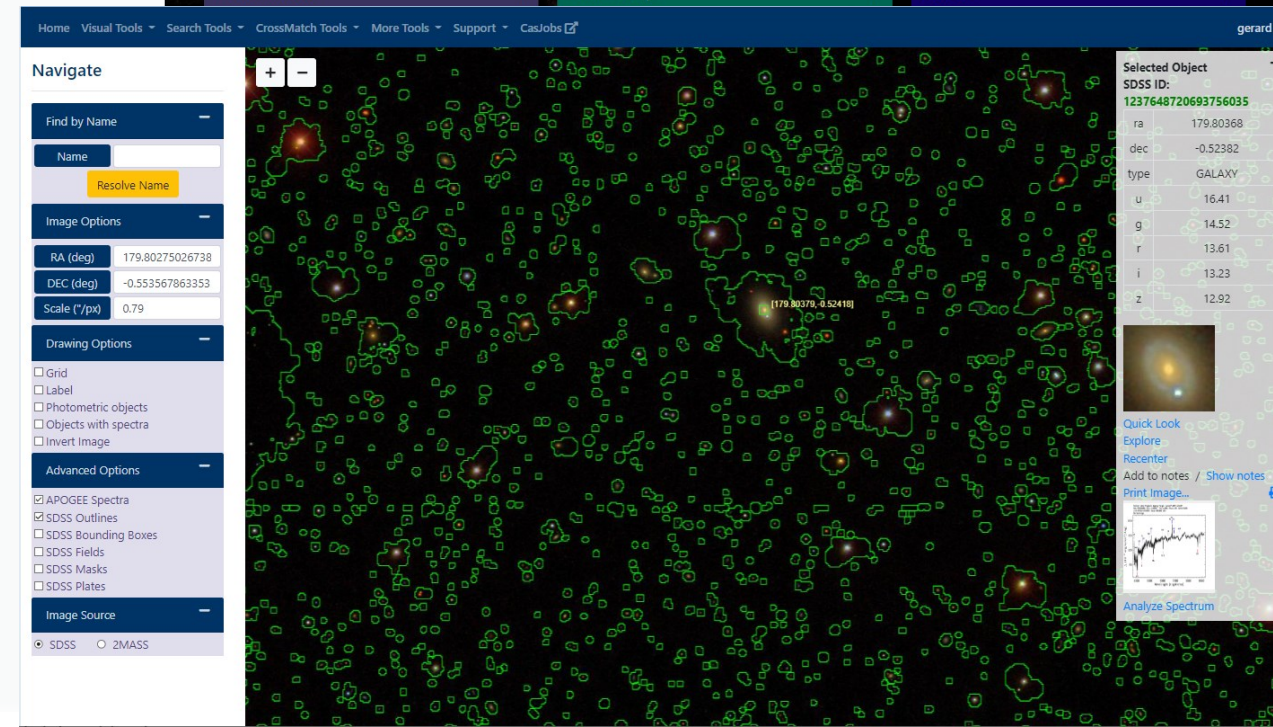
At JHU/IDIES we started with SkyServer
Disseminating data from Sloan Digital Sky Survey
Goal: instant access to rich content
Idea: bring the analysis to the data
Interactive access at the core



SLOAN DIGITAL SKY SURVEY: EARLY DATA RELEASE

Stoughton et al, 2002

<https://ui.adsabs.harvard.edu/abs/2002AJ....123..485S>



But need for longer/deeper searches...

SDSS Query / CasJobs

MyDB **Local Only**

PMStar
Contains ~21,157,643 rows (~2,069,208 kB)

Table Schema type [size]

objID	night	obsID	ccd	brightStar	match	pmRa
bigint [8]	int [4]	smallint [2]	tinyint [1]	tinyint [1]	smallint [2]	real [4]

Notes

ConeSearchClient 1/1/2

PMStar

Proper motion of a star (point source) in the Deep PM catalog.

Each SDSS DR7 star that falls within the survey area, and within a given r magnitude range, is listed in SDSS objects in the magnitude range $16 \leq r \leq 23$, while for the 1.3m survey we include objects in the limits are determined by the onset of saturation for stars, and the faint limits are set to one magnitude limits. There is a table entry for each SDSS object, whether or not it was detected in our survey. Unmat 'match' set to 0. Objects in our survey are matched to SDSS objects by searching in annuli using progr is found within an annulus, the nearest match is used.

name	unit	ucd enum	description
objID			Unique SDSS identifier for the SDSS object.
night			MJD number of the night the observation was obtained.
obsID			Observation number.
ccd			CCD on which the object was or should have been detected.
brightStar			1 if objects falls in a masked region around a bright star, else 0.
match			Multiplicity of matches between our survey and SDSS. One number of matching objects in SDSS, and 10s digit indicates matching objects in this survey. Thus, 11 indicates a one-to-one match. If there is no match for this SDSS object in our survey, columns for this object will be set to 0.

SDSS Query / CasJobs

Context Table (optional) Task Name
DR10 MyTable My Query

Samples Recent Clear [1 s] Query complete! Syntax Plan Quick Submit

```

1 select top 40 p.objid,p.ra,p.dec,p.u,p.g,p.r,p.i,p.z
2 from photoobj p
3 join dbo.fgetnearbyobjeq(44.41, 5.99, 40) n on p.objid=n.objid
4 where p.g between 14 and 18
    
```

40 row(s)

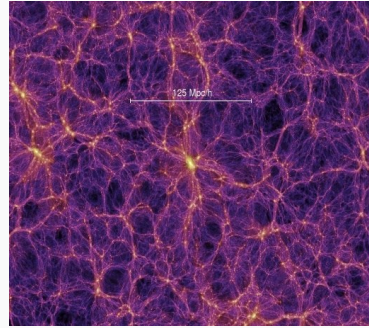
objid	ra	dec	u	g	r	i	z
1237667228764668018	44.4251887464174	5.99796533218519	17.45296	16.15438	15.62593	15.41414	15.32727
1237667228764668016	44.4136115331179	6.01166729833912	15.93901	14.92868	14.99346	14.86291	12.71118
1237667228764668020	44.3902630919155	5.97695081867117	16.96872	14.88728	13.86288	13.33199	12.91775
1237667228764668021	44.4400546420409	5.97533084978952	18.35162	16.20482	15.22672	14.71813	14.34378
1237667228764668022	44.4371469033969	5.98117747486914	19.06305	17.71112	17.12041	16.85705	16.74301
1237670016198443013	44.4085082530024	6.02169136812494	16.27055	15.73507	15.74194	15.6974	13.71063
1237670016198443017	44.4235465056448	6.0224376643852	16.90246	14.66859	13.60513	13.04366	12.63324
1237667228764668003	44.3721225178138	6.00777066389459	16.3297	14.95203	14.42902	16.75699	14.15577
1237667228764668019	44.4514249394326	5.9876249662911	15.91361	15.65941	15.36141	14.96405	13.73247
1237667228764668024	44.4465743657887	5.9847330396699	18.88211	16.72498	15.74462	15.21885	14.93993
1237667228764668244	44.4402786447451	5.95252910090616	16.85921	15.68296	15.02576	14.59115	14.3564
1237670016198443015	44.423185823296	6.02691535484237	16.87503	14.63161	13.5211	12.95495	12.57853
1237667228764667924	44.3467577913547	5.97438474657849	18.82665	16.70258	15.79071	15.4536	15.26757
1237667228764667925	44.3493849949577	5.98186248680855	18.20589	16.3327	15.39307	14.91184	14.54276

RESULTS Plot Save As HTML

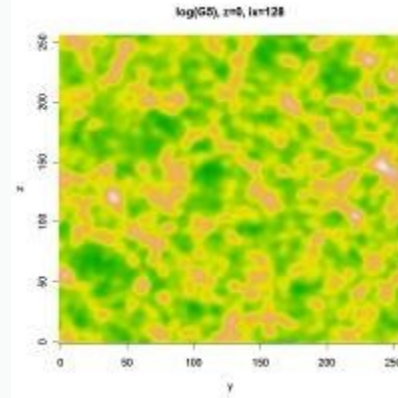
Contact v3_8_6_a

At MPA: Virgo-Millennium Database

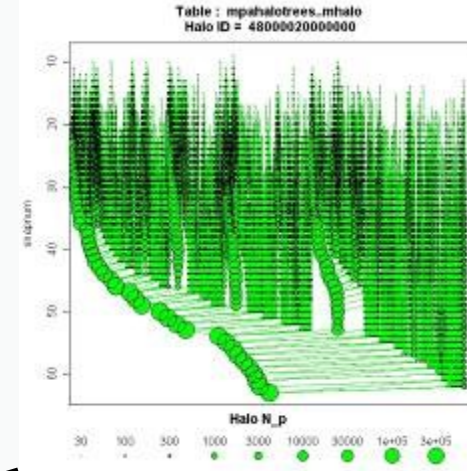
Raw data:
Particles



Density fields

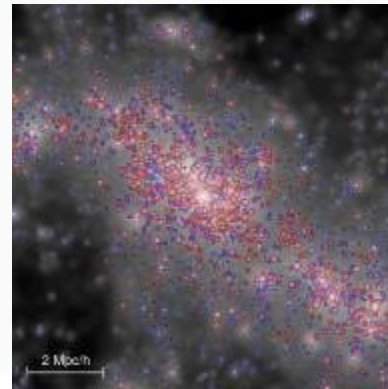
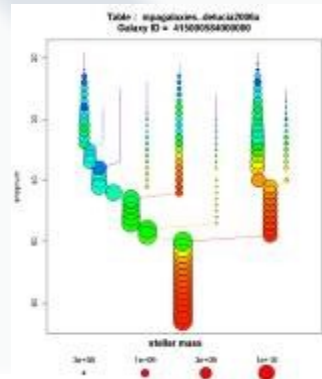


FOF groups
and Subhalos



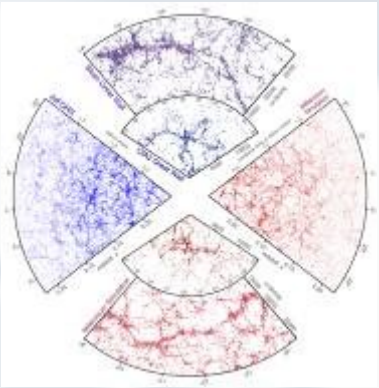
Subhalo merger trees

Synthetic galaxies (SAM)



<http://gavo.mpa-garching.mpg.de/Millennium/>

Mock catalogues



Halo and Galaxy Formation Histories from the Millennium Simulation

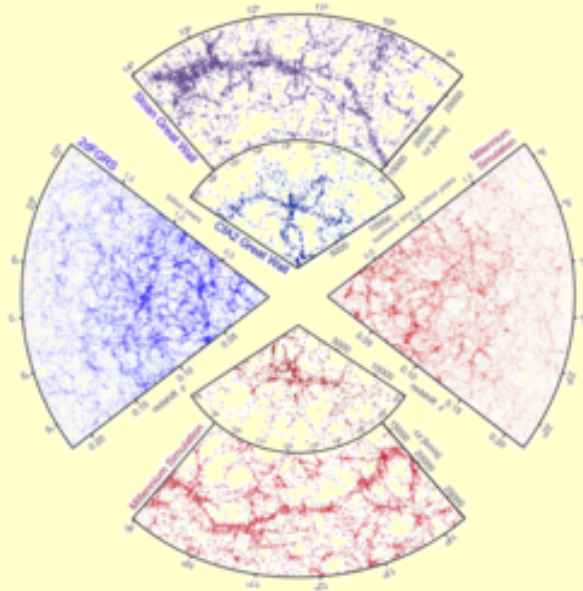
Public release of a VO-oriented and SQL-queryable database for studying the evolution of galaxies in the Λ CDM cosmogony

Gerard Lemson & the Virgo Consortium

[astro-ph/0608019](https://arxiv.org/abs/astro-ph/0608019)

[full description of release \(PDF\)](#)

[database mirror site at ICC, Durham University](#)



- ▶ [Database Access](#)
- ▶ [Visual Material](#)
- ▶ [Related Links](#)
- ▶ [Publications](#)

INTRODUCTION

When published in 2005, the Millennium Run was the largest ever simulation of the formation of structure within the Λ CDM cosmology. It uses 10^{10} particles to follow the dark matter distribution in a cubic region $500h^{-1}\text{Mpc}$ on a side, and has a spatial resolution of $5h^{-1}\text{kpc}$. Application of simplified modelling techniques to the stored output of this calculation allows the formation and evolution of the

<https://wwwmpa.mpa-garching.mpg.de/millennium/>

- Documentation
- CREDITS/Acknowledgments
- Registration
- News
- FAQ
- Public Databases
 - Ayromlou2021a
 - DGalaxies
 - DHaloTrees
 - Guo2010a
 - Guo2013a
 - Henriques2012a
 - Tables
 - wmap1.BC03_0ij
 - wmap1.BC03_AllSky_0ij
 - wmap1.M05_0ij
 - wmap1.M05_AllSky_0ij
 - wmap1_rest.BC03_0ij
 - wmap1_rest.M05_0ij
 - Henriques2015a
 - Tables
 - cones.AllSky_M05_001
 - cones.AllSky_M05_002
 - cones.MRscPlanck1_BC03_0ij
 - cones.MRscPlanck1_M05_0ij
 - MRIscPlanck1
 - MRscPlanck1
 - SFH_Times_MR
 - SFH_Times_MR11
 - Henriques2020a
 - Tables
 - MRscPlanck1
 - MRscPlanck1_Rings
 - MField
 - MillenniumII
 - millimil
 - miniMilII
 - MMSnapshots
 - MPAGalaxies
 - MPAHaloTrees
 - Tables
 - MHalo
 - MR
 - MR7
 - MR11
 - MRIscPlanck1
 - MRIscWMAP7
 - MRscPlanck1
 - MRscWMAP7
 - MPAMocks
 - Snapshots
- Private (MyDB) Databases

Welcome Gerard Lemson.
 Streaming queries return unlimited number of rows in CSV format and are cancelled after 420 seconds.
 Browser queries return maximum of 1000 rows in HTML format and are cancelled after 30 seconds.

There is a partial mirror of this database in Durham at <http://galaxy-catalogue.dur.ac.uk:8080/Millennium/>.
 The Durham database does not contain all the latest L-Galaxies models but does contain more recent GALFORM models.

```
select PROG.*
  from millimil..DeLucia2006a PROG,
       millimil..DeLucia2006a DES
 where DES.galaxyId = 1
       and PROG.galaxyId between DES.galaxyId and DES.lastprogenitorId
```

- Query (stream)
- Query (browser)
- Explain
- Help

Maximum number of rows to return to the query form:

Previous queries:
 List of all queries executed sofar in this session. Selecting a query will make it appear in the query window.
 The link will show all of them in a separate window. Refreshing that window will load the latest queries again.

[Show all previous queries](#)

Demo queries: click a button and the query will show in the query window.
 Holding the mouse over the button will give a short explanation of the goal of the query. These queries are described in some more detail on [this page](#).

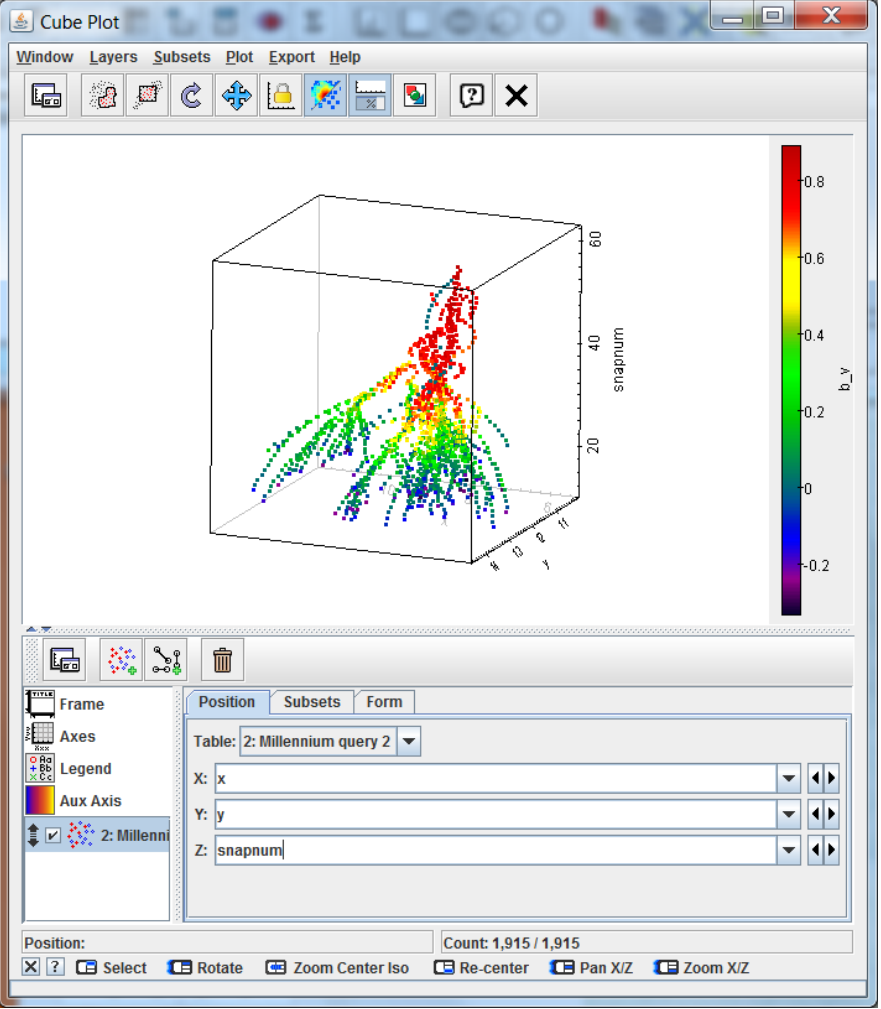
- Mainly Halos:
- Mainly Galaxies:

Metadata queries: The SQL statements under these buttons provide examples for querying and managing the state of a private database.
 Holding the mouse over the button will give a short explanation of the goal of the statement.

-
-
-
-
-
-
-
-
-

Query time (in millisec) = 181
 Number of rows retrieved from database = 1000

galaxyID	lastProgenitorId	descendantId	haloID	subHaloID	fofID	treeld	firstProgenitorId	nextProgenitorId	type	snapnum	redshift	centralMvir	phkey	x	y	z
1	1915	0	1	6200000000000000	6200000000000000	0	2	-1	0	62	0.019932542	3552.877	673	6.587909	13.099106	25
2	1889	1	2	6100000000000000	6100000000000000	0	3	1870	0	61	0.041403083	3494.8687	673	6.597178	13.111782	25
3	1400	2	3	6000000000000000	6000000000000000	0	4	1401	0	60	0.064493395	3318.8645	673	6.615912	13.121013	25
4	1337	3	4	5900000000000000	5900000000000000	0	5	1338	0	59	0.08928783	3285.643	673	6.6276503	13.1303835	25
5	1337	4	5	5800000000000000	5800000000000000	0	6	-1	0	58	0.11588337	3249.5815	673	6.6414022	13.1400175	25



<http://gavo.mpa-garching.mpg.de/Millennium/>

distribution through galaxy-galaxy lensing in the Hyper Suprime-Cam survey

Authors: [Wenting Wang](#), [Xiangchong Li](#), [Jingjing Shi](#), [Jiaxin Han](#), [Naoki Yasuda](#), [Yipeng Jing](#), [Surhud More](#), [Masahiro Takada](#), [Hironao Miyatake](#), [Atsushi J. Nishizawa](#)

Comments: submitted to ApJ - comments welcome - data available upon request

1088. [astro-ph/2104.07664](#) [[abs](#), [ps](#), [pdf](#), [other](#)]:

Title: Optimizing high redshift galaxy surveys for environmental information

Authors: [Tobias J. Looser](#), [Simon J. Lilly](#), [Larry P. T. Sin](#), [Bruno M. B. Henriques](#), [Roberto Maiolino](#), [Michele Cirasuolo](#)

Comments: 29 pages, 40 figures. Accepted for publication in MNRAS

1089. [astro-ph/2104.08295](#) [[abs](#), [ps](#), [pdf](#), [other](#)]:

Title: The evolution of the mass-metallicity relations from the VANDELS survey and the GAEA Semi-Analytic model

Authors: [Fabio Fontanot](#) (1,2), [Antonello Calabro](#) (3), [Margherita Talia](#) (4,5), [Filippo Mannucci](#) (6), [Marco Castellano](#) (3), [Giovanni Cresci](#) (6), [Gabriella De Lucia](#) (1), [Anna Gallazzi](#) (6), [Michaela Hirschmann](#) (7), [Laura Pentericci](#) (3), [Lizhi Xie](#) (8), [Ricardo Amorin](#) (10,11), [Micol Bolzonella](#) (5), [Angela Bongiorno](#) (3), [Olga Cucciati](#) (5), [Fergus Cullen](#) (12), [Johan P. U. Fynbo](#) (13), [Nimish Hathi](#) (14), [Pascale Hibon](#) (15), [Ross J. McLure](#) (12), [Lucia Pozzetti](#) (5) ((1) INAF - Astronomical Observatory of Trieste, Italy (2) IFPU - Institute for Fundamental Physics of the Universe, Trieste, Italy (3) INAF - Astronomical Observatory of Rome, Italy (4) Dipartimento di Fisica e Astronomia, Università di Bologna, Italy (5) INAF - Astronomical Observatory of Bologna, Italy (6) INAF - Astrophysical Observatory of Arcetri, Firenze, Italy (7) DARK, Niels Bohr Institute, University of Copenhagen, Denmark (8) Tianjin Astrophysics Center, Tianjin Normal University, China (10) Instituto de Investigación Multidisciplinar en Ciencia y Tecnología, Universidad de La Serena, Chile (11) Departamento de Física y Astronomía, Universidad de La Serena, Chile (12) SUPA Scottish Universities Physics Alliance, Institute for Astronomy, University of Edinburgh, Royal Observatory (13) The Cosmic Dawn Center, Niels Bohr Institute, University of Copenhagen, Denmark (14) Space Telescope Science Institute, Baltimore, USA (15) ESO-Chile, Santiago, Chile)

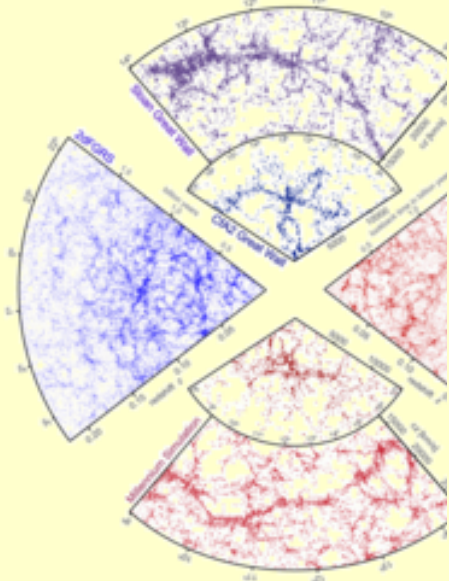
Comments: 12 pages, 8 figures, MNRAS submitted

1090. [astro-ph/2105.09126](#) [[abs](#), [ps](#), [pdf](#), [other](#)]:

Title: Dynamical analysis of clusters of galaxies from cosmological simulations

Authors: [Tania Aguirre Tagliaferro](#), [Andrea Biviano](#), [Gabriella De Lucia](#), [Emiliano Munari](#), [Diego Garcia Lambas](#)

Comments: 14 pages, 17 figures



- ▶ [Database Access](#)
- ▶ [Visual Material](#)
- ▶ [Related Links](#)
- ▶ [Publications](#)

the Millennium

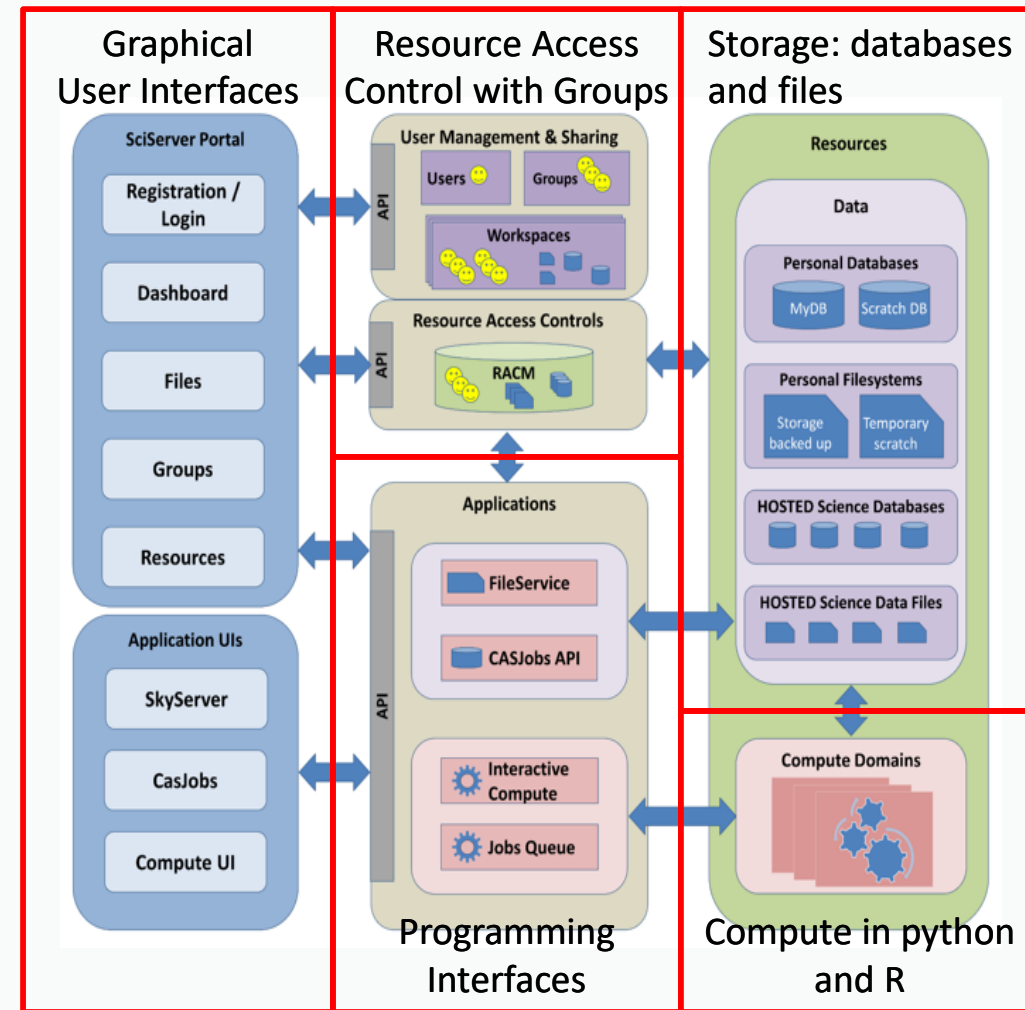
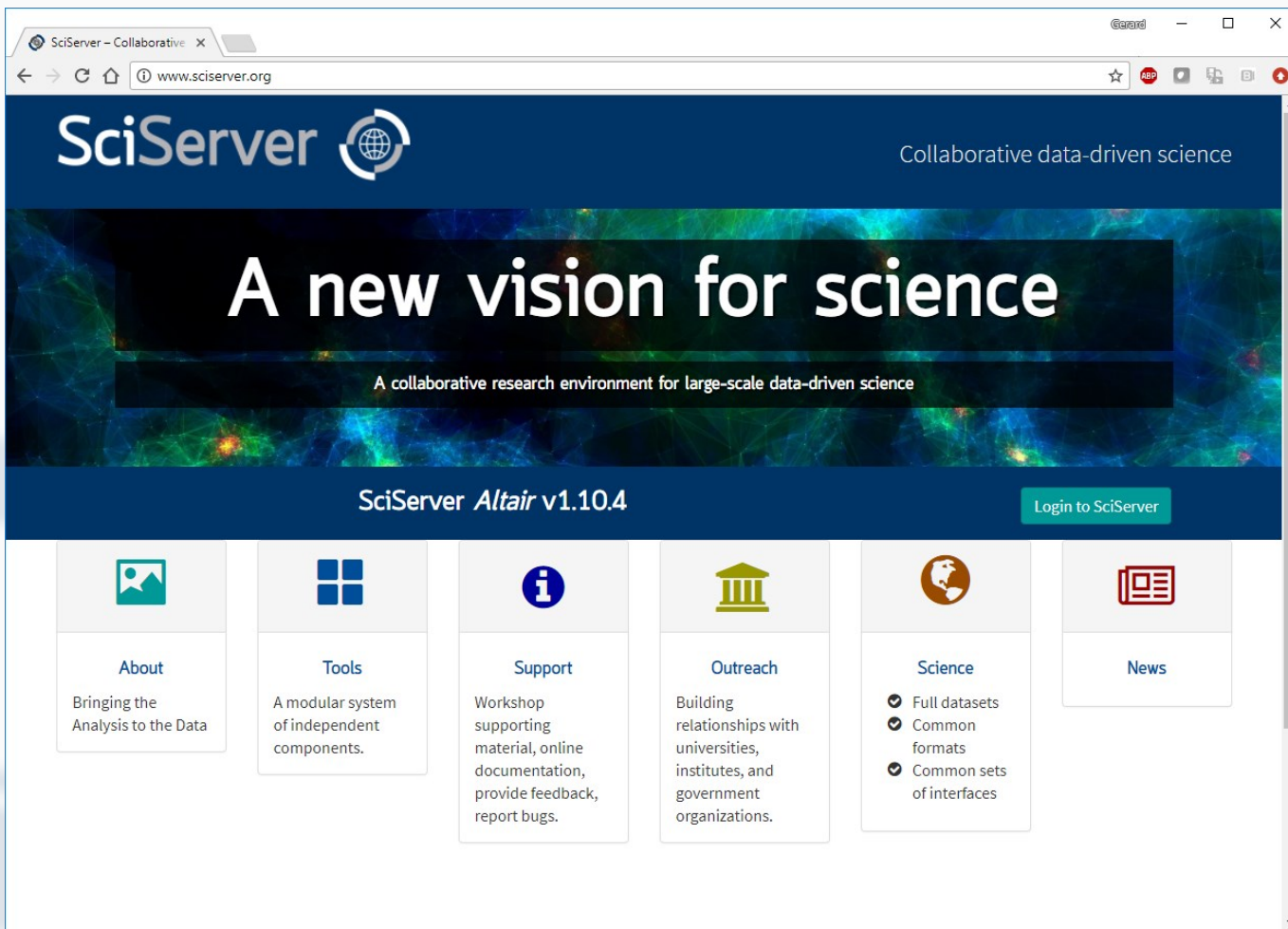
base for studying the evolution of the dark matter distribution in a kpc. Application of simplified formation and evolution of the

ly.

mulation of the formation of
dark matter distribution in a
kpc. Application of simplified
formation and evolution of the

Users want more than “just” SQL:

- ▶ Want full analysis, visualization near the data
 - Run Python, R, C++, etc
- ▶ Want access to data that does not fit in a relational database
 - Images, spectra, data cubes, custom data sets
- ▶ Want to upload own data and combine
 - Need workspace close to data, databases and file system
- ▶ Want to share work with collaborators
 - Data, Scripts, Results
- ▶ Support for data and libraries from different disciplines
 - No single data model, ontologies

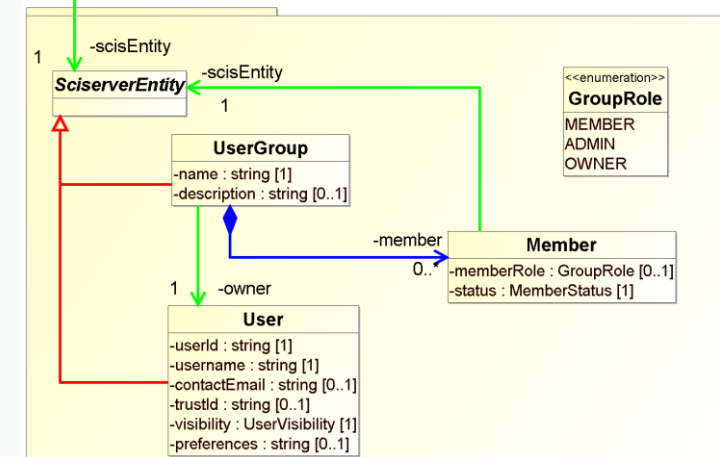
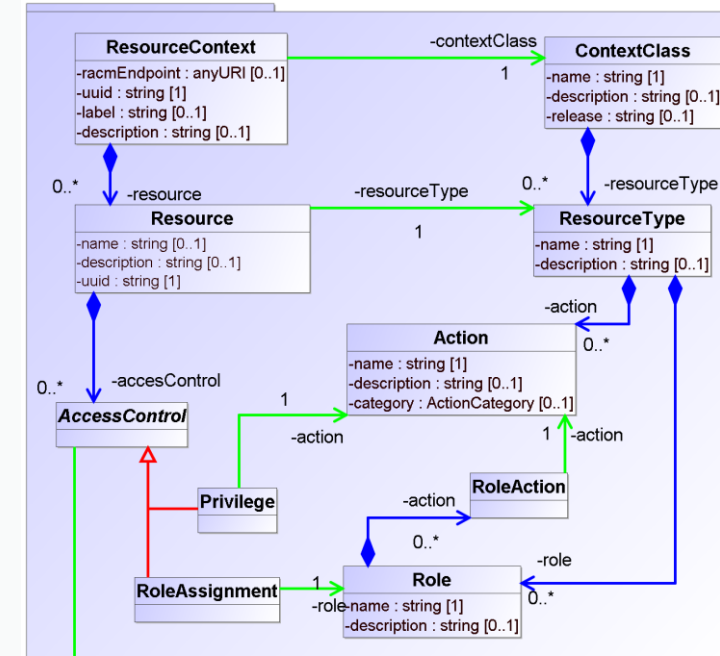


- ▶ **SciServer** is a system allowing Science Researchers across multiple domains to host and share their datasets, and provide query and analysis tools for collaborative research.
- ▶ **Core Services:**
 - Science Data Hosting (Files and Databases)
 - Query of hosted databases
 - Data Integration across hosted data sources
 - Computational analysis on hosted data
 - Collaboration and Sharing
 - Personal Storage (Files and Database)
 - API Integration

Resource sharing

The screenshot shows the SciServer interface for the ADASS2019Demo group. It includes a navigation bar with 'Home', 'Files', 'Groups', 'Admin', and 'Resources'. The main content is divided into three sections: 'Groups', 'Shared Files', and 'Members'.

- Groups:** A list of groups with a search bar containing 'as'. The 'ADASS2019Demo' group is selected, showing a 'Leave group' button. Other groups include AS171d205, AS4010_2018_S1, Astroinformatics2018-Students, baseball, Baseball Collaboration, GerardsDatabases, HEASARC software user group, Polybase, TamasClas, and TamasTest.
- Shared Files:** Shows two folders: 'ADASS2019Demo' and 'ADASS2019DemoD ata'.
- Shared Data Volume:** A message: 'Share Data Volume with this group to see them here.'
- Shared Compute Images:** A message: 'Share compute images with this group to see them here.'
- Shared Databases:** A message: 'Share database context with this group to see them here.'
- Members:** A list of users with their roles: jkim485 (OWNER), gerard (ADMIN), adi (MEMBER), bac29 (MEMBER), bfalck (MEMBER), jcg (MEMBER), raddick (MEMBER), and tjaffe (MEMBER). There is an 'Edit Member List' button at the bottom.



SciServer Compute: Jupyter (and more) in docker containers

SciServer Compute

Interactive Notebooks

Jobs



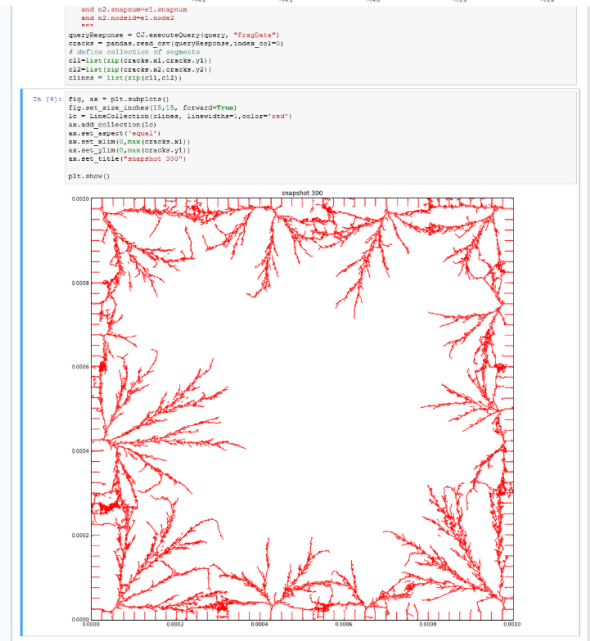
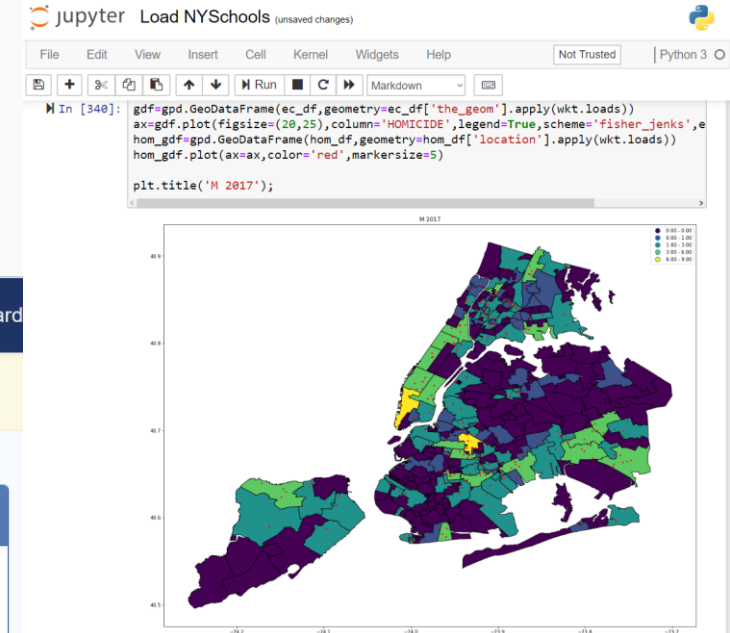
gerard

Now Available: JupyterLab and Classical Jupyter images are now combined. Containers default to the classical interface and will remember the last interface used.

Containers

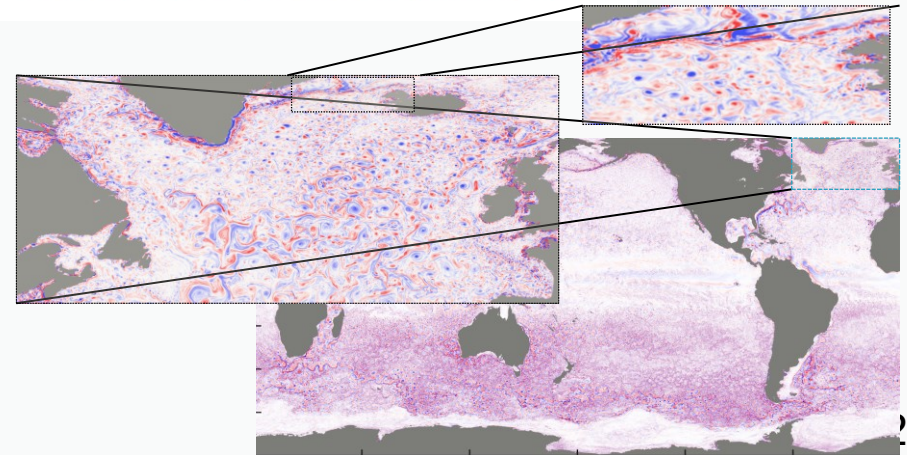
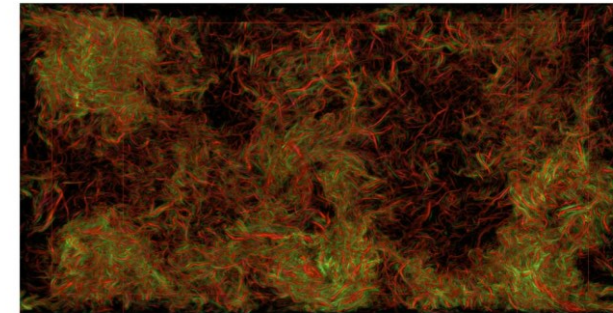
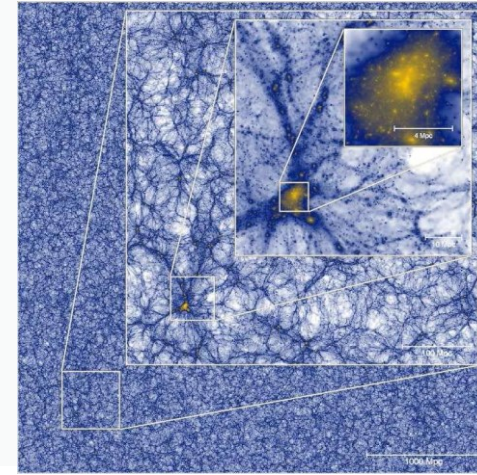
Created At	Name	Domain	Image	Status	
2019-10-06 07:35:21.0	ADASS 2019 interactive GPU domain	GPU Interactive	TensorFlow	running	■ ⓘ ✕
2019-10-02 08:36:11.0	heasoft	Interactive	Heasarc	stopped	▶ ⓘ ✕
2019-10-01 16:24:04.0	montage	Interactive	Montage	stopped	▶ ⓘ ✕
2019-09-30 14:44:39.0	ADASS2019-pytorch	Interactive	PyTorch	running	■ ⓘ ✕
2019-07-09 17:45:08.0	indra (filedb)	Interactive	Python + R	stopped	▶ ⓘ ✕
2019-05-02 10:33:40.0	polybase and aris	Interactive	Python + R	stopped	▶ ⓘ ✕
2019-04-25 16:28:15.0	old english	Interactive	BeakerX	stopped	▶ ⓘ ✕
2019-03-19 11:46:22.0	WFIRST	Interactive	WFIRST-Archive SIT	stopped	▶ ⓘ ✕
2018-10-28 15:51:49.0	Sociology + III + ukbiobank	Interactive	Python + R	stopped	▶ ⓘ ✕

Create container



Simulations at SciServer

- ▶ Cosmology
 - Virgo Millennium suite ~40TB DB, ~60TB raw
 - Indra ~0.8PB
 - ApogeeFire DB ~29TB
 - Jason Hunt 2021 ~100TB
- ▶ Turbulence database, soon >2PB
- ▶ Ocean circulation, soon ~2PB



Science Domains

Cosmological Simulations

Genomics

Cosmological Simulations



Join

Cosmological Simulations

If you join this science domain (by clicking the "join" button), you will have access to close to a peta-byte of data from a number of cosmological simulations.

The theoretical input for galaxy-scale physics and large-scale galaxy redshift surveys is provided by cosmological numerical simulations, which model the evolution of the Milky Way or the entire Universe. SciServer hosts several sets of these simulations, offering unified access and collaborative visualization tools for hundreds of TB of data. The unique advantage of accessing these simulations through SciServer is the ease with which they can be compared with real astronomical observations.

Some of the data sets made accessible in this domain are stored in relational databases and will be accessible through the CasJobs web page and the CasJobs modules in python and R. Others are accessible as files on data volumes. To use these data in SciServer Compute, create a container with the **Cosmological Simulations** compute image and mount the appropriate volumes.

Indra

Indra is a suite of large-volume cosmological N-body simulations, all with the same cosmology and different initial phases. Each simulation has 64 snapshots of particle data and halo catalogs and 505 time steps of Fourier-space density fields. The Indra data volume contains the full 750 TB dataset, and the Indra database contains the halo catalogs plus the Spatial3d library of search tools.

Notebooks giving examples of how to access Indra data are found in the Indra data volume and the [indra-tools git repository](#).

- For an overview see: <https://www.sciserver.org/datasets>
- For detailed information see the data release paper: <https://arxiv.org/abs/2101.03631>
- Contact: Bridget Falck (bridget.falck@jhu.edu)

Virgo Consortium, Millennium suite of simulations

This science domain also gives access to a variety of simulations and related products produced by the [Virgo Consortium for Cosmological Supercomputer Simulations](#). In particular SciServer contains a copy of the [Millennium Database](#), which contains halo catalogues and catalogues of galaxies simulated using semi-analytical galaxy formation algorithms.

SciServer extends the functionality of that site by providing online analysis capabilities in Jupyter notebooks. But in particular SciServer allows user to access the underlying raw simulation data, i.e. particles for the N-body simulations. Notebooks in the [Getting Started](#) data volume give examples how to access these.

See also:

- For more details of the Millennium simulations and their usage by the community see <https://wwwmpa.mpa-garching.mpg.de/millennium/>
- For more details of the Millennium Database please see the original [web site](#).

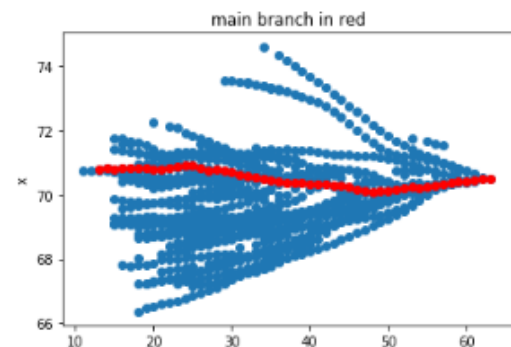
References

- Millennium Simulation: [Springel et al 2005](#)
- Millennium Database: [Lemson & the Virgo Consortium 2006](#)

```
In [1]: import SciServer.CasJobs as cj
import numpy as np
import matplotlib.pyplot as plt

In [2]: # select a random halo at z=0 (snapnum=63) and mass between 500 and 5000x10^10 Msun/h
# then find merger tree rooted in that halo
sql="""
with descendants as (
select top 1 haloid , lastprogenitorid, mainleafid, m_crit200, d.np
from MR d
where snapnum=63 -- reshift=0
and firsthaloinfofgroupid=d.haloid -- centers of FOF groups
and m_crit200 between 500 and 5000
order by newid() -- random sub sample
)
select d.haloid as d_id, d.m_crit200 as d_mass, d.np as d_np
, p.haloId,p.x,p.y,p.z,p.snapnum,p.mainLeafId,p.halfMassRadius
from descendants d
inner join MR p
on p.haloid between d.haloid and d.lastprogenitorid
order by d.haloid, p.haloid
"""
tree=cj.executeQuery(sql,"MPAHaloTrees") # send query to the MPAHaloTrees database context
```

```
In [4]: X='snapnum'
Y='x'
i=0
#for d_id,tree in gdf:
i1,i2=tree[['haloid','mainLeafId']].iloc[0]
tree.set_index('haloid',inplace=True)
main=tree.loc[i1:i2]
plt.scatter(tree[X],tree[Y])
plt.scatter(main[X],main[Y],color='red')
plt.xlabel(X)
plt.ylabel(Y)
plt.title('main branch in red');
```



Raw data, accessible from inside and outside the database

jupyter

Select items to perform actions on them.

0 / virgo

- ..
- Eagle
- EagleFITS
- Illustris
- Millennium**
- Millennium2
- millimil
- MRObs
- ScaleFree
- sdss_ml
- yt_samples

jupyter

- snapdir_032
- snapdir_033
- snapdir_034
- snapdir_035
- snapdir_036
- snapdir_037
- snapdir_038
- snapdir_039
- snapdir_040**
- snapdir_041
- snapdir_042

a year ago

jupyter

snap_millennium_040.173	a year ago	569 MB
snap_millennium_040.174	a year ago	666 MB
snap_millennium_040.175	a year ago	529 MB
snap_millennium_040.176	a year ago	564 MB
snap_millennium_040.177	a year ago	572 MB
snap_millennium_040.178	a year ago	525 MB
snap_millennium_040.179	a year ago	597 MB
snap_millennium_040.18	a year ago	644 MB
snap_millennium_040.180	a year ago	649 MB
snap_millennium_040.181	a year ago	639 MB
snap_millennium_040.182	a year ago	766 MB
snap_millennium_040.183	a year ago	622 MB
snap_millennium_040.184	a year ago	595 MB
snap_millennium_040.185	a year ago	595 MB

jupyter CosmoUC_5_PlotHalo Last Checkpoint: 11 minutes ago (autosaved)

File Edit View Insert Cell Kernel Help Notebook saved

Code CellToolbar

Dark-matter Halos in a Cosmological Simulation

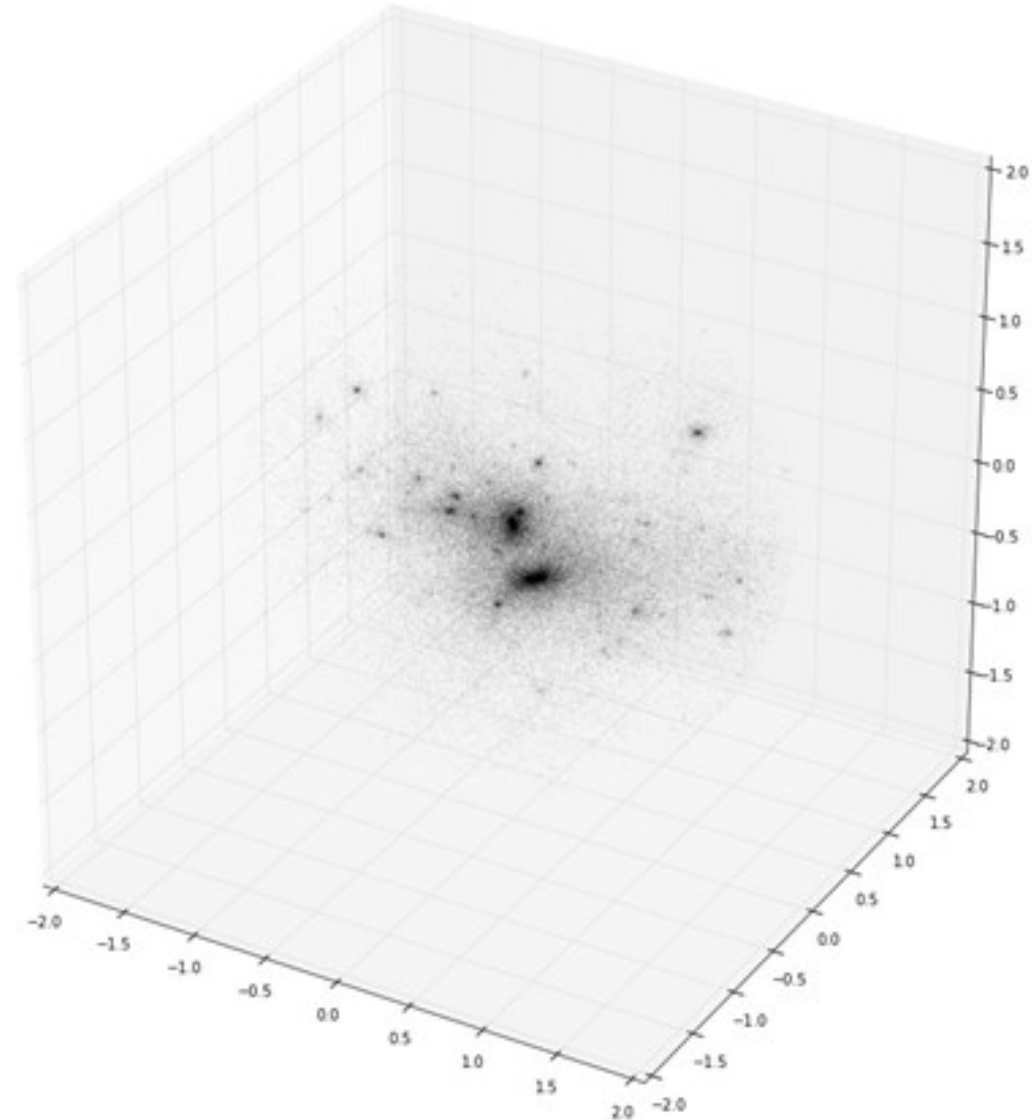
```
In [3]: import SciServer.LoginPortal as Login
token = Login.getToken()
import SciServer.CasJobs
import pandas
import tables
import numpy as np
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
```

```
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 mpahalotrees.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=84000007000000 order by newid()
"""
responseStream = SciServer.CasJobs.executeQuery(queryString, token=token,context="SimulationDB")
df = pandas.read_csv(responseStream, index_col=None)

CPU times: user 351 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>
```



jupyter CosmoUC_5_PlotHalo Last Checkpoint: 11 minutes ago (autosaved)

File Edit View Insert Cell Kernel Help Notebook saved

Code CellToolbar

Dark-matter Halos i

```
In [3]: import SciServer.LoginPortal
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

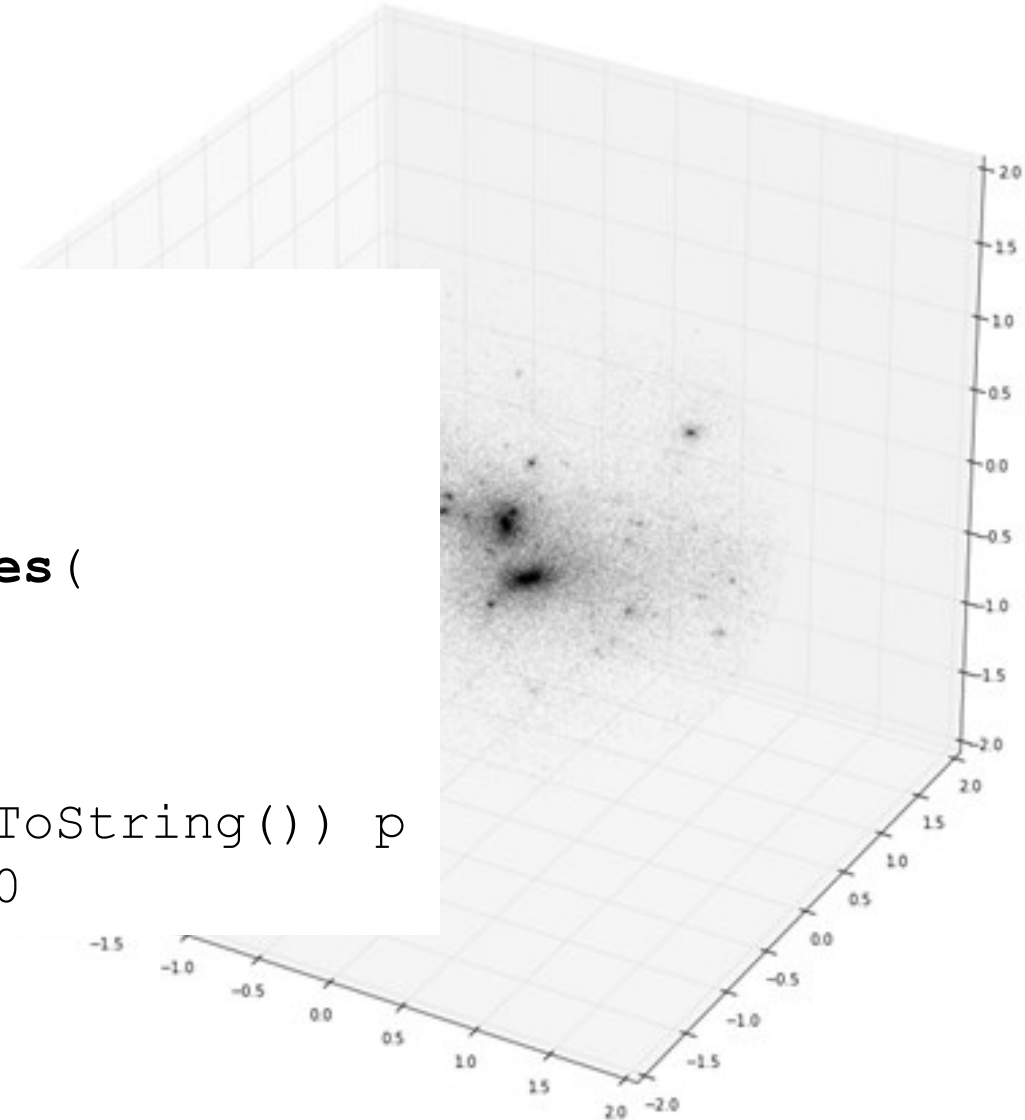
In [10]: %%time
queryString = """
select top 100000 p.x-hh.x as x
,
p.y-hh.y as y
,
p.z-hh.z as z, hh.np
from mpahalotrees.mr hh
cross apply MillenniumParticles(
hh.snapnum,
dbo.Sphere::New(
hh.x, hh.y, hh.z,
3*hh.halfmassradius).ToString()) p
where hh.haloid=8400000700000000

responseStream = SciServer.CasJobs.GetResponse(queryString)
df = pandas.read_csv(responseStream)

CPU times: user 351 ms, sys: 0 ms, total: 351 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>
```



Indra: a public computationally accessible suite of cosmological N-body simulations

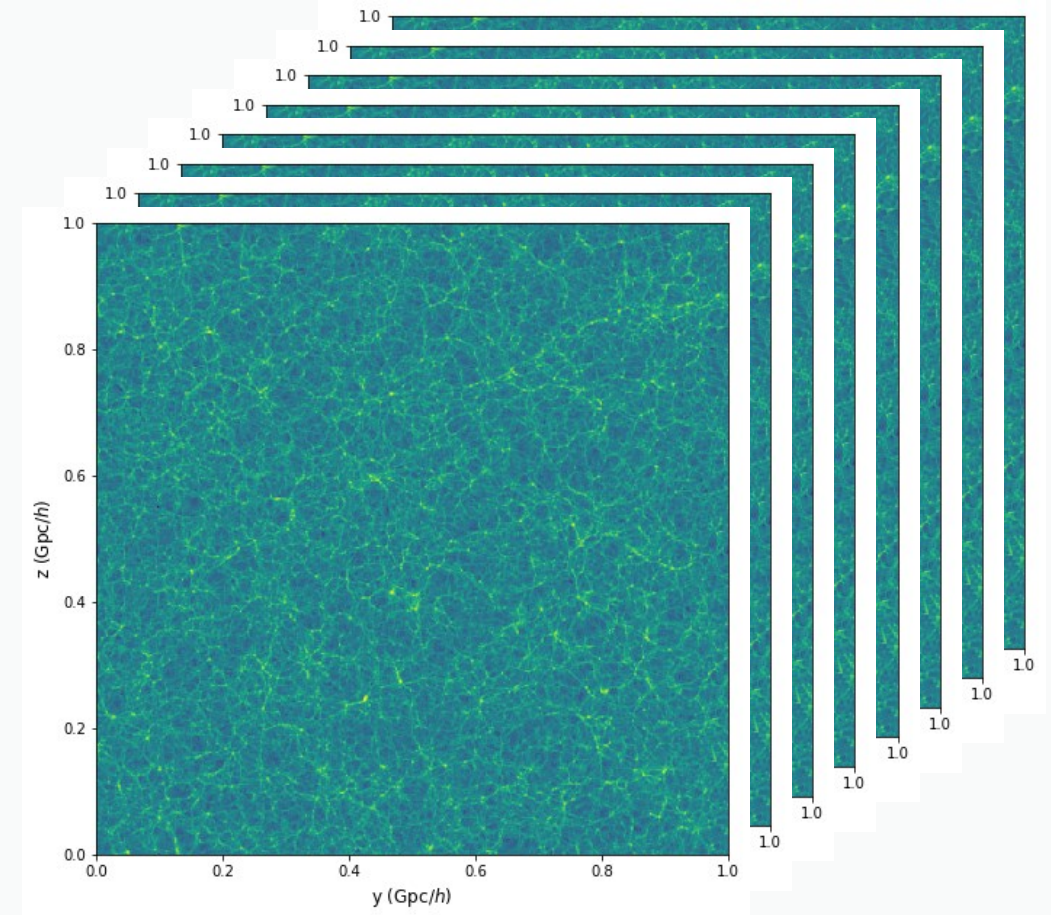
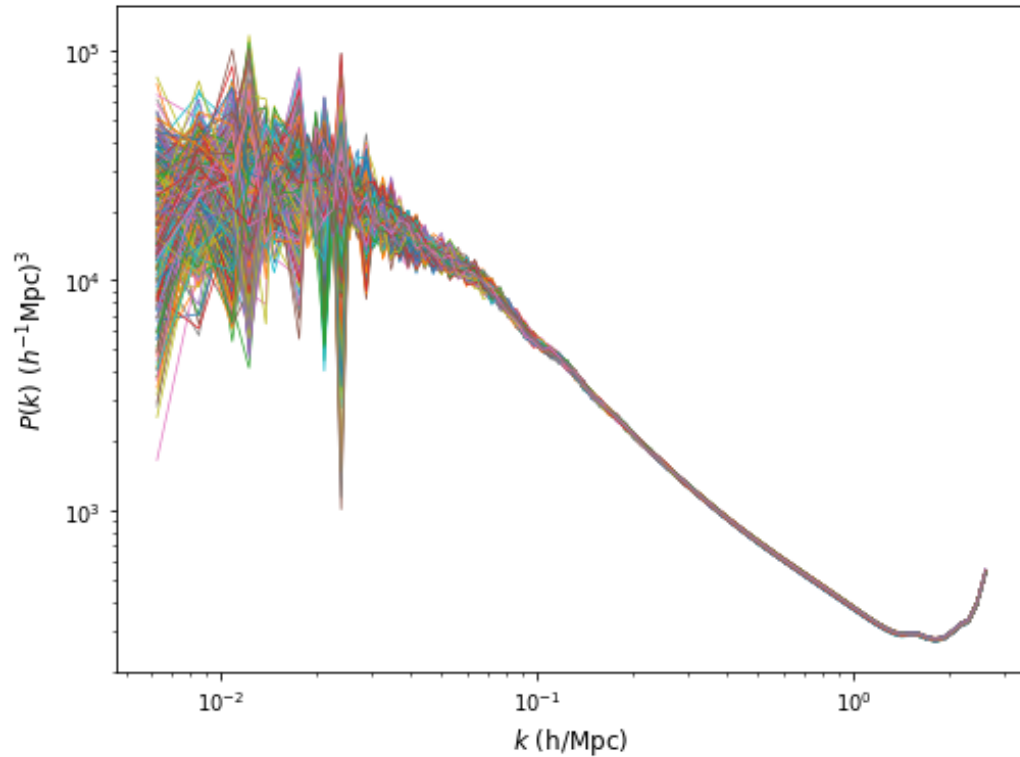
Falck et al 2021, <https://ui.adsabs.harvard.edu/abs/2021MNRAS.506.2659F/abstract>

Power spectra with DASK

448 simulations, 1 billion particles each.

Same initial conditions, different seed.

448 Cloud-In-Cell density grids created and FFT-ed in 2 hours, using 8 DASK workers on distributed file system

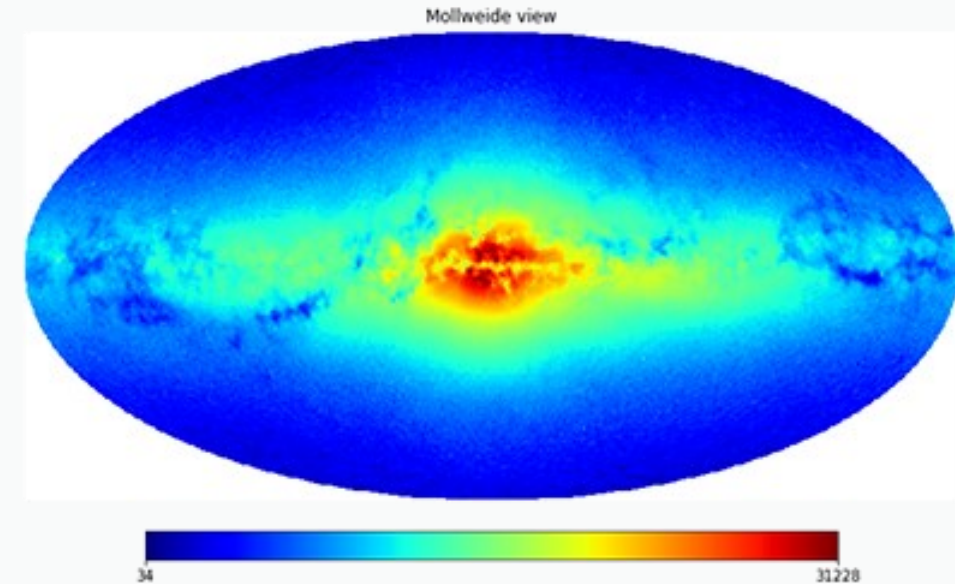


APOGEE-centric Ananke Simulations in a SciServer SQL Database

Beaton et al 2022, <https://ui.adsabs.harvard.edu/abs/2022RNAAS...6..125B>

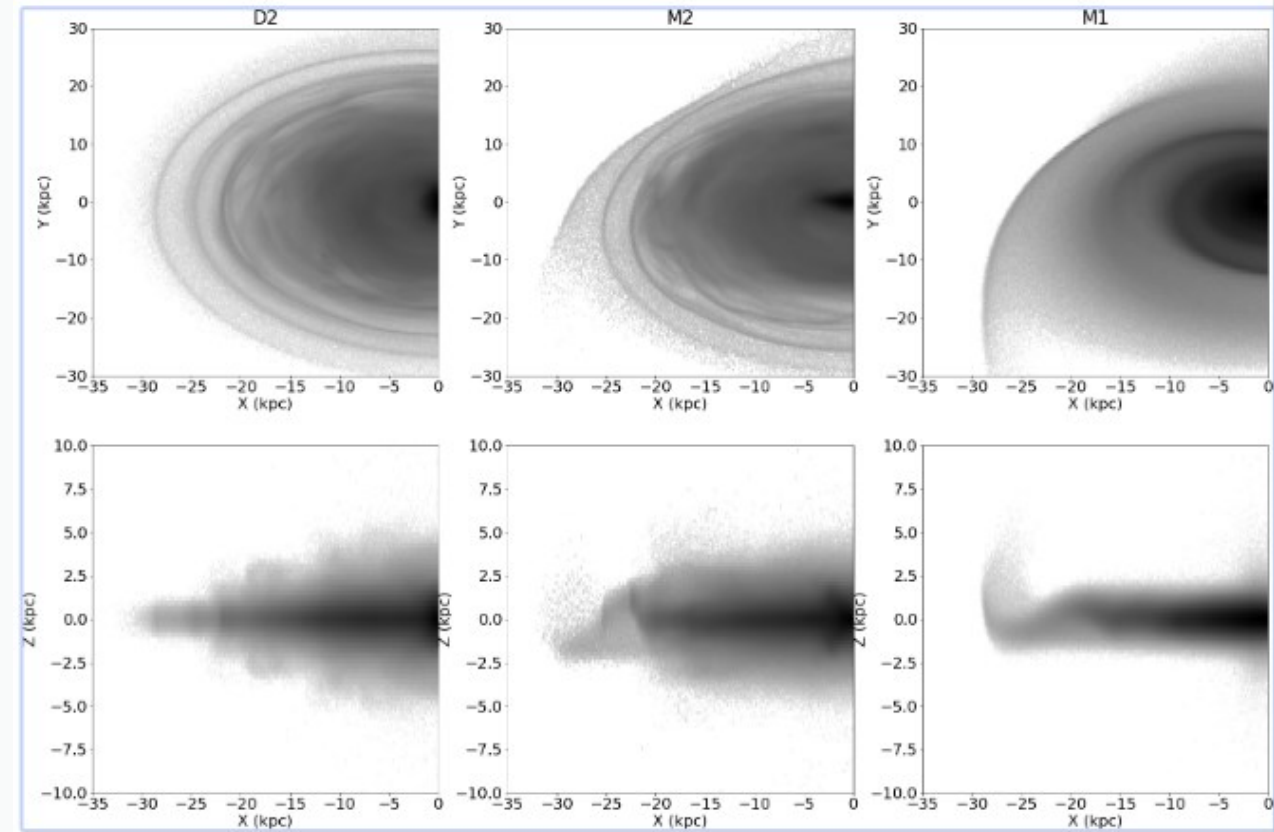
▶ ELT process on Spark:

- Globus transfer of data
- Spark/HDFS used for transforming FITS→parquet
- Addition of special columns:
 - htm20Id, heal20Id, positions in different coordinate systems, ...
- MS libraries used for loading into RDB column store
- Registration in SciServer for access through CasJobs and from within Jupyter notebooks



Resolving local and global kinematic signatures of satellite mergers with billion particle simulations Hunt et al 2021, <https://ui.adsabs.harvard.edu/abs/2021MNRAS.508.1459H/abstract>

- ▶ Data transfer through Globus tasks scripted in Jupyter notebooks on SciServer to distribute files over 12x3 storage volumes of *FileDB cluster*
- ▶ Analysis on sciserver in progress (C. Fillion)



VNC KDE compute image: linux desktop

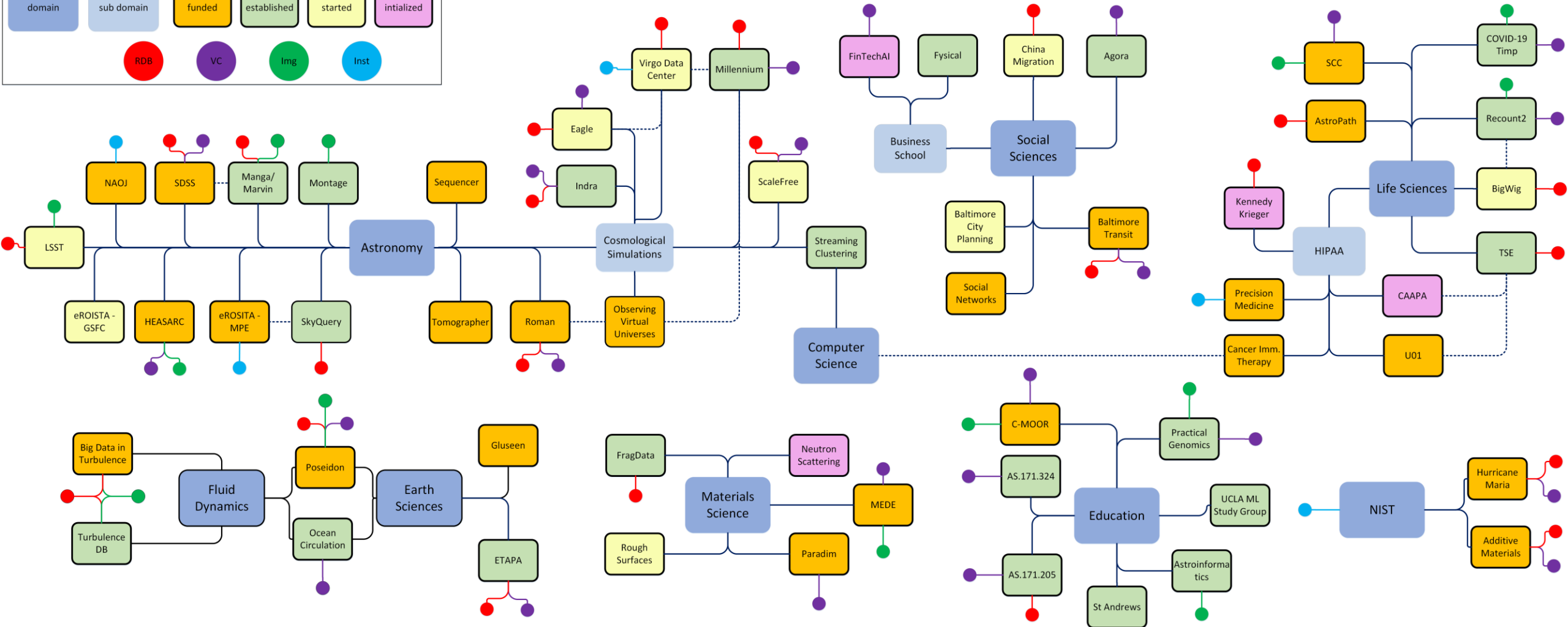
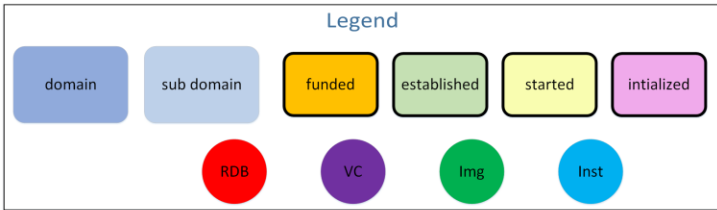
The screenshot displays a Linux desktop environment with the TOPCAT software running. The interface includes a menu bar (File, Views, Graphics, Joins, Windows, Interop, Help), a toolbar, and several panels:

- Table List:** Shows a list of tables, including '1: Millennium query 1' and '3: 3d_halo.csv'.
- Current Table Properties:** Displays details for the selected table, such as 'Label: 3d_halo.csv', 'Location: /home/idies/w...', 'Name:', 'Rows: 100,000', 'Columns: 5', 'Sort Order:', 'Row Subset: All', and 'Activation Actions: 0 / 1'.
- Table Parameters for 1: Millennium query 1:** A table with columns 'Name' and 'Value'. The 'SQL' field contains the query:

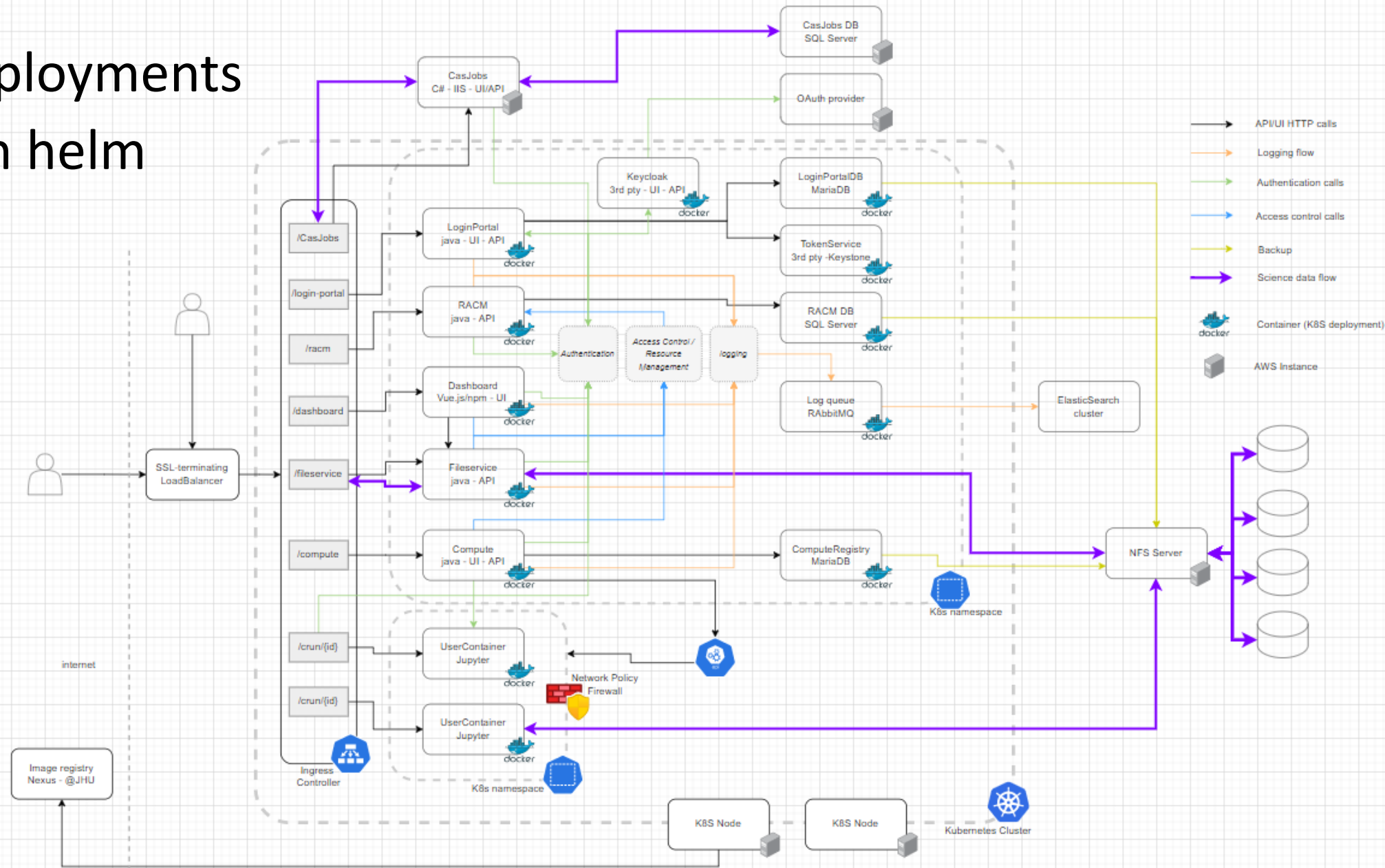

```
select PROG.* from millimil..DeLucia2006a PROG, millimil..DeLucia2006a DES where DES.galaxyId = 1 and PROG.galaxyId between DES.galaxyId and DES.la
```
- 3D Plot (3):** A 3D scatter plot showing a dense cloud of points in a purple-to-red color gradient. The axes are labeled x, y, and z, with a color scale on the right ranging from 0 to 2.0.
- 3D Plot (1):** A 3D scatter plot showing a sparse, branching structure of red points. The axes are labeled x, y, and z, with a color scale on the right ranging from 0 to 2.0.
- Form and Subsets Panels:** These panels allow for configuring the appearance and data subsets of the plots. The 'Form' panel includes options for 'Shading' (Mode: weighted, Weight, Combine: mean) and 'Global Style'. The 'Subsets' panel shows '1: Miller' as an active subset.

The desktop taskbar at the bottom shows several open windows: 'jars : java — Konsole', 'TOPCAT', 'Cube Plot (1)', 'TOPCAT(1): Table Parameters', and 'Cube Plot (3)'. The system tray on the right indicates the time as 12:29.

Beyond astrophysics only




- ▶ External deployments
- ▶ On K8S with helm




SciServer Dashboard

Data, Collaboration, Compute


Your Activities



Files
You have 3 Shared User Volumes.
You have 5 Owned User Volumes.




Groups
You have 1 Group Invitation.
You have 2 Owned Groups.




Compute Jobs
You have 0 Jobs Running.
You have 0 Jobs Completed in 24 hours.


SciServer Apps



CasJobs
Search online big relational databases collections, store the results online, and share them.



Compute
Analyze data with interactive Jupyter notebooks in Python, R and MATLAB.




Compute Jobs
Asynchronously run Jupyter notebooks in Python, R and MATLAB or commands.

NIST


crunchr

Home Files Groups Compute Resources Help shandy3


Data . Collaboration . Compute




Files
Upload and access data volumes
3 Owned User Volumes
1 Shared User Volume



Groups
Share files and notebooks with a study team
0 Groups
0 Group Invitations



Compute
Analyze data with interactive Jupyter notebooks




Jobs
Submit non-interactive jobs in Python or R
0 Jobs Running
0 Jobs Recently Completed

PMAP


SciServer Dashboard

Data, Collaboration, Compute


Your Activities



Files
You have 1 Shared User Volume.
You have 2 Owned User Volumes.



Groups
You have 0 Group Invitations.
You have 0 Owned Groups.



Compute Jobs
You have 0 Jobs Running.
You have 0 Jobs Completed in 24 hours.

SciServer Apps



CasJobs
Search online big relational databases collections, store the results online, and share them.



Compute
Analyze data with interactive Jupyter notebooks in Python, R and MATLAB.




Compute Jobs
Asynchronously run Jupyter notebooks in Python, R and MATLAB or commands.

MPE


SciServer Dashboard

Data, Collaboration, Compute


Your Activities



Files
You have 0 Shared User Volumes.
You have 2 Owned User Volumes.




Groups
You have 0 Group Invitations.
You have 0 Owned Groups.




Compute Jobs
You have 0 Jobs Running.
You have 0 Jobs Completed in 24 hours.


SciServer Apps




Compute
Analyze data with interactive Jupyter notebooks in Python.




Online SQL Editor
A multi-featured editor to check and submit queries



Schema Browser
Your database navigator



hscMap
Your navigator around the deep sky



File Search Tool
Looking for flat files?





Image Cutouts
Get postage stamps of your objects



PSF Picker
Want to know PSF around your objects?

NAOJ

SciServer, a Scalable Data Integrator

- ▶ Difficult to aggregate large, geographically distributed data sets:
 - joint analysis requires co-location.
 - SciServer facilitates creating ETL pipelines
- ▶ Joint analysis requires integration: aggregate data from various sources in a common context
 - A Science Platform like SciServer provides such a context
 - It allows users with diverse skills to collaborate on a single data set
- ▶ Most frequent mistake: trying to create the “mother of all databases”
 - Building ontologies and data models is hard
 - We learned an enormous amount during the Virtual Observatory project
- ▶ Real life uses require interactive exploration before big analysis

The SciServer philosophy:

- ▶ Create Data Contexts, each possibly with their own data model and ontology, self documenting
- ▶ These are secure and read-only, under access control
- ▶ User get their own databases/user volumes and resources to create value added results
- ▶ These can be shared at will with authenticated users at owner's discretion
- ▶ We can bring in new datasets in isolation very quickly
- ▶ **Reproducible science on integrated system**

Thank you

Registration is free at

<https://www.sciserver.org>

<https://apps.sciserver.org>

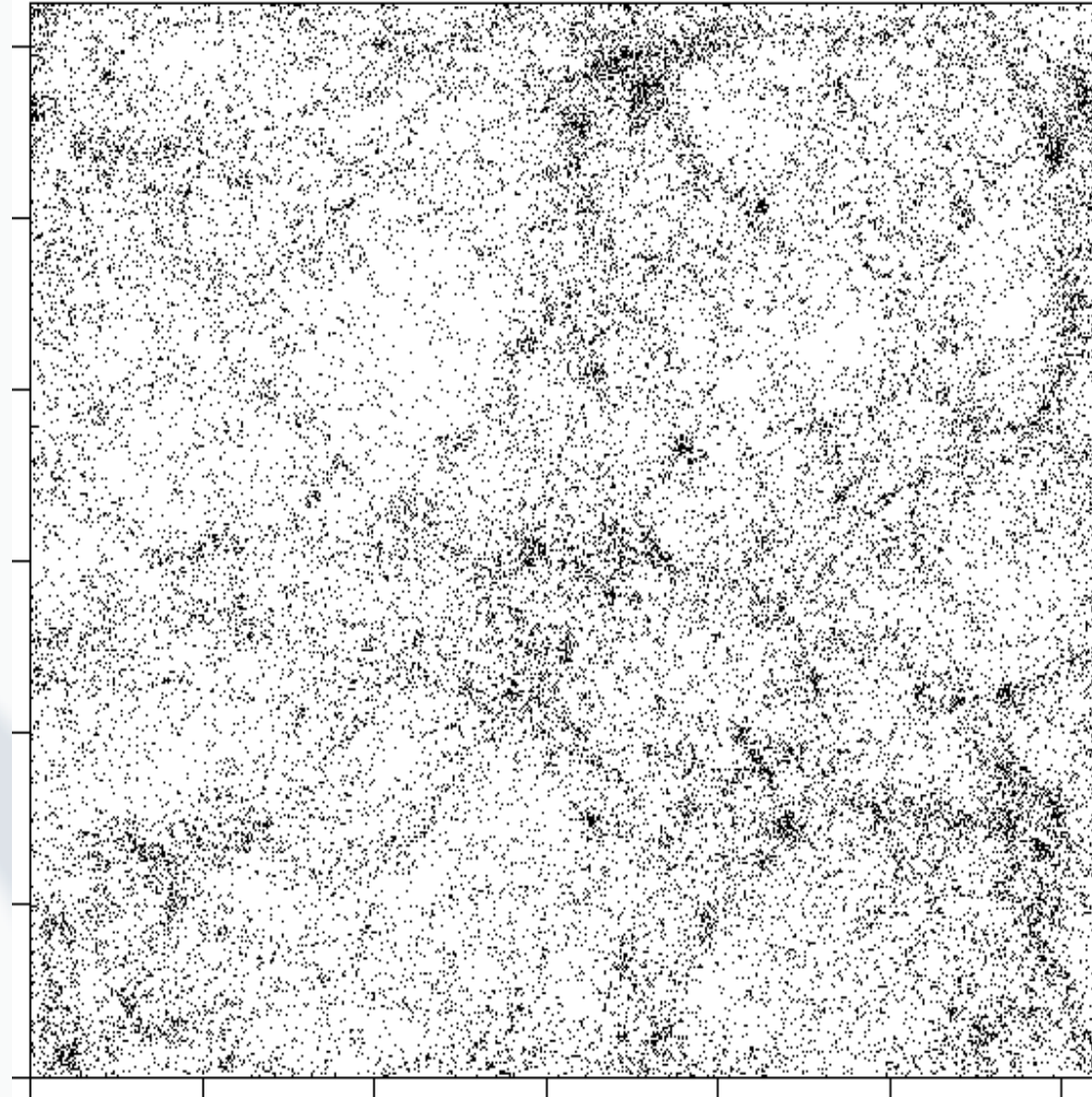
Appendix:

Efficiently querying 3D sub-volumes

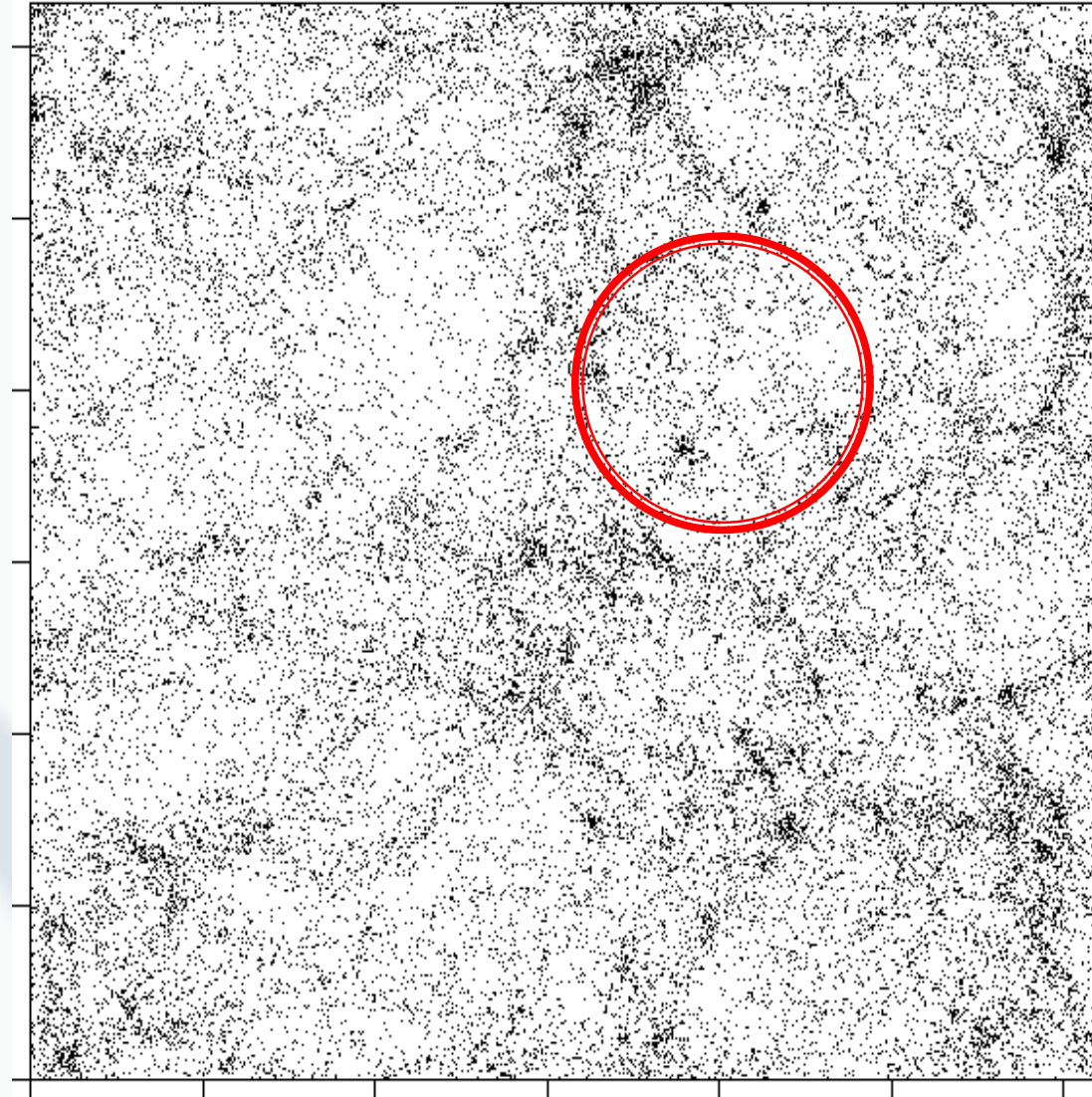
Implementing a General Spatial Indexing Library for Relational Databases of Large Numerical Simulations

Lemson, Budavari & Szalay, 2012

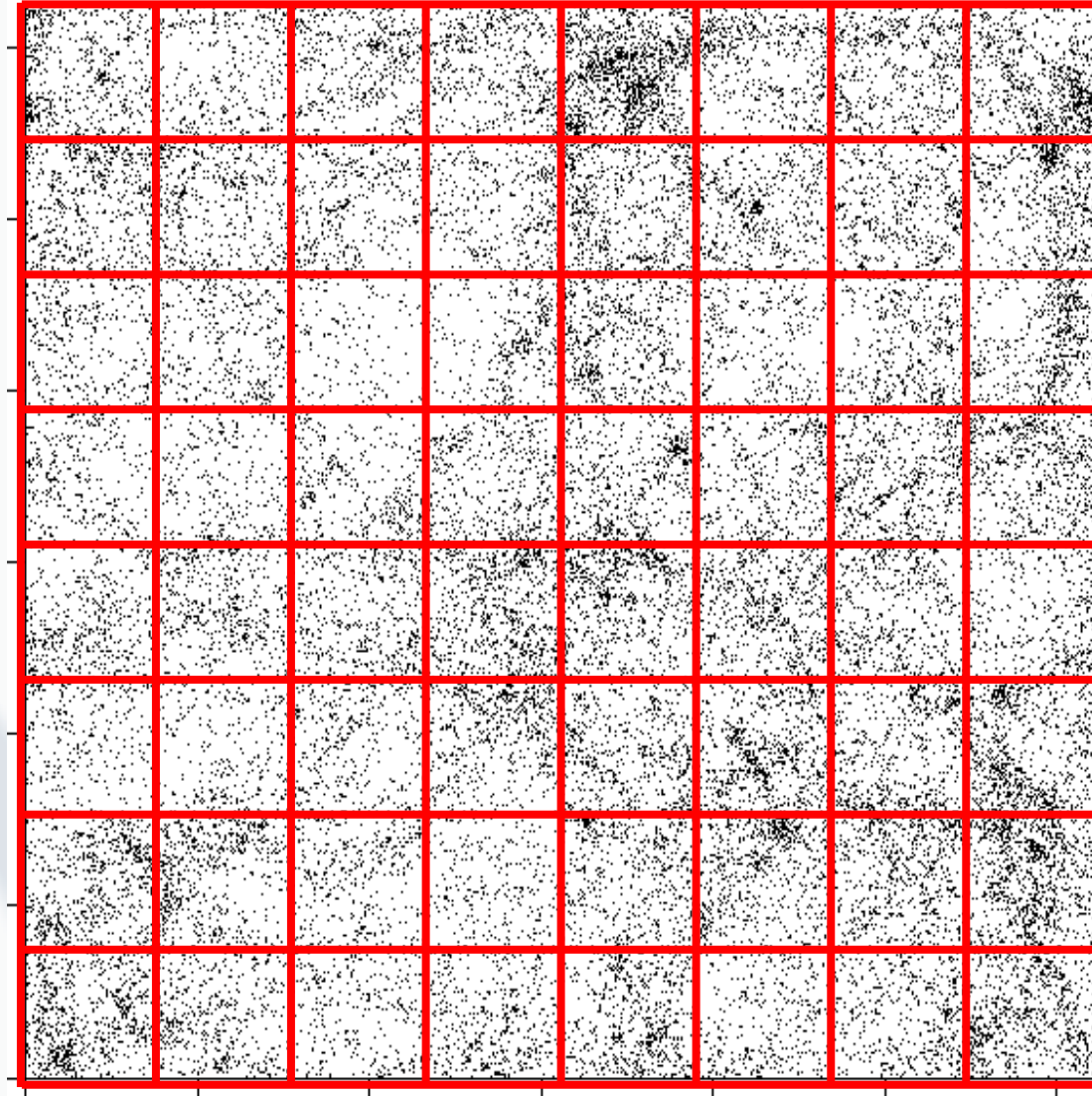
https://link.springer.com/chapter/10.1007/978-3-642-22351-8_34



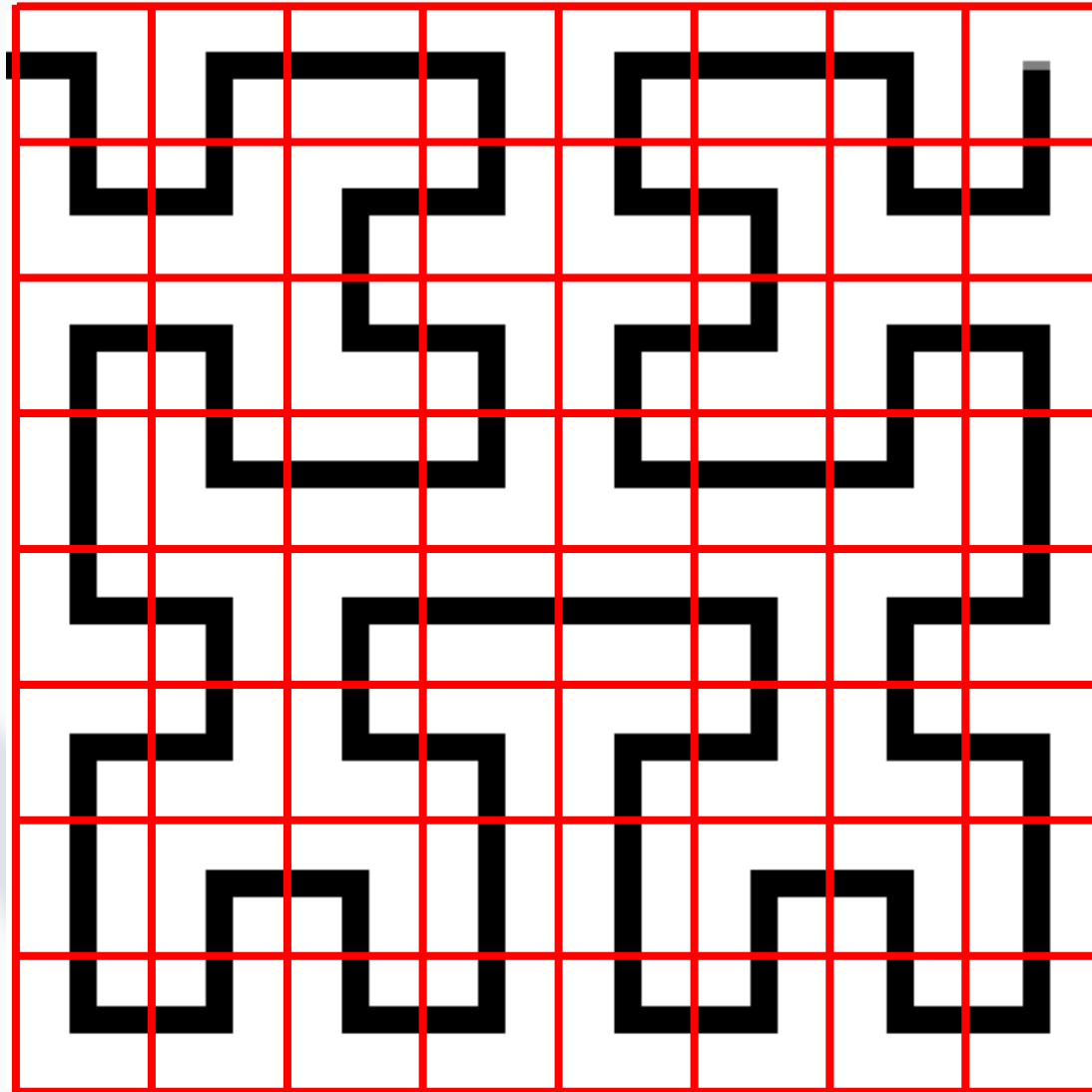
Find particles in some sub-volume



Divide simulation volume in regular grid



Index cells using space filling curve



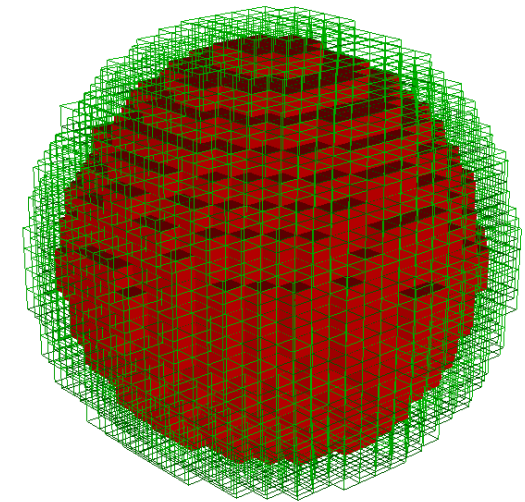
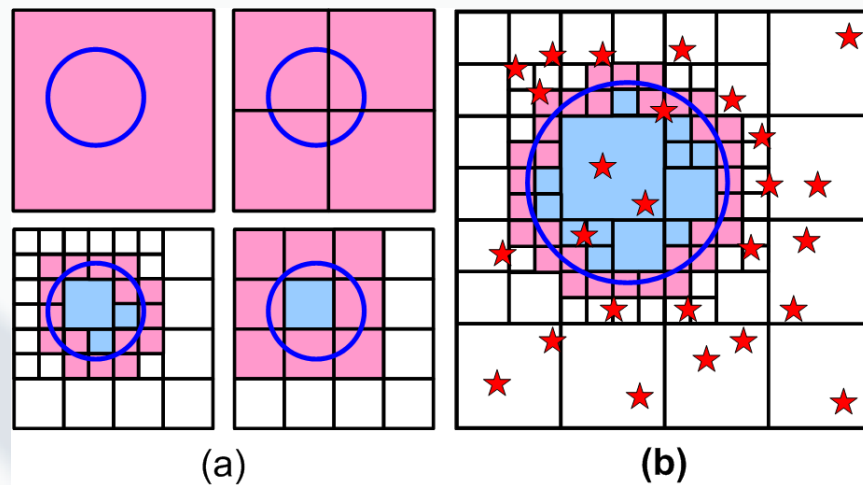
Sort data by index, in DB or files

1	4	5	6	59	60	61	64
2	3	8	7	58	57	62	63
15	14	9	10	55	56	51	50
16	13	12	11	54	53	52	49
17	18	31	32	33	34	47	48
20	19	30	29	36	35	46	45
21	24	25	28	37	40	41	44
22	23	26	27	38	39	42	43

Lemson, Budavari & Szalay (2011)

https://link.springer.com/chapter/10.1007/978-3-642-22351-8_34

- ▶ Divide simulation volume in regular grid
- ▶ Index using space filling curve (Peano-Hilbert, Morton)
- ▶ Calculate overlap space filling curve with query volume
 - Iterate from root volume down, stopping at fully contained boxes



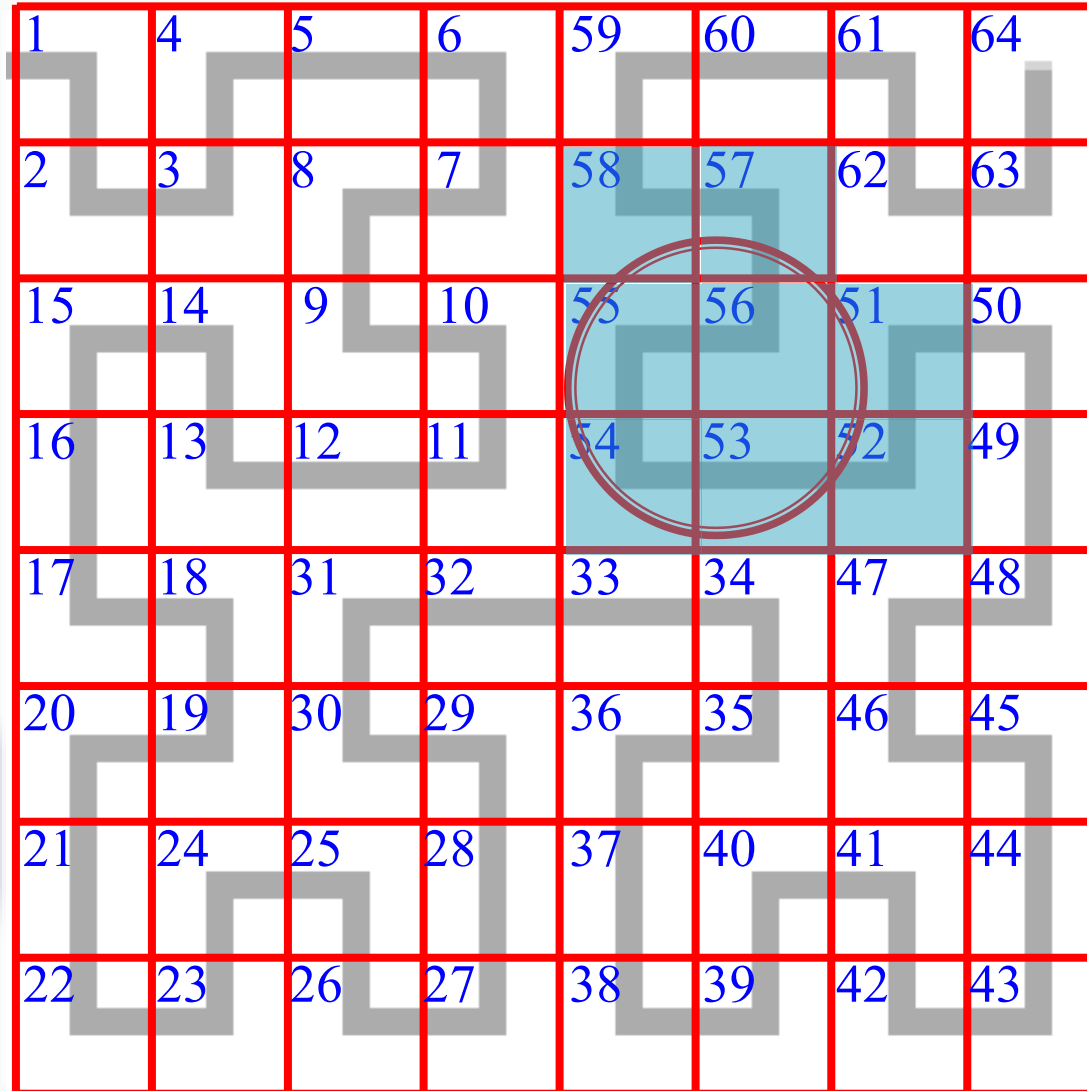
- ▶ Find particles in overlap ranges,
 - Only for those filter further on exact 3D volume
- ▶ Execute as table-valued function from database

Space-filling ordering co-locates near-by data better

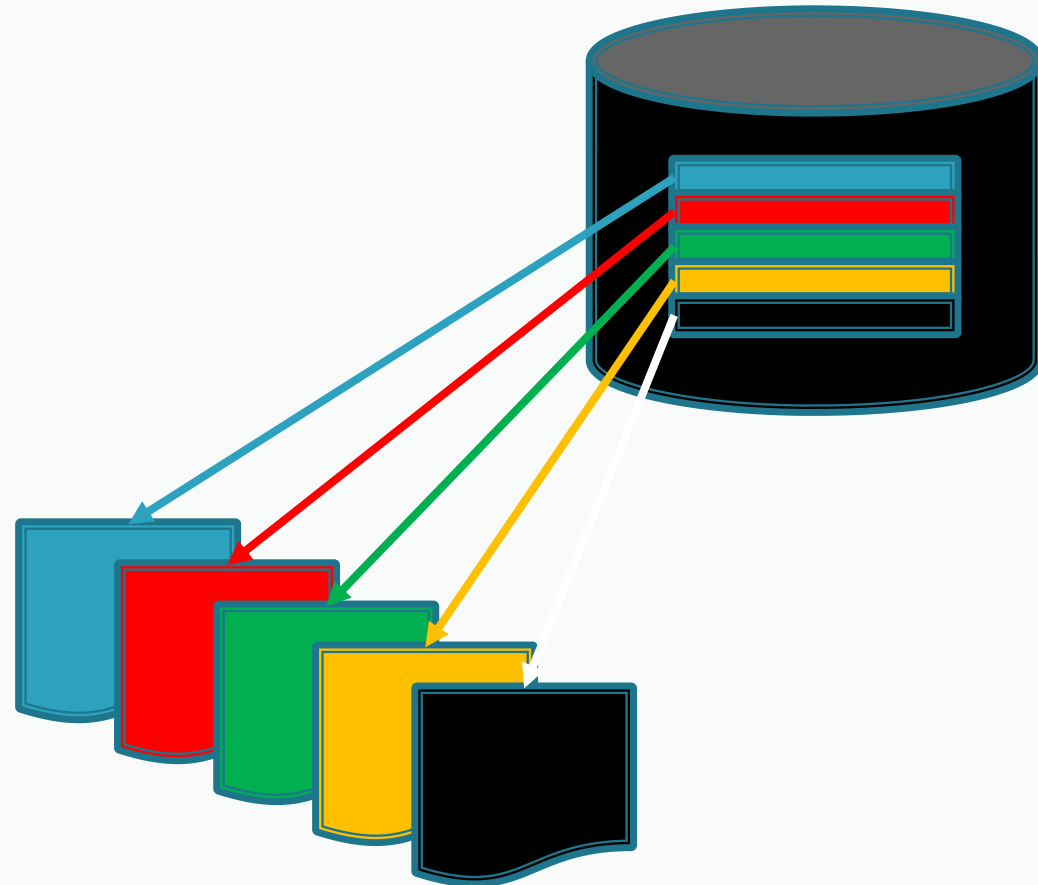
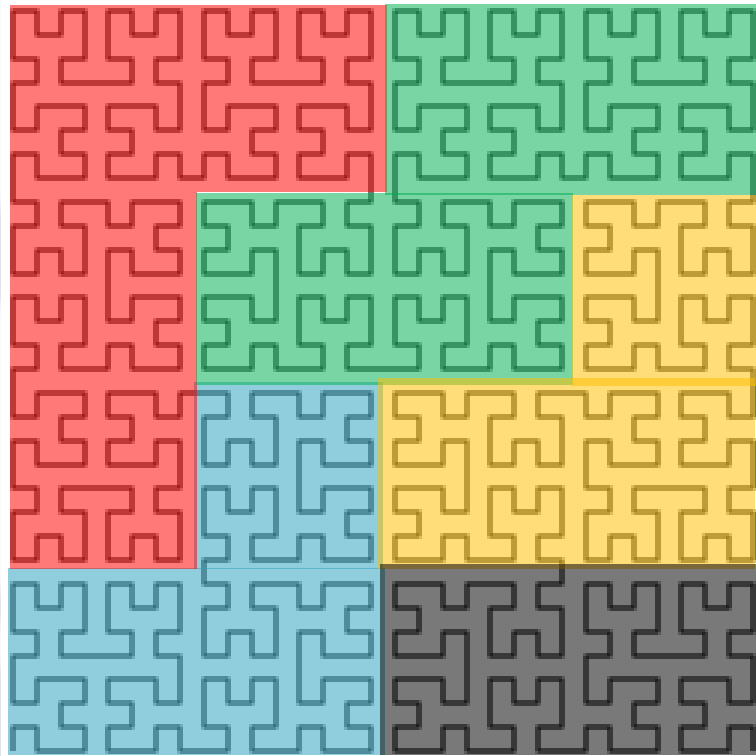
51 52 53 54 55 56 57

vs

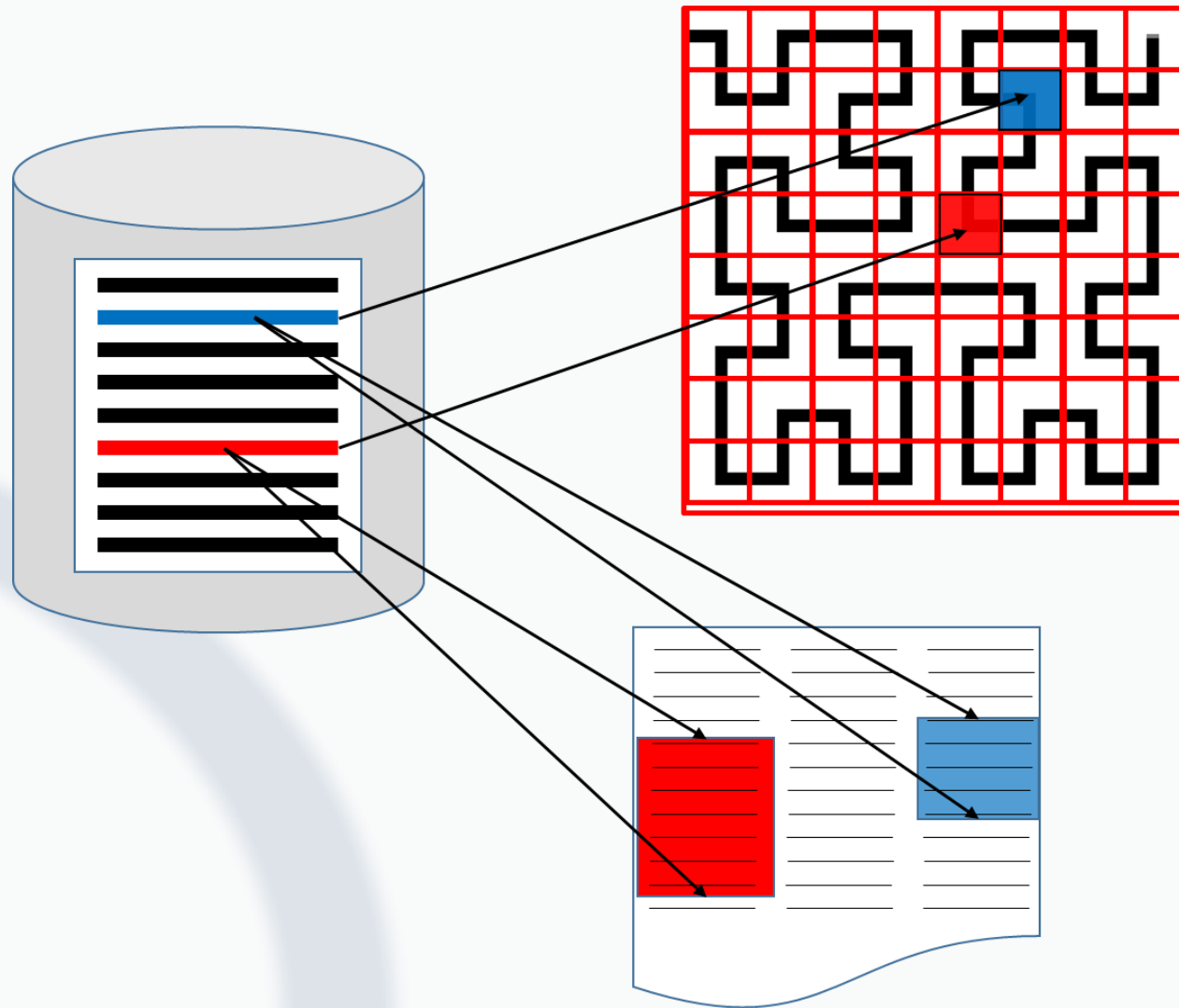
15 16
23 24 25
31 32 33



Individual files store intervals on curve; indexed in DB



DB also stores locations of “buckets” inside files.



```

jupyter CosmoUC_5_PlotHalo Last Checkpoint: 11 minutes ago (autosaved)
File Edit View Insert Cell Kernel Help Notebook saved Python 3
+ -> ↺ ↻ ↶ ↷ ⌂ Code CellToolBar

Dark-matter Halos in a Cosmological Simulation

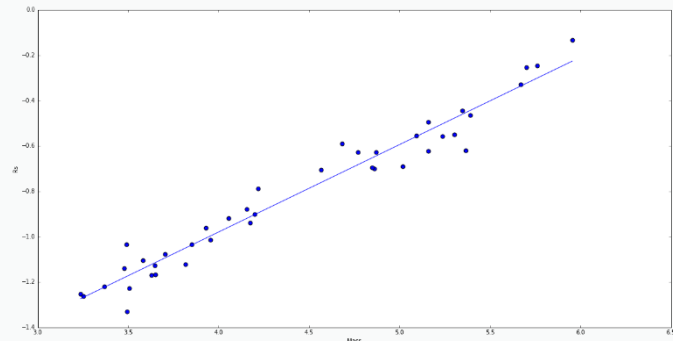
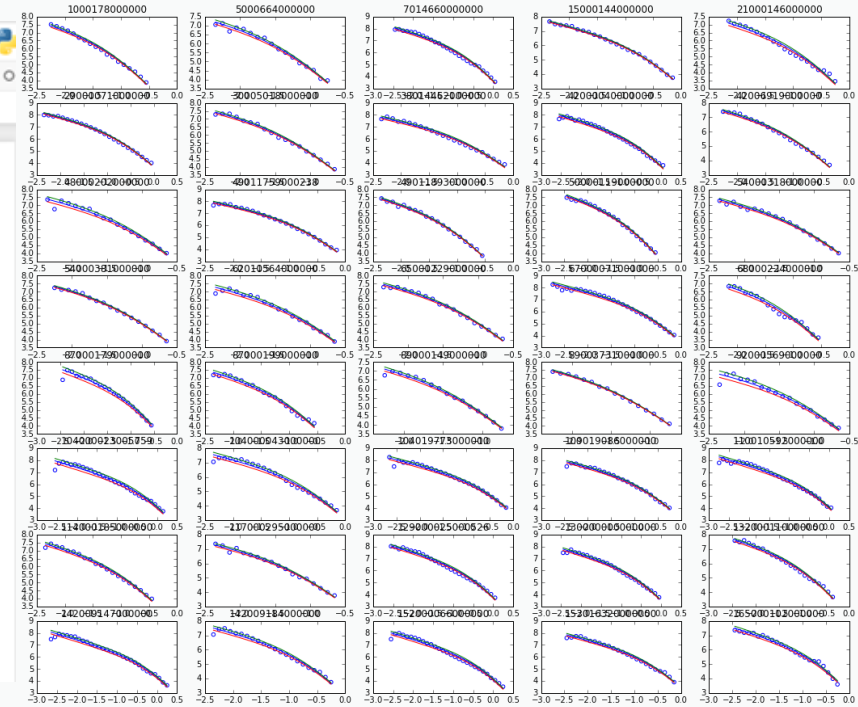
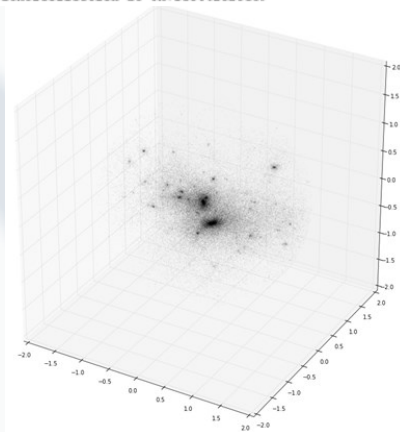
In [3]: import SciServer.LoginPortal as Login
        token = Login.getToken()
        import SciServer.CasJobs
        import pandas
        import tables
        import numpy as np
        import matplotlib.pyplot as plt
        from mpl_toolkits.mplot3d import Axes3D

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 mpahalotrees.nr hh
        cross apply dbo.MillenniumParticles(hh.snapnum,
        dbo.Sphere:New(hh.x,hh.y,hh.z,3*hh.halfmassradius).ToString()) p
        where hh.haloid=84000007000000 order by newid()
        """
        responseStream = SciServer.CasJobs.executeQuery(queryString, token=token,context="SimulationDB")
        df = pandas.read_csv(responseStream, index_col=None)

        CPU times: user 351 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 0x7fe8428b9e8>
    
```




```
ax.set_title(title)
divider = make_axes_locatable(ax)
# Append axes to the right of ax3, with 20% width of ax3
cax = divider.append_axes("right", size="5%", pad=0.05)
# Create colorbar in the appended axes
# Tick locations can be set with the kwarg `ticks`
# and the format of the ticklabels with kwarg `format`
cbar = plt.colorbar(im, cax=cax, ticks=MultipleLocator(0.2), format="%0.2f")
```

CPU times: user 17.4 s, sys: 3.87 s, total: 21.3 s
Wall time: 27 s

