



EARTH CUBE



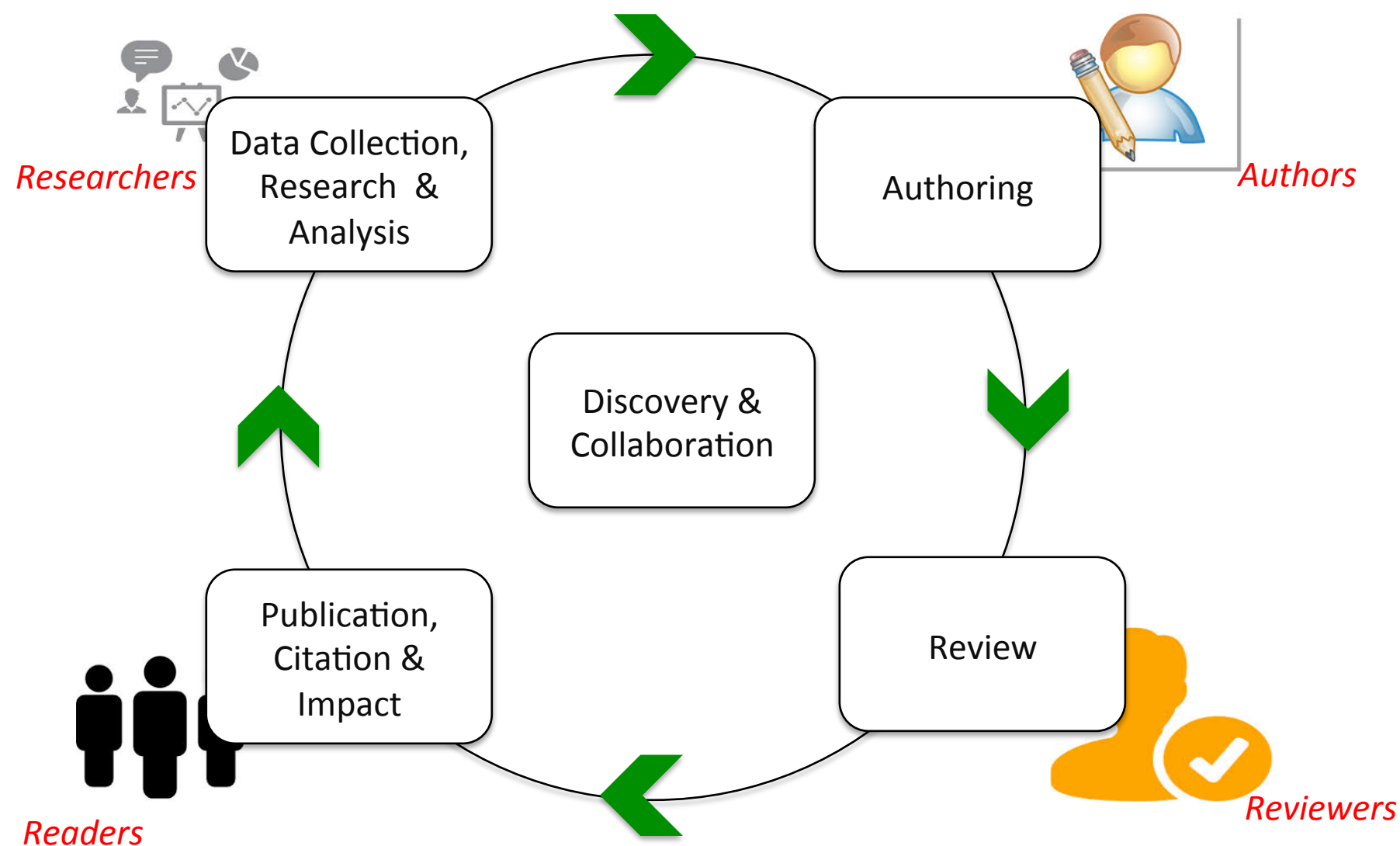
Science Dataspaces for Data Management and Reproducibility

Tanu Malik, Ian Foster, Kyle Chard

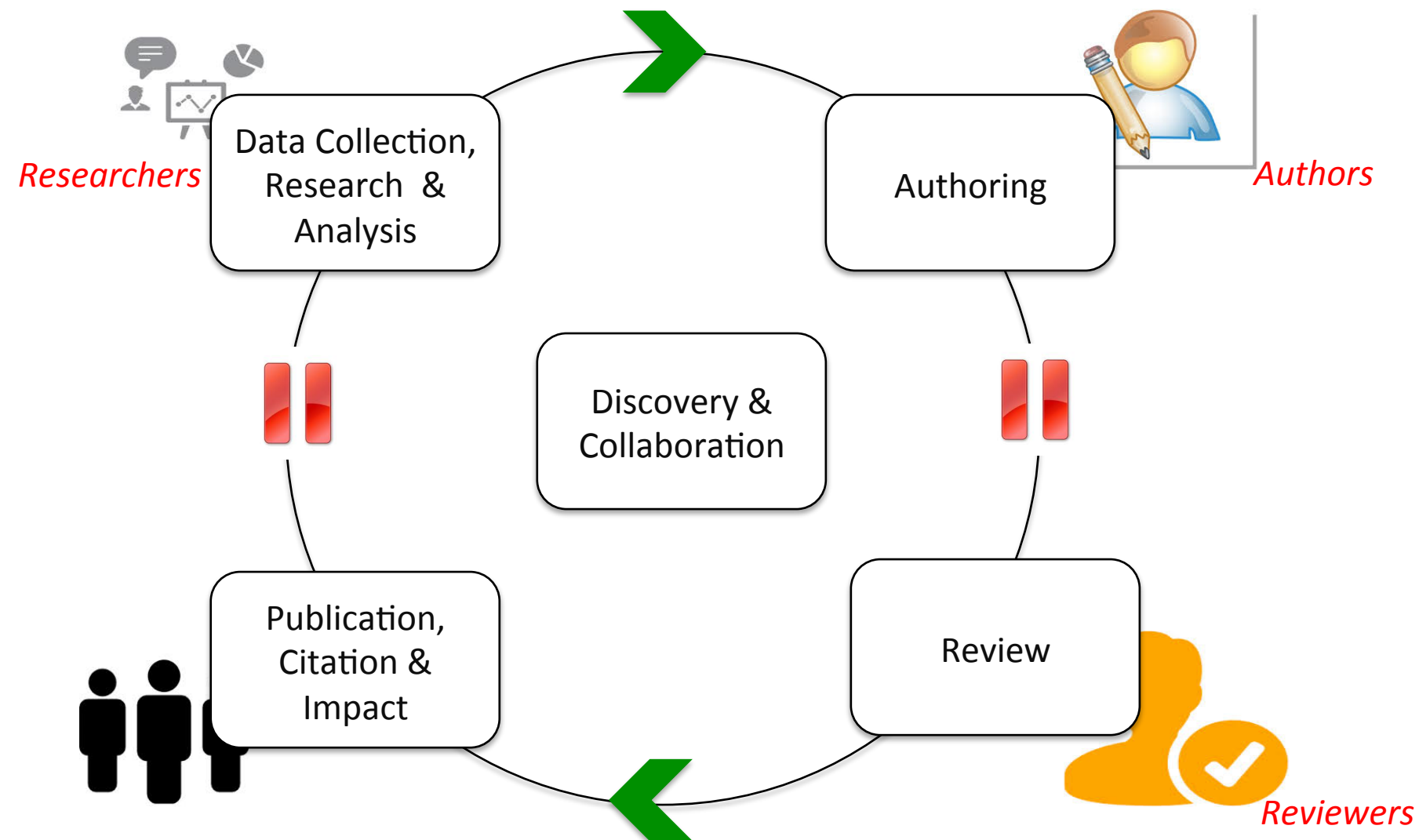
Jonathan Goodall, Scott Peckham, Joseph Baker, Mike Gurnis



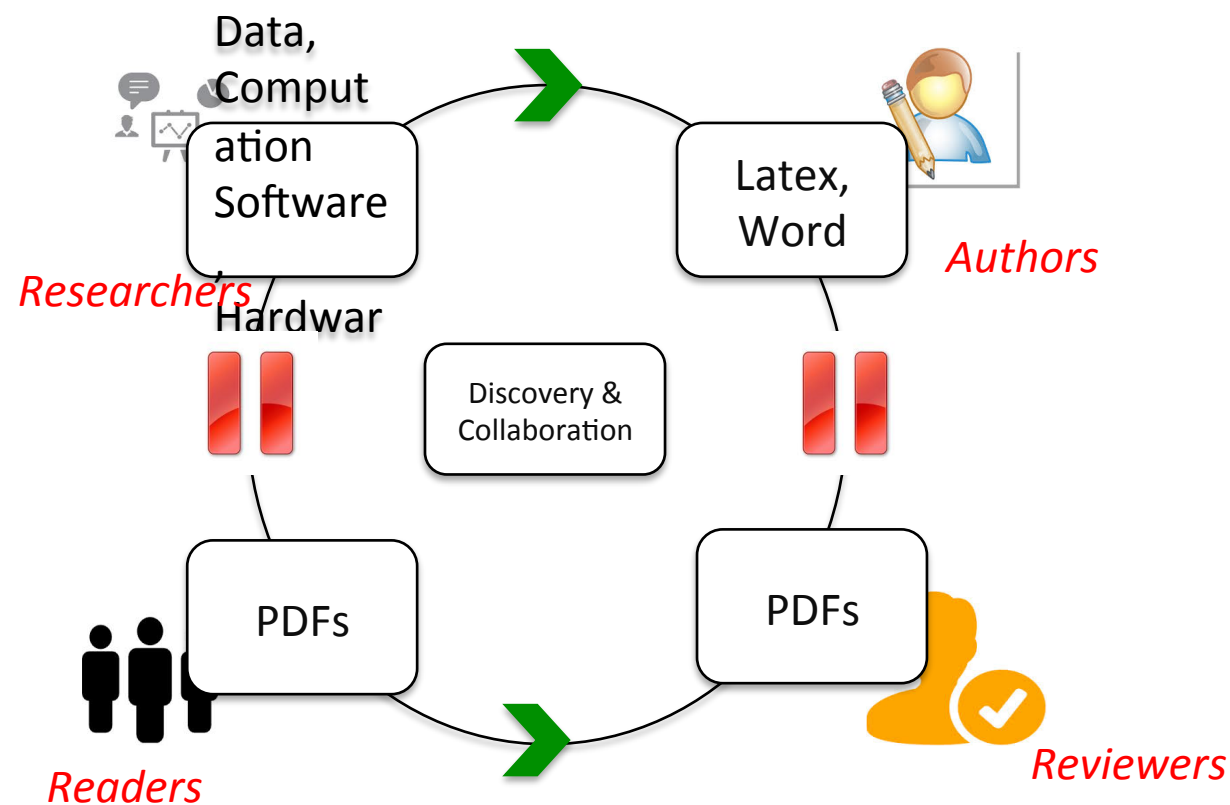
The scientific method is self-correcting



But...Scientific Publishing is Not



Computational Scientific Publishing is Broken



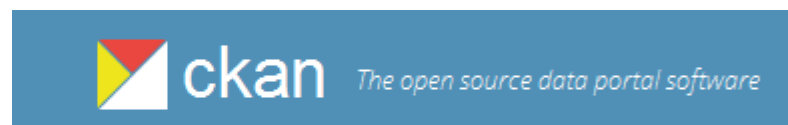
Computational science inputs are not linked with outputs.

- *Inputs:* Large quantities of data, complex data manipulation and/or numerical simulation use of large and often distributed software stacks, etc. (software, data, execution, environment)
- *Outputs:* Research papers (text-based, non-interactive)

Encourage Open-X



- X = access/code/data/design/standard
- Use the Internet and be more social



Two examples of Open-X



Share
with
collaborators

GitHub



Share publicly
and
get social credit



- Why aren't Github + DropBox sufficient for data management and reproducibility?

Shared Github Repo

The screenshot shows the GitHub interface for the repository 'uva-hydroinformatics-lab / VIC_Pre-Processing_Rules'. The repository has 4 watchers, 0 stars, and 0 forks. It contains 62 commits, 1 branch, 0 releases, and 1 contributor. The main branch is 'master'. A list of files is displayed, including .gitattributes, .gitignore, Main_Shell_Script.scr, combine_wind, combine_wind.c, convertPrp.cpp, convertTmax.cpp, convertTmin.cpp, convert_tif_ascii.py, create_LDAS_soil_nearest.c, create_LDAS_veg_param.c, and oet_prism.c. The latest commit is 581afe7 on Oct 13.

Repository: uva-hydroinformatics-lab / VIC_Pre-Processing_Rules

Watch 4 Star 0 Fork 0

Code Issues 0 Pull requests 0 Wiki Pulse Graphs

Preprocessing Rules

62 commits 1 branch 0 releases 1 contributor

Branch: master New pull request

New file Find file HTTPS https://github.com/uva-hydr Download ZIP

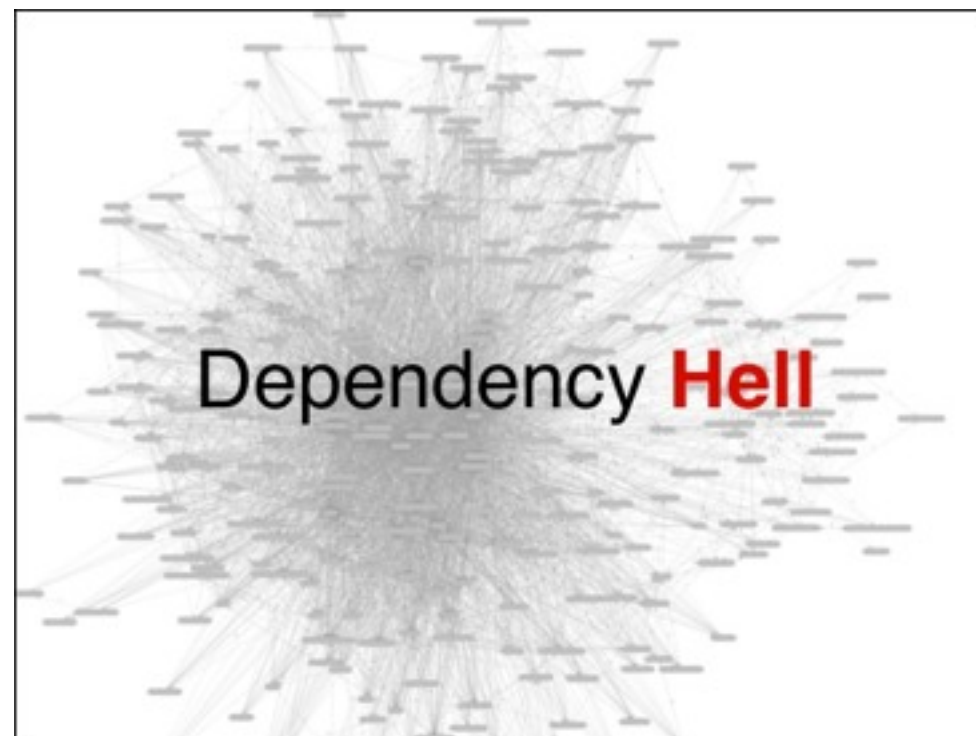
Bakinam Shell script Latest commit 581afe7 on Oct 13

.gitattributes	Added .gitattributes & .gitignore files	2 months ago
.gitignore	Added .gitattributes & .gitignore files	2 months ago
Main_Shell_Script.scr	Shell script	2 months ago
combine_wind	combine_wind.c	2 months ago
combine_wind.c	Combine_wind.c	2 months ago
convertPrp.cpp	convertPrp.cpp	2 months ago
convertTmax.cpp	Source code	2 months ago
convertTmin.cpp	source code	2 months ago
convert_tif_ascii.py	convert_tif_ascii.py	2 months ago
create_LDAS_soil_nearest.c	source code	2 months ago
create_LDAS_veg_param.c	shell script	2 months ago
oet_prism.c	Source code	2 months ago



Missing Dependencies

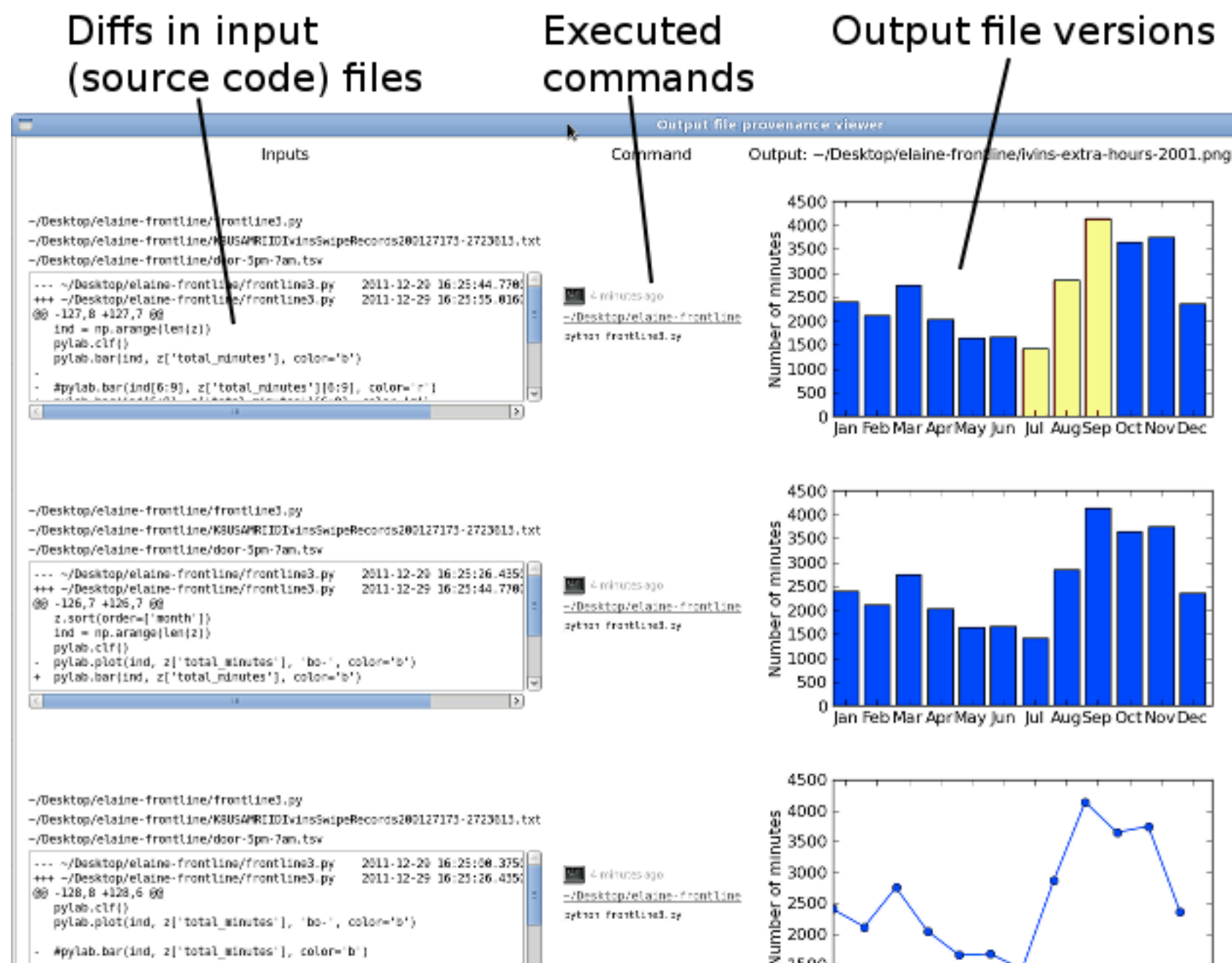
- Has the author shared everything, including code, data, and environment?
- Will this code, if downloaded, run in my computational environment?





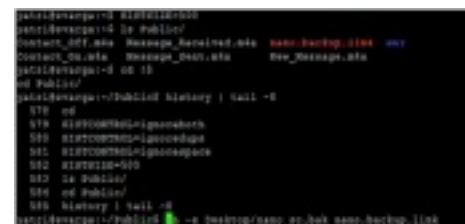
Missing Lineage

- Which version of the data produced this result?





- 





SciDataspace

Personalized, Shareable Dataspace
for

Data Management and Reproducibility



GeoDataspace

Username :

Password :

[Sign in](#)

The GeoDataspace framework is assists scientists and communities to create and maintain collections of geounits that pertain to a specific research project. To create and use GeoDataspace, please sign up and download the client.



EARTH CUBE



GeoDataspace client



- Git and DropBox-like Python client for Linux and Mac OS X
 - annotate
 - provide semantic annotations
- package
 - code, data, environment into Docker containers
- track
 - tracks provenance of the scientific program

GeoDataspace Tools

Tool	Repeatability Feature	Users
SPADE ¹	Provenance (V)	*
CDE ²	Packaging (P)	*
PTU ³	P+V	
SciDataspace⁴	P+V+Usability	Geoscience NSF EarthCube
PTU-NFS ⁵	P'+V' (in Network File System)	DOE High-Energy Physics
LDV: Light-weight Database Virtualization ⁶	P'+V' (in databases)	Urban Science

1. <https://github.com/ashish-gehani/SPADE>

2. <http://www.pgbovine.net/cde.html>

3. <https://gitlab.com/quanpt/provenance-to-use>

4. <https://bitbucket.org/tanum/scids-client>

Preserved in GeoDataspace

The screenshot displays the GeoDataSpace web interface. At the top, the header includes the GeoDataSpace logo and the text "Manage geounits | tanum". Below the header, there are navigation links: "manage geounits", "transfer geounits", and "dashboard".

On the left side, there is a "GeoDataSpaces" section with a dropdown menu showing "Test". Below this is a "Filter by Annotation" section with a list of filters: "date", "not working", "owner", "special", "testagain", and "working".

The main content area shows a dataset named "Dataset1" with a date of "2014-09-24". The dataset is owned by "u:tanum" and has a label. Below the dataset name, there are tabs for "Overview", "Tags", "Sharing", "Files", "Packages", and "Commands". The "Files" tab is selected, showing a list of files and directories:

- bin
- etc
- home
 - ubuntu
 - .cache
 - .config
 - default
 - monthlySoilMoistureEcohydro.csv
 - run_psp_vic_soilmoisture.scr
 - run_psp_vic_soilmoisture.scr.cde
 - spatiotempSoilMoistureEcohydro.csv
 - spatiotempdatabase.py
 - spatiotempdatabase.pyc
 - uploadToS3.py
 - vicSoilMoistureEcohydro.pdf
 - vic_calc_mnth_mc.py
 - vic_monthly_soilmoisture.py
 - vic_soil_moisture.py
- lib
- lib64
- usr

At the bottom of the dataset view, there are two more entries: "test 2" and "test again", both dated "14-09-30".

An arrow points to the "Edit" link next to the "Dataset1" entry, with the word "Edit" written next to it.

The URL in the browser address bar is "local.gl.com/datasets/index.html#".

Manipulate Geounits



Save

Graph

Packages

bin

etc

home

ubuntu

.cache

.config

default

monthlySoilMoistureEcohydro.csv

run_psp_vic_soilmoisture.scr

run_psp_vic_soilmoisture.scr.cde

spatiotempSoilMoistureEcohydro.csv

spatiotempdatabase.py

spatiotempdatabase.pyc

uploadToS3.py

vicSoilMoistureEcohydro.pdf

vic_calc_mnth_mc.py

vic_monthly_soilmoisture.py

vic_soil_moisture.py

lib

lib64

usr

/home/karthik/Desktop/Extra_Files/FLP/Tanu/Packages_tmp/667c7d58b2fb1fae46bd92ab75333f34743718ec/cde-package/cde-root/home/ubuntu/monthlySoilMoistureEcohydro.csv

98

mc,200601,0.23,0.24,0.28

99

mc,200602,0.22,0.23,0.29

100

mc,200603,0.20,0.21,0.27

101

mc,200604,0.19,0.19,0.23

102

mc,200605,0.21,0.20,0.18

103

mc,200606,0.21,0.21,0.15

104

mc,200607,0.20,0.20,0.13

105

mc,200608,0.21,0.20,0.11

106

mc,200609,0.21,0.23,0.16

107

mc,200610,0.20,0.21,0.14

108

mc,200611,0.23,0.24,0.19

109

mc,200612,0.21,0.22,0.24

110

mc,200701,0.23,0.24,0.29

111

mc,200702,0.22,0.23,0.31

112

mc,200703,0.21,0.22,0.32

113

mc,200704,0.21,0.21,0.27

114

mc,200705,0.19,0.21,0.22

115

mc,200706,0.20,0.17,0.15

116

mc,200707,0.20,0.17,0.12

117

mc,200708,0.16,0.13,0.10

118

mc,200709,0.18,0.14,0.10

119

mc,200710,0.18,0.14,0.09

120

mc,200711,0.18,0.17,0.10

121

mc,200712,0.21,0.20,0.10

122

Shell Access:

karthik@murray:~\$

karthik@murray:~\$

Track the workflow

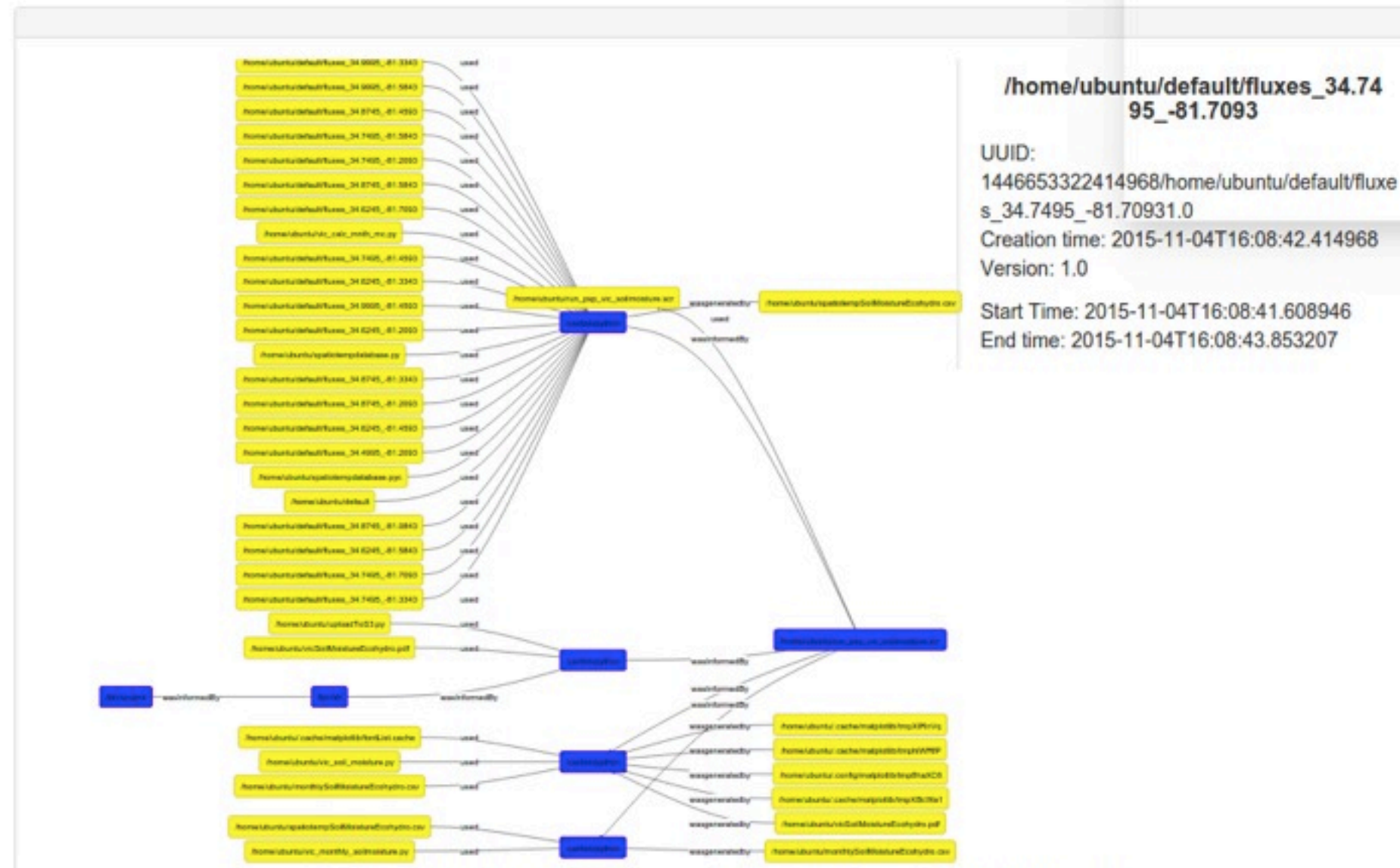


GeoDataSpace

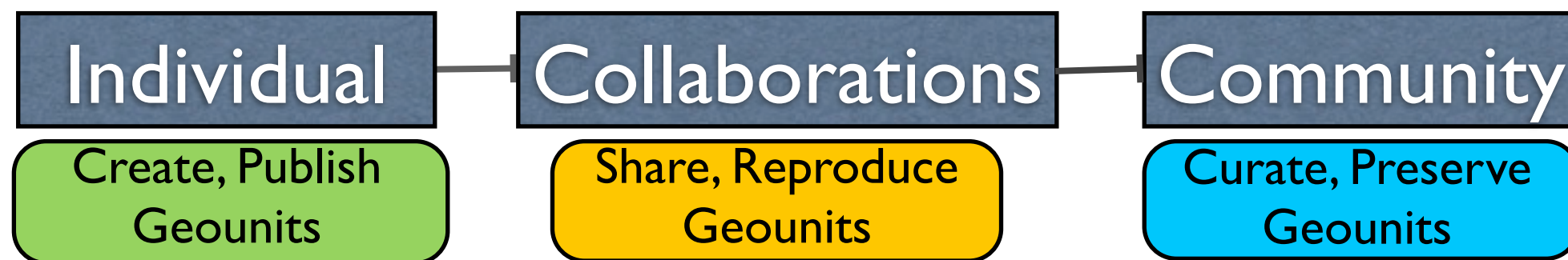
Manage geounits | Groups | News & Events | About | Support | Log In | Sign Up

Graph

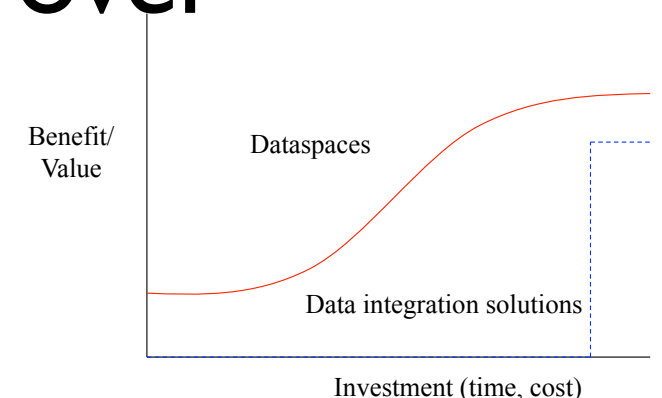
Packages



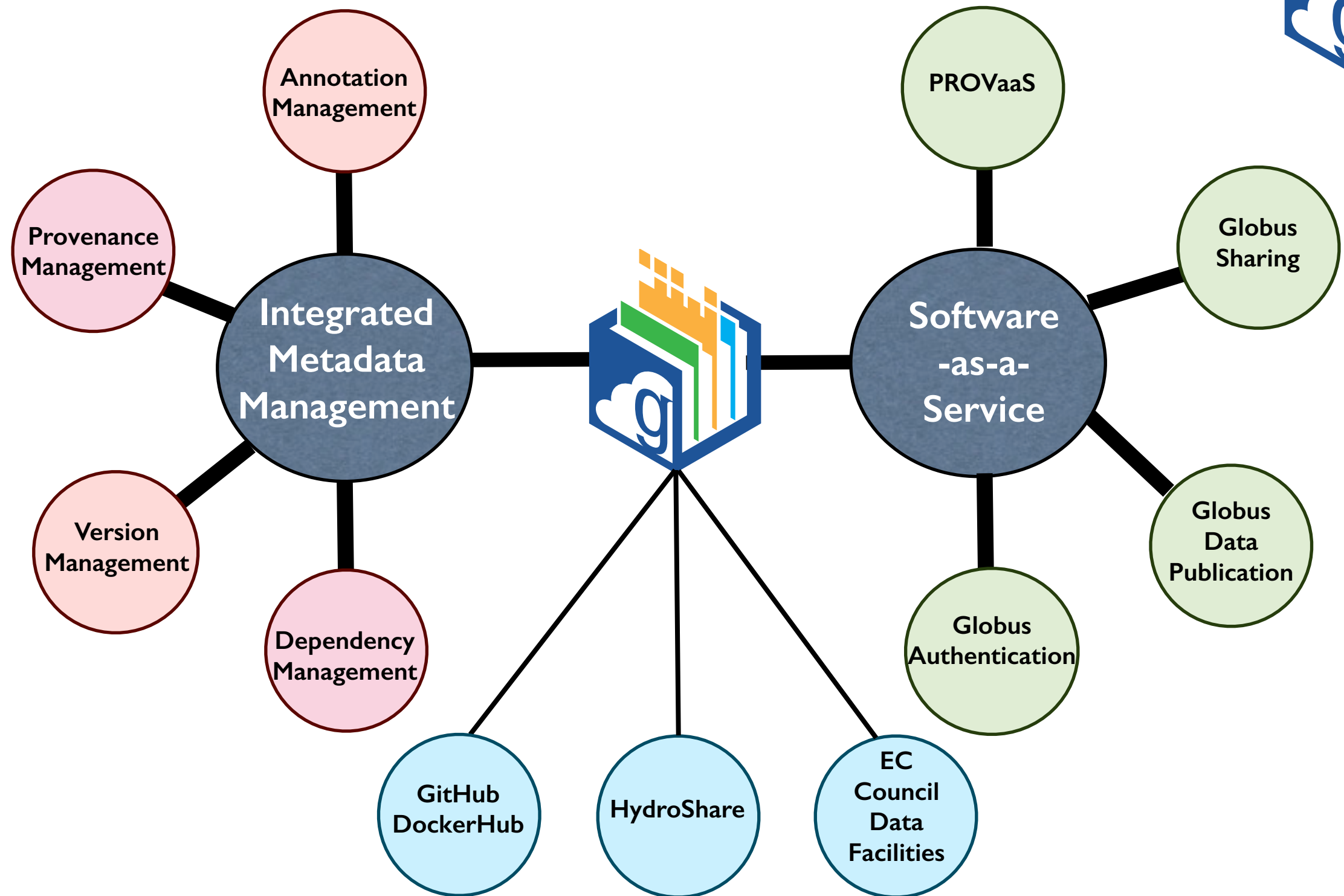
Preserved in a Dataspace



- Cloud-hosted Dataspace that can be shared with the collaboration with the help of services, and becomes standardized over time.



GeoDataspace Architecture





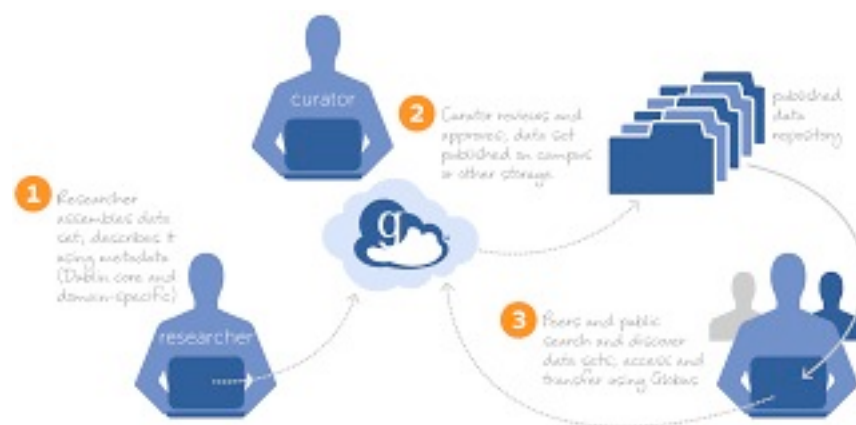
Globus Services



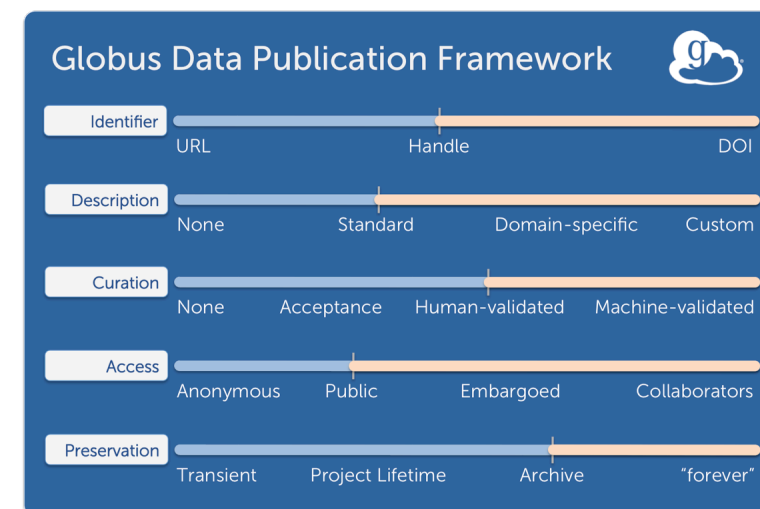
Authentication



Transfer

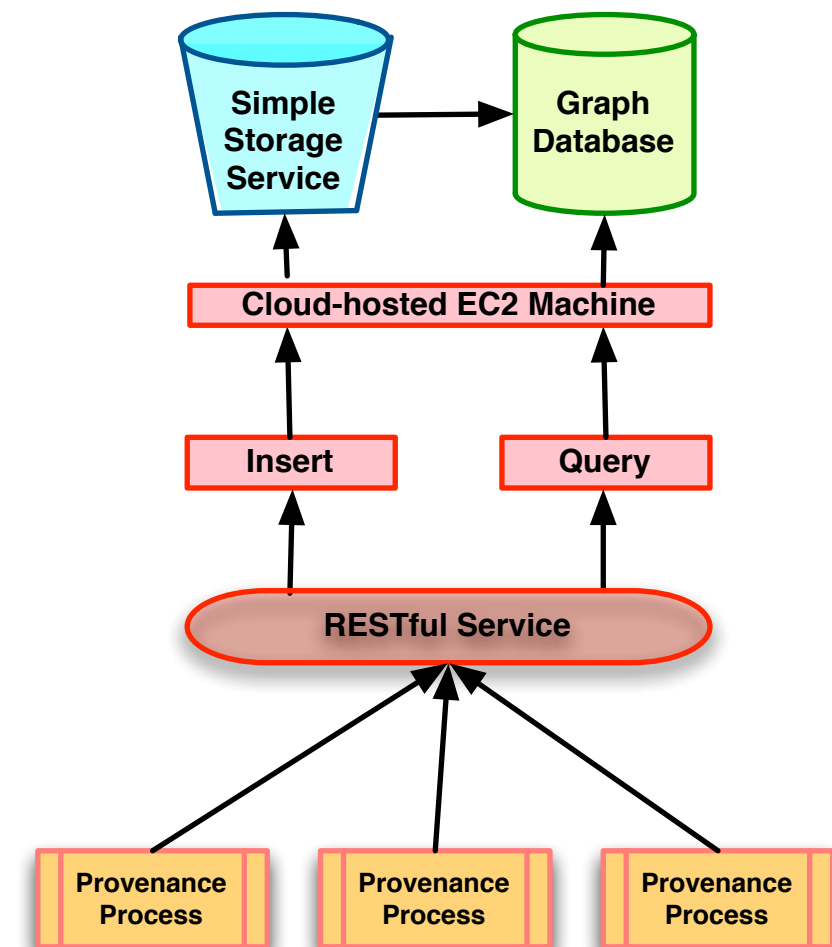
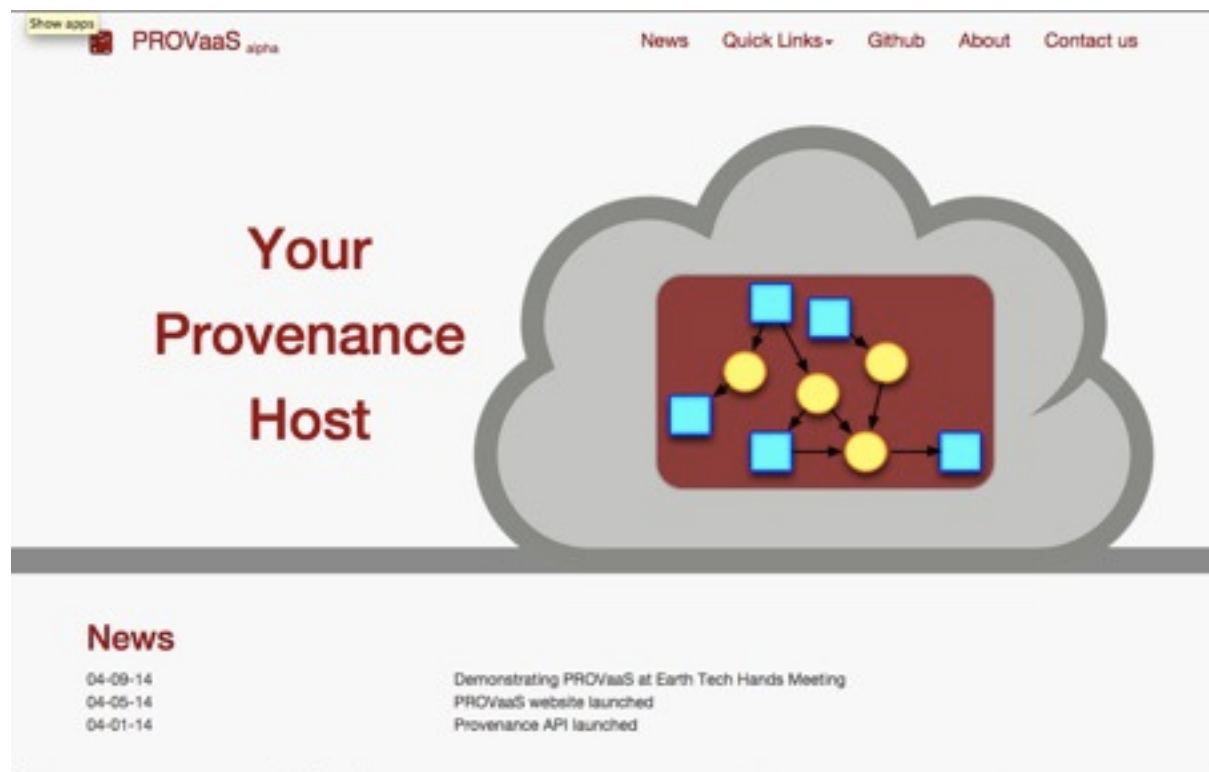


Sharing



Publication

PROVaaS

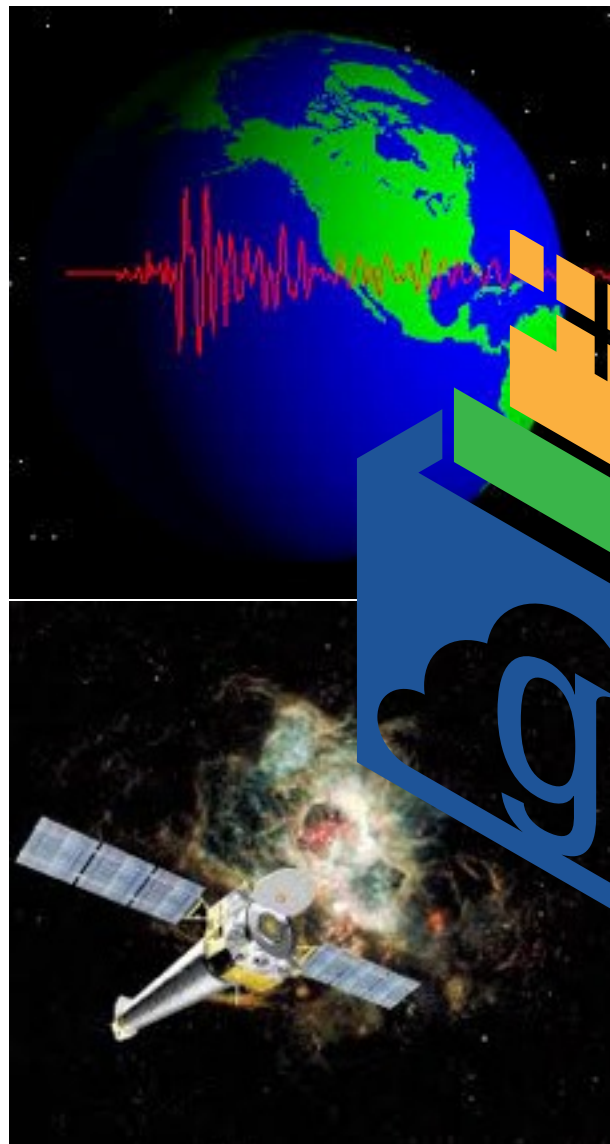


Bringing Order to Chaos in Science Use Cases



Solid Earth

Hydrology



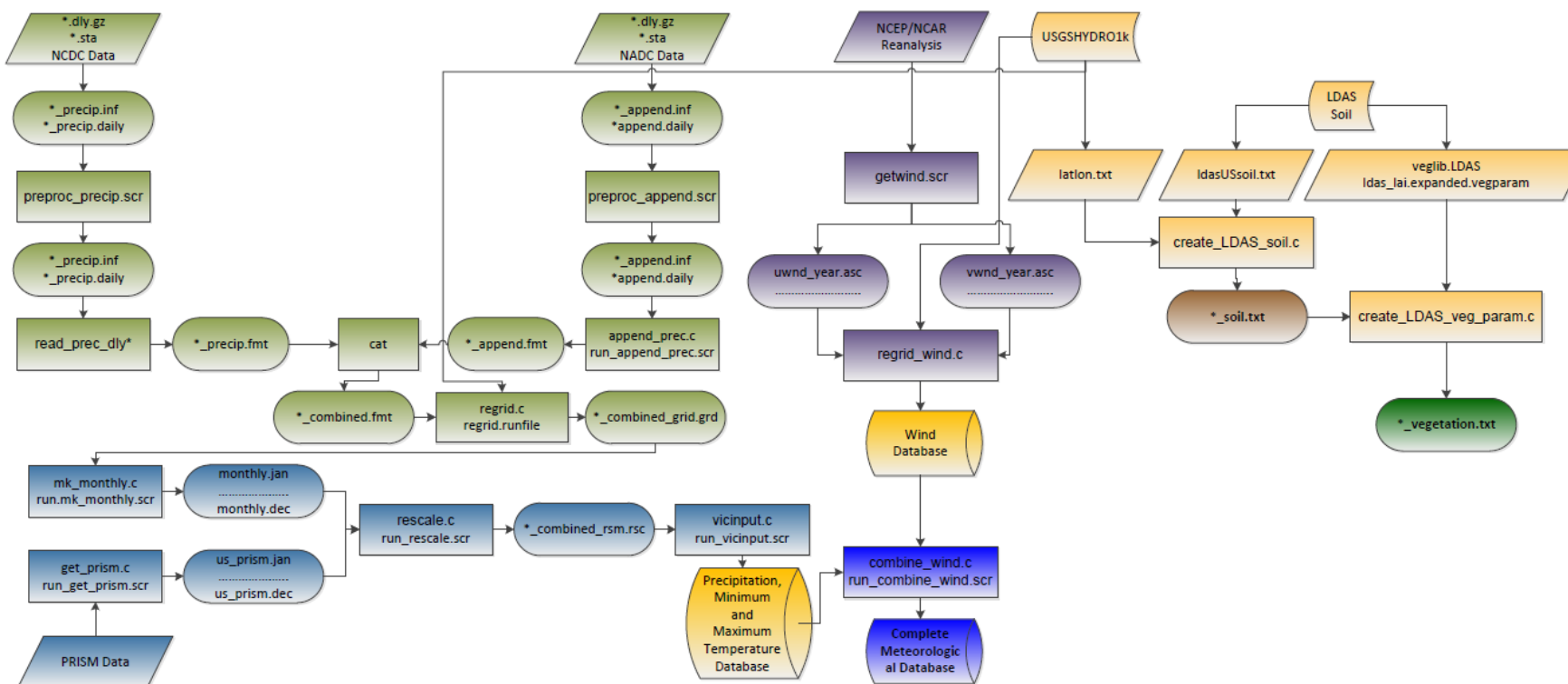
Space Science



CSDMS

Hydrology

(Curation, Encapsulation, Reuse)



- Data processing steps for the VIC model

The reality



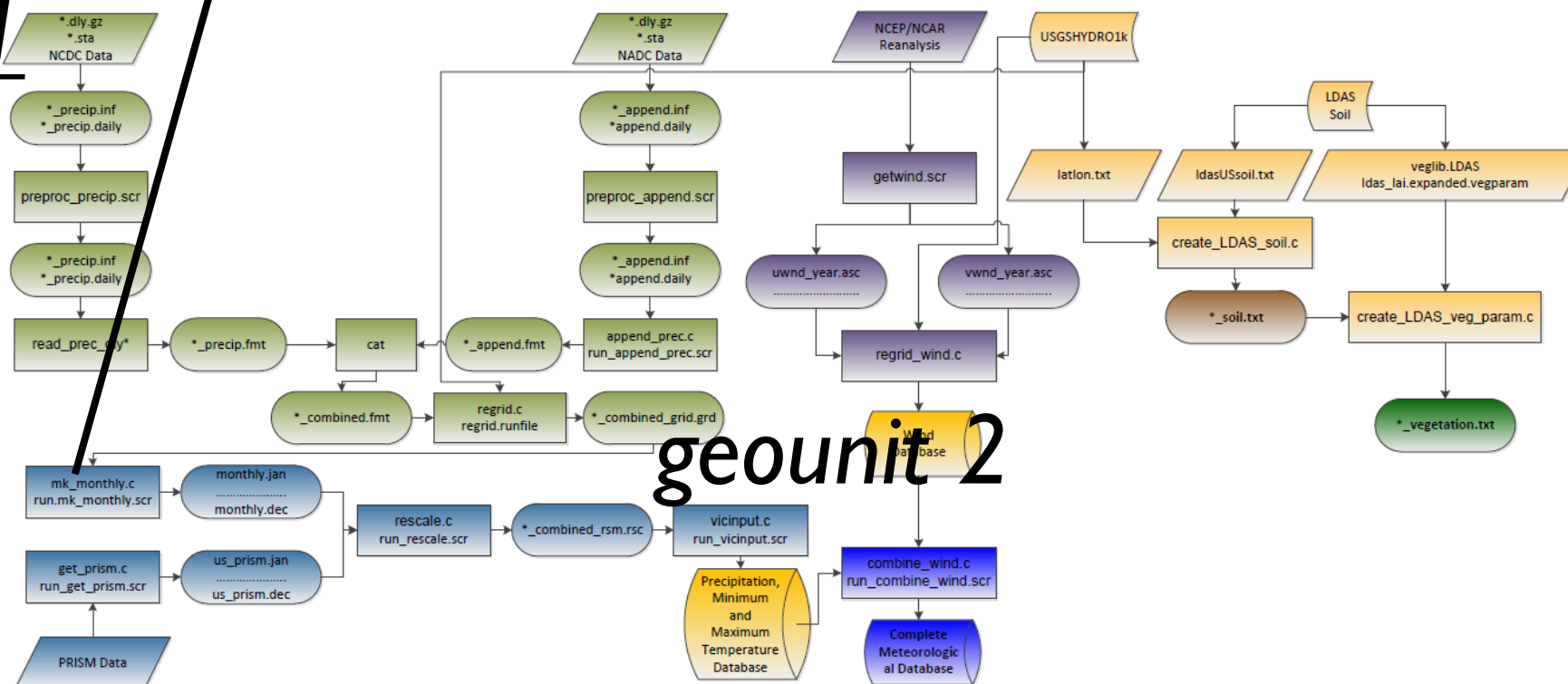
```
irods@ec2-54-86-215-185:~/iRODS/server/bin$ cd cmd/
irods@ec2-54-86-215-185:~/iRODS/server/bin/cmd$ ls
adjprcp_rsm.rsc      forcingData          old                  rescale_ir.scr      stationinfo.txt
adjtmin_grid.grd    gauss_t62_lat.list  osgeo.bak           run_combine_wind_ir.scr  temp_wind
amqprecv.py          gauss_t62_lon.list  output_1.txt        run_combine_wind.scr  TerraPopCountieswzVicOutput.py
amqpseend.py         gdal.py             PamAuthCheck        run_convert_prcp.scr  test
Basin                gdal.pyc            PopulationVsSoilMoisture.scr  run_convert_tif_ascii.scr  test3.py
basin_prcp_adj.fmt   get_prism            popvssm.run         run_convert_tmax.scr  test_execstream.py
basin_prcp.fmt       get_prism.c          prcp                run_convert_tmin.scr  test.py
boundaries_US_SLAD_2010.dbf  get_prism_ir.scr    prcp.daily          run_get_prism.scr     test.scr
boundaries_US_SLAD_2010.prj  getwind.scr         prcp.inf            run.mk_monthly.scr    tiff2ascii.py
boundaries_US_SLAD_2010.sbn  hello               prcp_tobAdj.scr     run_psp_vic_evapotranspiration.scr  Tmax
boundaries_US_SLAD_2010.sbx  hello.scr           prec_tob_adj        run_psp_vic_sm.scr    tmax_tob_adj
boundaries_US_SLAD_2010.shp  inputPrcp.scr       prec_tob_adj.f      run_psp_vic_soilmoisture1.scr  tmax_tob_adj.f
boundaries_US_SLAD_2010.shp.xml  inputTmax.scr      prec_tob_adj.input  run_psp_vic_soilmoisture2.scr  tmax_tobAdj.scr
boundaries_US_SLAD_2010.shx  inputTmin.scr      preproc_precip.scr  run_psp_vic_soilmoisture3.scr  Tmin
build                  irodsAgent          prism               run_psp_vicSoilMoistureComparison.scr  tmin_tob_adj
catchment.dbf          irodsReServer       prism-rawdata       run_psp_vic_soilmoisture.scr  tmin_tob_adj.f
catchment.prj          irodsServer         python_script1.py   run_regrid_wind_ir.scr  tmin_tobAdj.scr
catchment.shx          irodsServerMonPerf  raw_wind            run_regrid_wind.scr    tob
checkData              irodsXmsgServer     read_prec_dly       run_rescale.scr        univMSSInterface.sh
climate                latlong99.txt       read_prec_dly.f     run_vicinput_ir.scr    uploadToS3.py
cmd                    latlon.txt          readRodsLog.py      run_vicinput.scr       vacuumdb.pl
combine_wind           LDAS                read_temp_dly.f     script_1.scr            vegetation
combine_wind.c         ldas_lai.expanded.vegparams  read_tempn_dly      script2.scr             vic_calc_mnth_mc1.py
convertPrcp            ldas_latlon2.scr    read_tempx_dly      script3.scr             vic_calc_mnth_mc2.py
convertPrcp.cpp         ldas_latlon3.scr    regrdPrcp           script4.scr             vic_calc_mnth_mc.py
convert_tif_ascii.py    ldas_latlon4.scr    regrd_prcp.scr      script5.scr             vic_calc_mnth_peb.py
convertTmax            ldas_latlon.scr     regrdTmax           script6.scr             vic_evapotranspiration.py
convertTmax.cpp         ldas_soil.scr       regrd_tmax.scr      script7.scr             vicinput
convertTmin            ldas_veg.scr        regrd_tmin.scr      script.save             vicinput.c
convertTmin.cpp        LeastSoilMoistureduration.py  regrd_tmin.scr      script.scr              vic_monthlyPEBestimates.py
coop_tob.his           list.pl             regrdPrcp           script_test.scr         vic_monthly_soilmoisture2.py
create_LDAS_soil_nearest  mask2latlon        regrdPrcp.runfile   script_vic_pre-processing.scr  vic_monthly_soilmoisture.py
create_LDAS_soil_nearest.c  mask2latlon.c      regrdTmax           sm_comparison.py        VIC_Pre_processing_script.scr
create_LDAS_veg_param     metadata.txt        regrdTmin           sm_seasonal.py          vic_soil_moisture3.py
create_LDAS_veg_param.c   meteoCombined       regrdTmin.runfile   smseasonal.scr          vic_soil_moisture.py
data_US_SLAD_2010.csv     mk_mnth             regrd_wind          soil                    vic_spatiotempdatabase.py
default                  mk_monthly          regrd_wind.c        spatiotempdatabase.py   wind_latlong.txt
DEM.tif                 mk_monthly.c        rescale             spatiotempdatabase.pyc  write.py
DSMwithpopulation.py     mk_monthly_ir.scr   rescale.c
```


Organized in space



```
irods@ec2-54-86-215-185:~/IRODS/server/bin$ cd cmd/
irods@ec2-54-86-215-185:~/IRODS/server/bin/cmd$ ls
adjprcp_rsm.rsc      forcingData      old              rescale_ir.scr      stationinfo.txt
adjtmin_grid.grd    gauss_t62_lat.list  osgeo.bak       run_combine_wind_ir.scr  temp_wind
adjtmin_grid.grd    gauss_t62_lon.list  output_1.txt    run_combine_wind.scr   TerraPopCountieswVicOutput.py
adjtmin_grid.grd    hello.py          PopulationVsSoilMoisture.scr  run_convert_prcp.scr  test
adjtmin_grid.grd    inputPrism.c      popvssm.run     run_convert_tif_ascii.scr  test3.py
adjtmin_grid.grd    inputTmax.scr     prism           run_convert_tmax.scr    test_execstream.py
adjtmin_grid.grd    inputTmin.scr     prism-rawdata   run_convert_tmin.scr    test.py
adjtmin_grid.grd    irodsAgent        python_script1.py  run_get_prism.scr      test.scr
adjtmin_grid.grd    irodsReServer     read_prec_dly    run_mk_monthly.scr      tiff2ascii.py
adjtmin_grid.grd    irodsServer        read_prec_dly.f  run_psp_vic_sm.scr      tmax
adjtmin_grid.grd    irodsServerMonPerf  read_prec_dly.f  run_psp_vic_soilmoisture1.scr  tmax_tob_adj
adjtmin_grid.grd    irodsXmsgServer    read_prec_dly.f  run_psp_vic_soilmoisture2.scr  tmax_tob_adj.f
adjtmin_grid.grd    latlong90.txt      read_prec_dly.f  run_psp_vic_soilmoisture3.scr  tmax_tobAdj.scr
adjtmin_grid.grd    latlon.txt         read_prec_dly.f  run_psp_vicSoilMoistureComparison.scr  Tmin
adjtmin_grid.grd    LDAS               read_prec_dly.f  run_psp_vic_soilmoisture.scr  tmin_tob_adj
adjtmin_grid.grd    LDAS               read_prec_dly.f  run_psp_vic_soilmoisture.scr  tmin_tob_adj.f
adjtmin_grid.grd    LDAS               read_prec_dly.f  run_regrid_wind_ir.scr      tmin_tobAdj.scr
adjtmin_grid.grd    LDAS               read_prec_dly.f  run_regrid_wind.scr        tob
adjtmin_grid.grd    LDAS               read_prec_dly.f  run_rescale.scr            univMSSInterface.sh
adjtmin_grid.grd    LDAS               read_prec_dly.f  run_vicinput_ir.scr        uploadTo53.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  run_vicinput.scr          vacuumdb.pl
adjtmin_grid.grd    LDAS               read_prec_dly.f  script_1.scr              vegetation
adjtmin_grid.grd    LDAS               read_prec_dly.f  script1.scr               vic_calc_mnth_mcl.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  script2.scr               vic_calc_mnth_mc2.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  script3.scr               vic_calc_mnth_mc.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  script4.scr               vic_calc_mnth_pcb.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  script5.scr               vic_evapotranspiration.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  script6.scr               vicinput
adjtmin_grid.grd    LDAS               read_prec_dly.f  script7.scr               vicinput.c
adjtmin_grid.grd    LDAS               read_prec_dly.f  script_save               vic_monthlyPEBestimates.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  script_test               vic_monthly_soilmoisture2.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  script_vic_pre-processing.scr  vic_monthly_soilmoisture.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  sm_comparison.py          VIC_Pre_processing_script.scr
adjtmin_grid.grd    LDAS               read_prec_dly.f  sm_seasonal.py            vic_soil_moisture3.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  sm_seasonal.py            vic_soil_moisture.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  sm_seasonal.py            vic_spatialtempdatabase.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  soil                      wind_lat long.txt
adjtmin_grid.grd    LDAS               read_prec_dly.f  spatiotempdatabase.py     write.py
adjtmin_grid.grd    LDAS               read_prec_dly.f  spatiotempdatabase.pyc
```

geounit 1

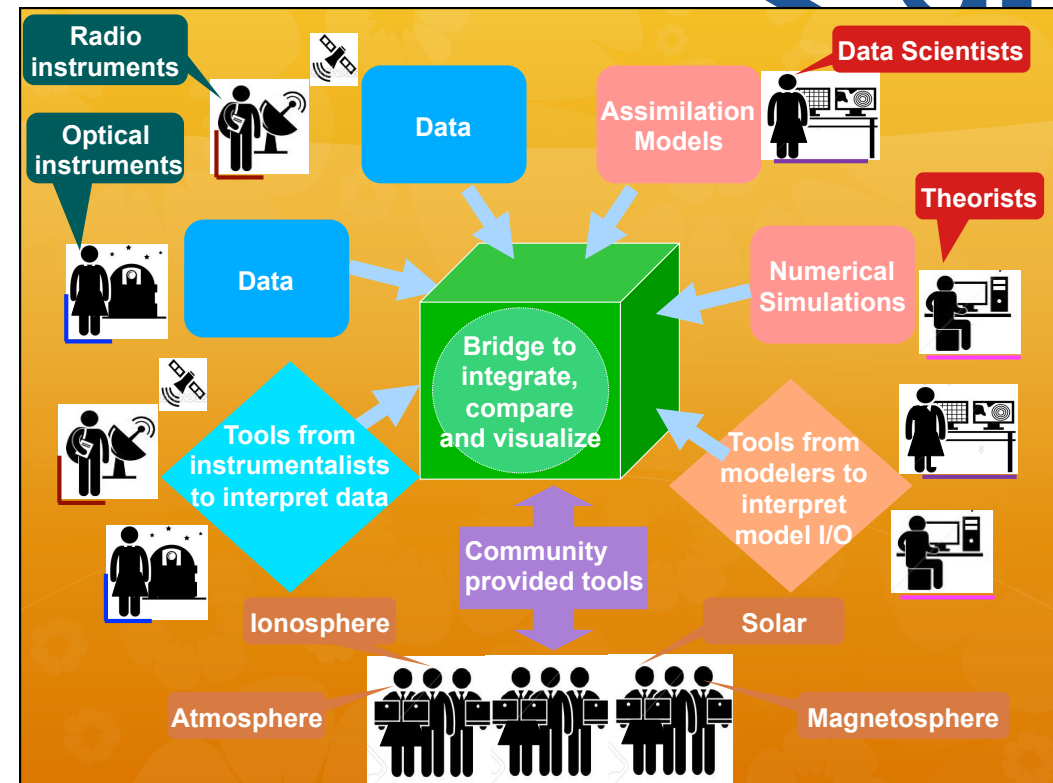
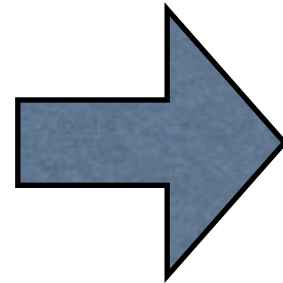
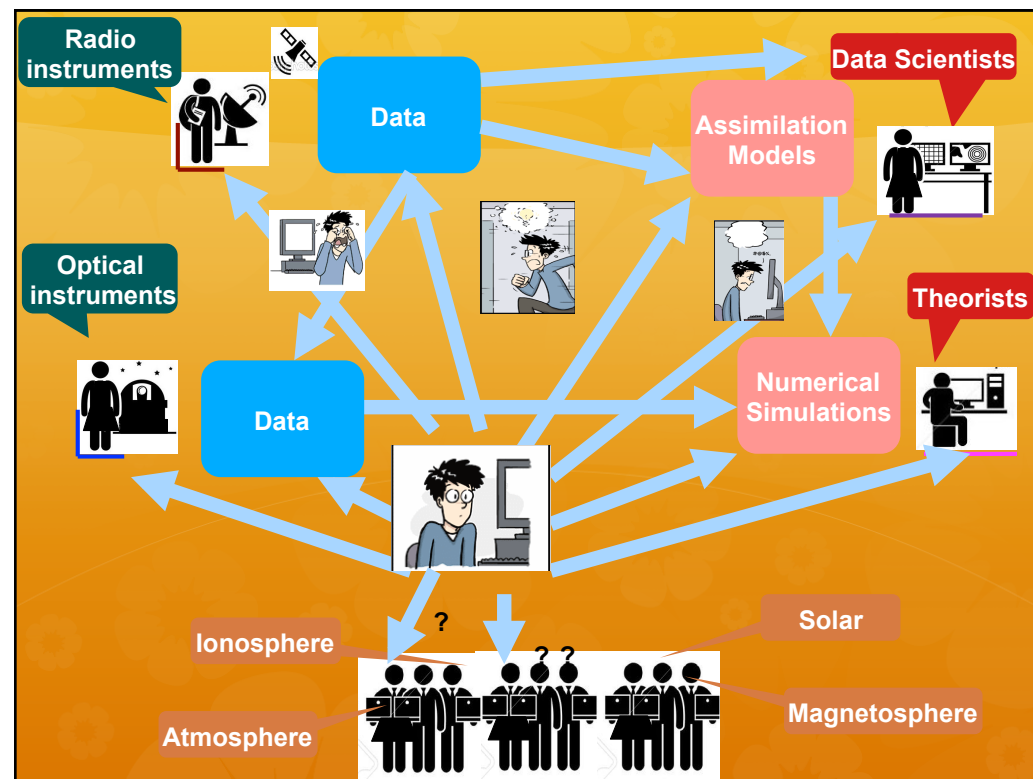


geounit 2



Space Science

(Model and Data Integration)



InGeo: Integrative GeoScience Observatory

EarthCube Integrative Activity: Asti Bhatt,
Russell Cosgrove (SRI International)

Summary

- Capability-rich tools
 - Self-curation
 - Tracking cause and effect between data and models
 - Ensuring reproducibility
 - Support services: discovery, sharing, publication
- **Simplify** model and data management for computational/data geoscientists
- **Faster and reproducible hand-shakes** among faculty and students
- **Cuts down IT support** for modeling centers and real science



Track it!

- Science Usecases, Reports, Presentations, News
 - <http://workspace.earthcube.org/geodataspace>
- Source Code (Public Release Pending)
 - <http://github.com/TanuMalik/SciDataspace/Geo>
- The GeoDataspace (Internal URL)
 - <https://scidataspace.org/geo>

Acknowledgements



EARTH CUBE



globus



- National Science Foundation
- EarthCube Community
- Globus team
- CI team