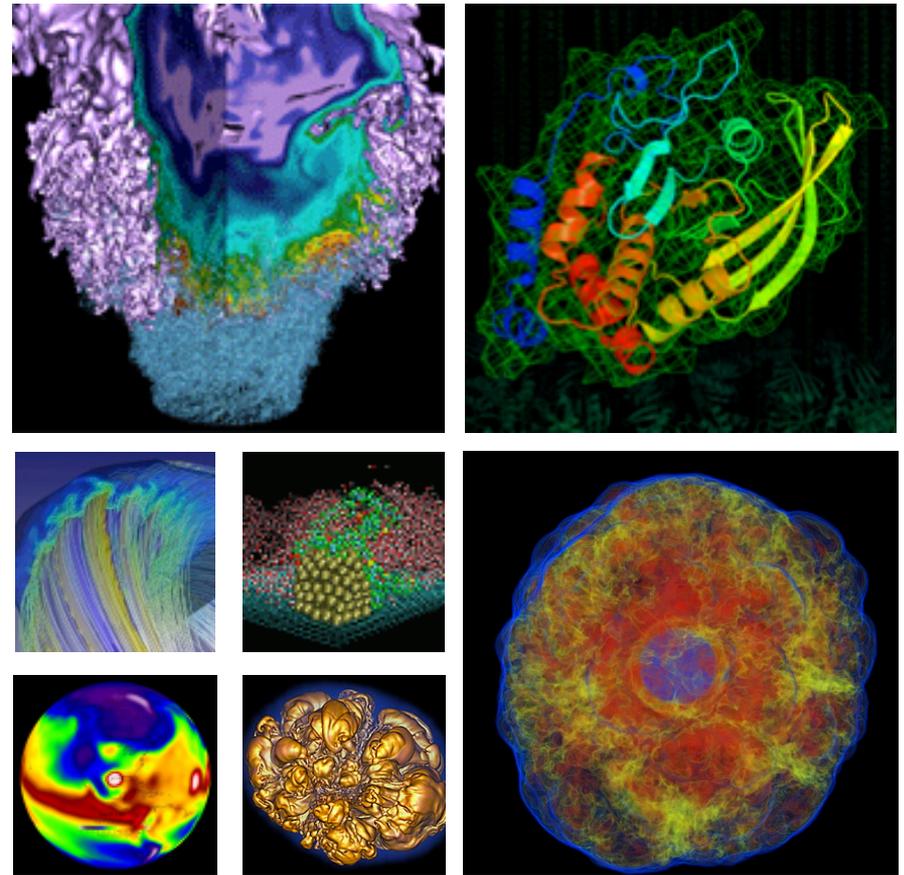


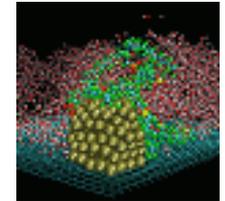
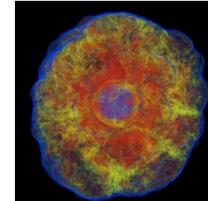
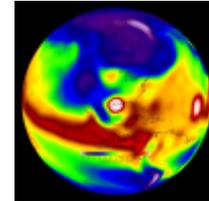
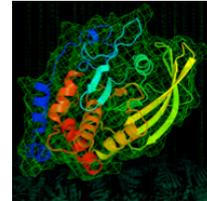
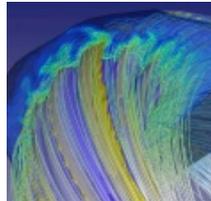
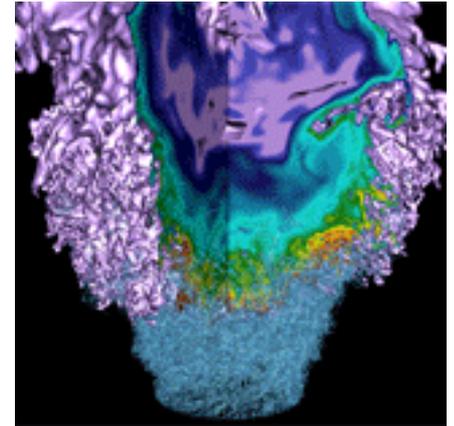
Design Patterns in Web Gateways for Scientific Data



Shreyas Cholia
Data and Analytics Services, NERSC

Joint Facilities User Forum on Data-Intensive
Computing, Oakland, June 17, 2014

Introduction



- **Motivation**
- **Patterns**
 - Web Data Sharing
 - End-to-End Frameworks
 - API-driven access
 - Web IDEs
 - Federated Identity
 - Social Data
- **Conclusions**

Motivation for Science Gateways



- People expect web interfaces and applications for usability

```
ssh carver.nersc.gov  
qsub -I -q matgen_reg -l nodes=2:ppn=16 -l walltime=00:30:00  
cd /project/projectdirs/matgen/vasp_test  
mpirun -n 32 vasp|
```

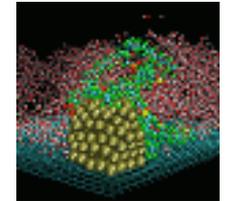
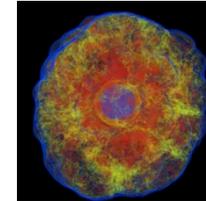
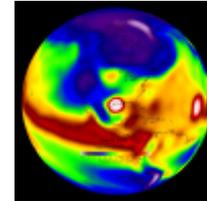
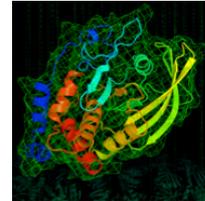
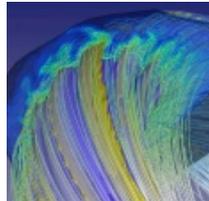
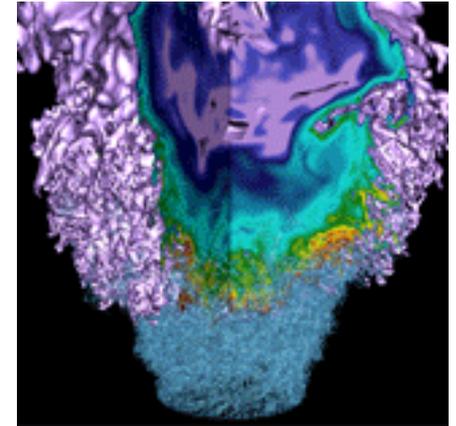
- ✗ don't want generic options/tools not applicable to your science
- ✗ don't want to deal with backend, middleware, UNIX CLI etc.

Science is Collaborative



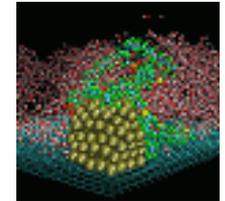
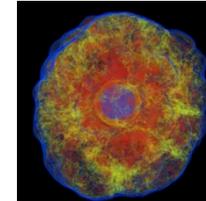
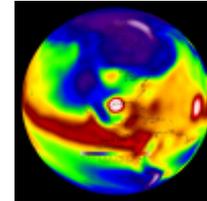
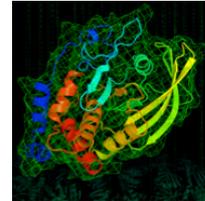
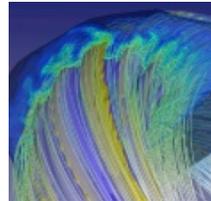
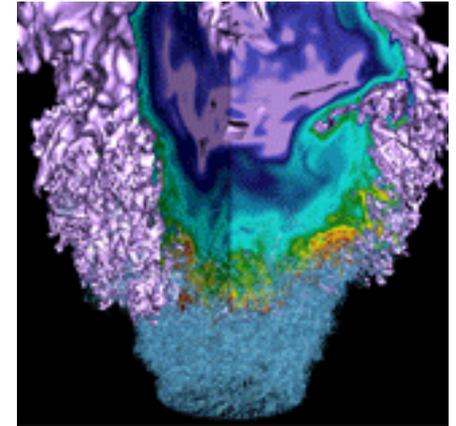
- Science is now a collaborative effort
- Large teams of people
- Requires shared access to compute and data resources

Patterns in Science Gateways



Attempt to catalog the interesting features of science gateways that span science domains in order to gain some insight on how to best expose scientific data over the web.

Data Sharing Over The Web

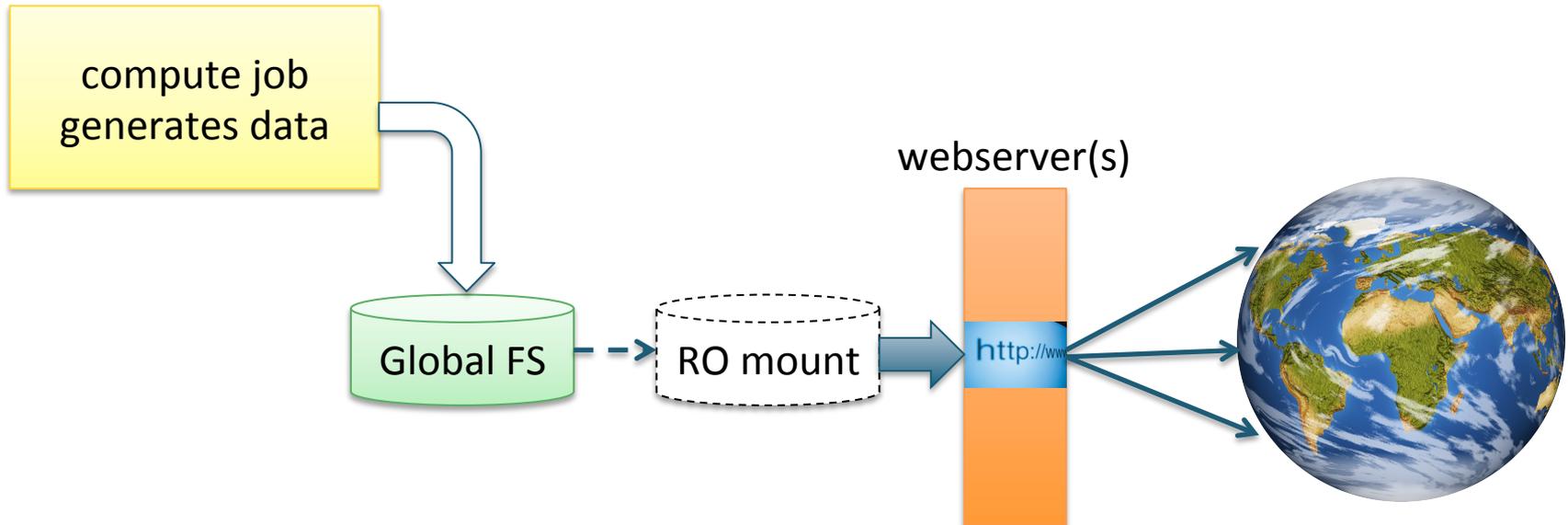


Data Sharing Over The Web



- **OSTP requirements around data management**
- **What is the basic infrastructure required to make this happen?**

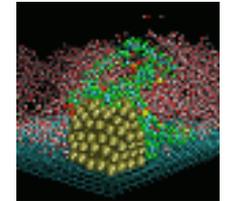
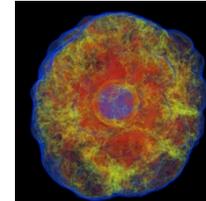
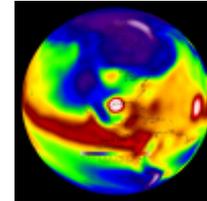
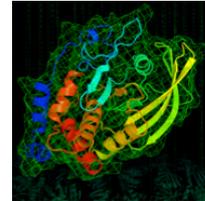
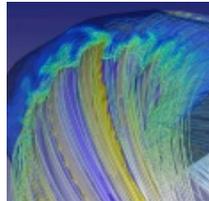
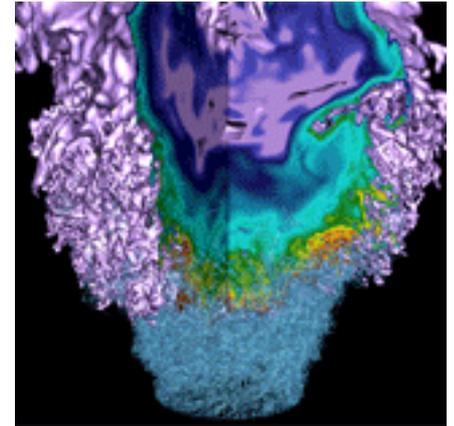
Data Sharing Over The Web



- **Store data on persistent global filesystem**
- **Designate an area for export**
- **Mount export area RO on a web node using a shared FS**
- **Run a webserver on web node to share with world**
- **Add AuthNZ to web server configuration as needed**

- **Scale out – add more nodes + loadbalancer**
- **Use sophisticated interfaces to search/browse data**
 - HTML+JS, PHP, Web Frameworks
- **High performance data transfer tools**
 - Globus
 - Xrootd + webdav
 - Custom downloaders

End-to-End Frameworks



- **Open source frameworks that can be customized for specific science use cases:**
 - HubZero
 - Galaxy
 - OGCE
- **Takes a plugin + configuration based approach to domain specific customizations.**
- **Very useful for managing canned workflows.**
- **Community and management tools often baked-in**

The screenshot shows the Galaxy workflow editor interface. The browser address bar displays `http://localhost:8080/workflow/editor?id=adb5f5c93f827949`. The main navigation bar includes 'Analyze Data', 'Workflow', 'Data Libraries', 'Help', and 'User'. On the left, a 'Tools' sidebar lists various bioinformatics tools such as 'GMOD 2010 Course Tools', 'Get Data', 'Send Data', 'ENCODE Tools', 'Lift-Over', 'Text Manipulation', 'Filter and Sort', 'Join, Subtract and Group', 'Convert Formats', 'Extract Features', 'Fetch Sequences', 'Fetch Alignments', 'Get Genomic Scores', 'Operate on Genomic Intervals', 'Statistics', 'Graph/Display Data', 'Regional Variation', 'Multiple regression', 'Multivariate Analysis', 'Evolution', 'Metagenomic analyses', and 'FASTA manipulation'. The central 'Workflow Canvas' shows a workflow titled 'Workflow constructed from history 'Unnamed history''. It consists of three steps: two 'Input dataset' steps, a 'Map with BWA' step, and a 'SAM Filter' step. The 'Map with BWA' step has two inputs: 'Select a reference from history' and 'FASTQ file'. The 'SAM Filter' step has one input: 'File to filter'. The 'SAM Filter' step is currently selected, and its configuration is shown in the 'Details' panel on the right. The 'Details' panel for the 'SAM Filter' tool includes the following options: 'File to filter' (Data input 'input1' (sam)), 'Optional field to filter on' (set to 'Edit Distance'), and 'Value to require for flag' (set to '1'). Below the 'Details' panel is the 'Edit Step Attributes' section, which includes an 'Annotation / Notes' text area and a prompt: 'Add an annotation or notes to this step; annotations are available when a workflow is viewed.'



National Integrated Drought Information System

U.S. Drought Portal

www.drought.gov

[Contact Us](#) | [Log In](#) | [Text-Only](#)

Search:

HOME
WHAT IS NIDIS?
CURRENT DROUGHT
FORECASTING
IMPACTS
PLANNING
EDUCATION
RESEARCH
RECOVERY
REPORTS

Area Drought Information

Select State...

Select Region...

Maps & Tools

- [Map & Data Viewer - new!](#)
- [Geodata Portal](#)
- [Drought Monitor Graphics](#)
- [CRN Soil Data](#)

Events & Announcements

- [Carolinas DEWS Scoping Workshop July 2012](#)
- [April 2012 Southern Plains Drought Assessment & Outlook Forum](#)
- [Risk Management Meeting 11/2011](#)
- [ACF Climate Outlook Forum and Pilot Review Meeting 2011](#)
- [Engaging Preparedness Communities Webinar, Dec 13th 1-2 PM EST](#)

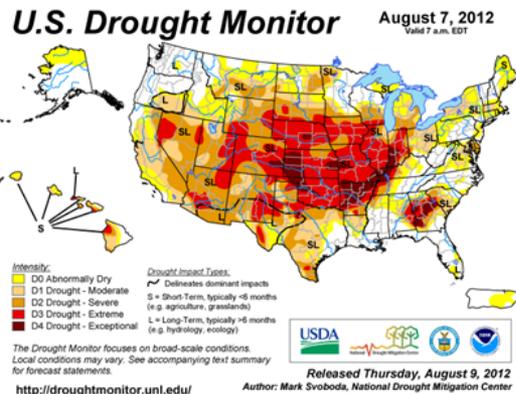
[View Archive](#) | [Portal Release Notes](#)

Regional Drought Webinars

- [ACFRB Drought Assessment Webinar - July 3rd, 2012](#)
- [ACFRB Drought Assessment Webinar - June 19th, 2012](#)
- [ACFRB Drought Assessment Webinar - May 22nd, 2012](#)
- [ACFRB Drought Assessment Webinar - February 28th, 2012](#)
- [ACFRB Briefing Presentation - January 17th, 2012](#)
- [Managing Drought in the Southern Plains](#)

Featured Products

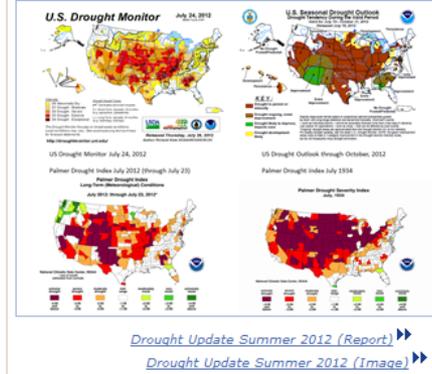
- [Where are Drought Conditions Now?](#)
[How is the Drought Affecting Me?](#)
[Will the Drought Continue?](#)



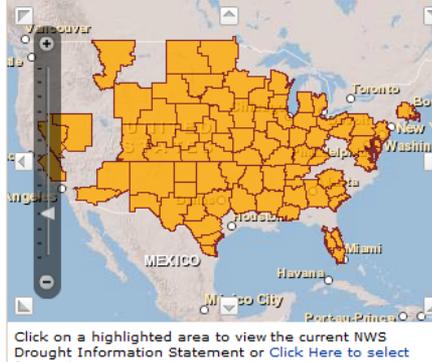
Regional Drought Early Warning Systems (DEWS)



NIDIS Feature

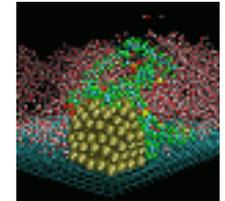
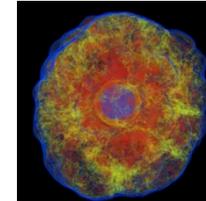
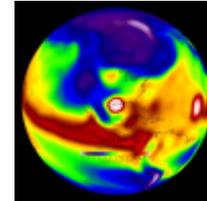
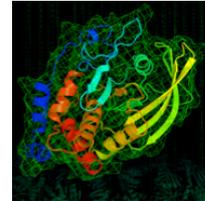
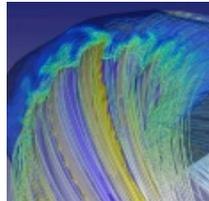
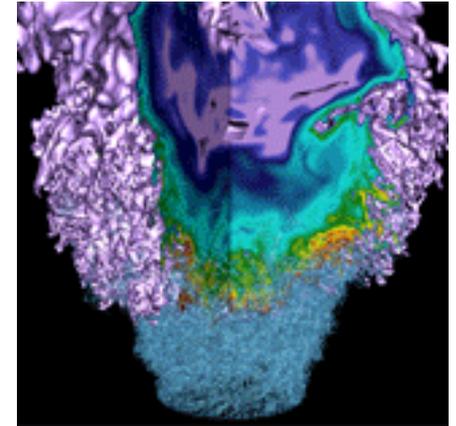


Drought Information Statements



- **Plugin/Portlet approach to development requires developers to learn the vernacular of the framework.**
- **Applications tend to be built around framework constraints rather than science use-cases. Rich set of front-end technologies in HTML5 and JavaScript much harder to leverage in this context.**
- **If you want a blank-slate that starts with science use cases and builds up interface from here, you might try a different approach.**

API Driven Access



- **Web APIs provide programmatic access to data and resources to developers over the web**
- **Access to data as well-defined objects allows users to develop their own custom applications and code**
- **Build web application around this API using the software stack of your choice**

The 1-minute intro to REST APIs



- Use HTTP verbs in conjunction with URLs
 - GET, POST, PUT, DELETE
- Verb + URL + parameters = function call
- Return structured data
- For instance (from NEWT API):

VERB	RESOURCE	DESCRIPTION
POST	/queue/R	Submits POST data to queue on R; returns job id
GET	/file/R/path/	Returns directory listing for /path/ on R
GET	/account/user/U	Returns user account info for U
DELETE	/queue/R/id	Deletes jobs from Queue

- Structured Output (JSON):

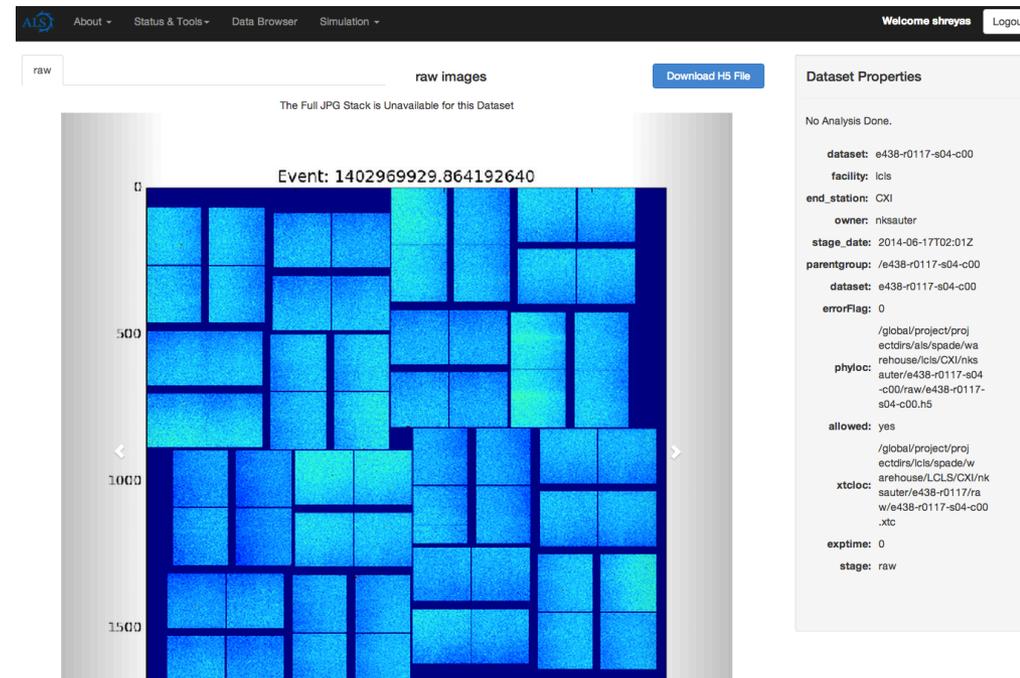
```
{"status": "OK",  
"output": [{"status": "up", "system": "hopper"}, {"status": "up", "system": "carver"}]  
"error": ""}
```

APIs Drive Programmatic Access to Data

The NERSC logo is a dark blue rounded rectangle with the word "NERSC" in white, bold, sans-serif font. Behind the text, there are several light blue lines radiating outwards from the center, creating a starburst or sunburst effect.

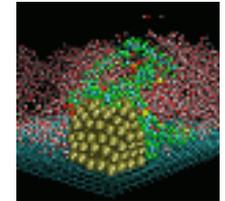
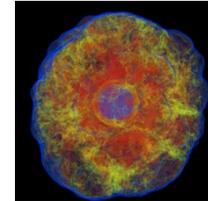
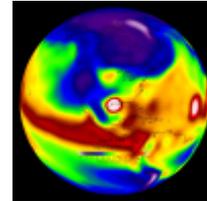
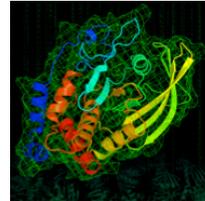
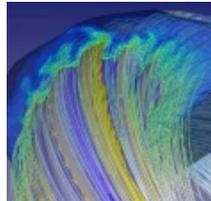
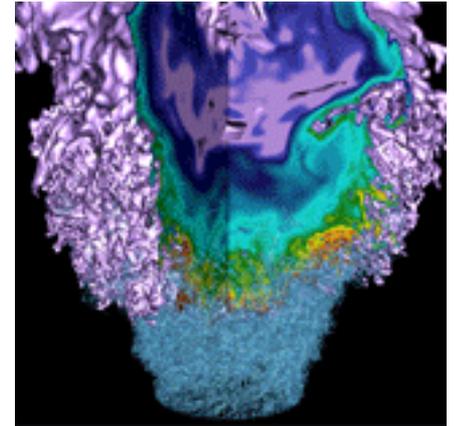
- Enable heavy duty server side operations and long running tasks
- Perform sub-selection of data – pull down only what is needed in client and render it in the browser
 - eg. OpenDAP connects HDF5 datasets with web clients

- All the client side logic resides in HTML5/JavaScript
- Make callouts to REST API for all server-side information

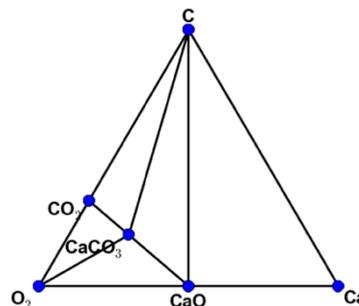
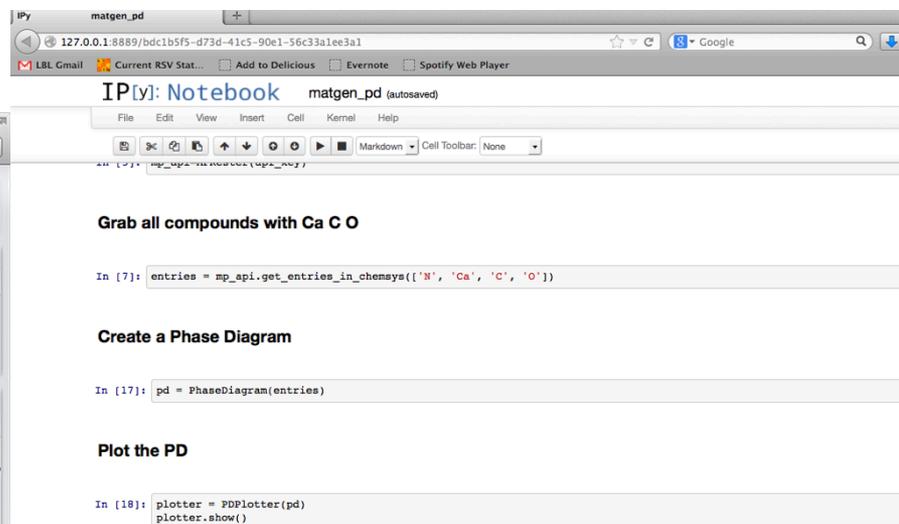
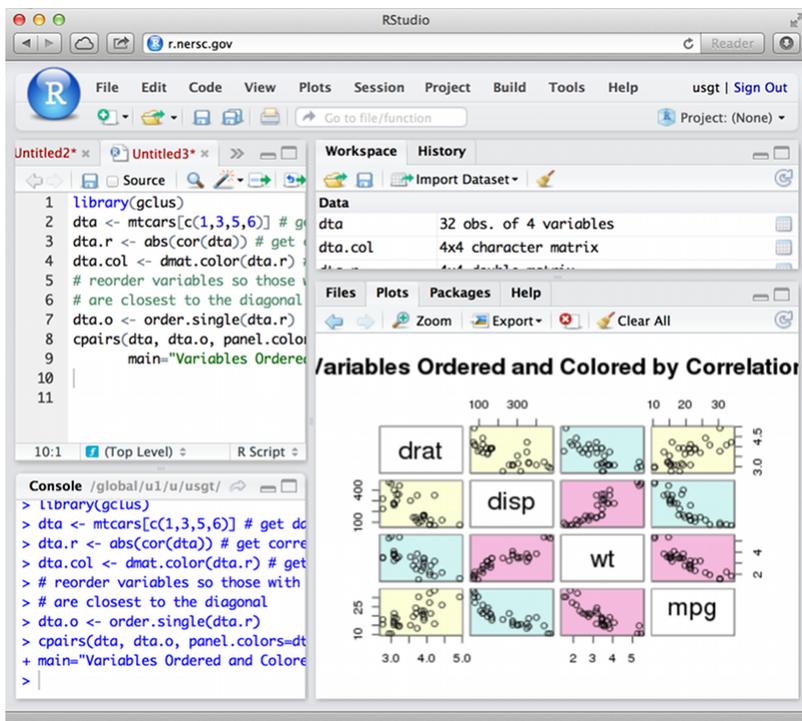


- **Increasingly, web frameworks are following Model-View-Controller pattern**
 - eg. Django, Ruby-On-Rails, Backbone.js
- **Enables separation of user interfaces (V) from backend data models (M) and scientific libraries for interacting with data (C)**
- **Provide you with just enough scaffolding to build out application without imposing front-end application constraints.**
- **Web APIs lend themselves nicely to this model**

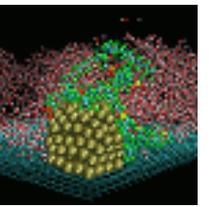
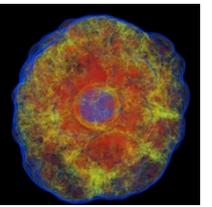
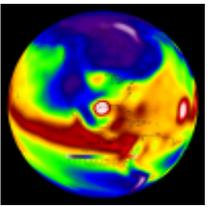
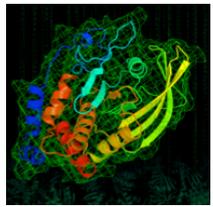
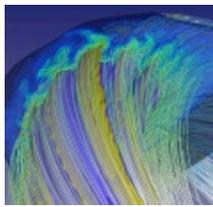
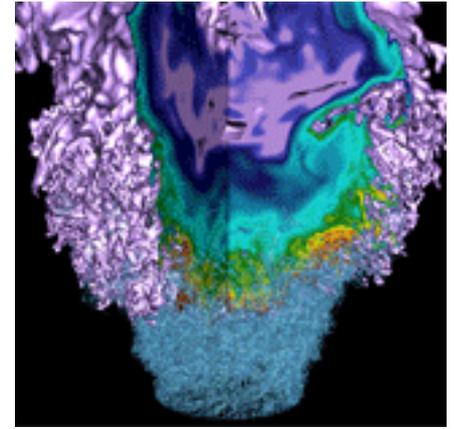
Web IDEs



- Combine the expressiveness of command line tools with web GUIs
- RStudio, Mathematica, iPython Notebook



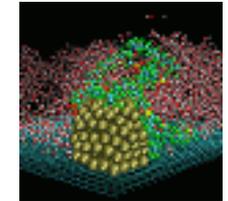
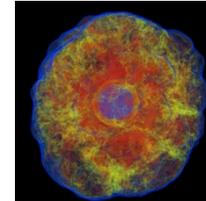
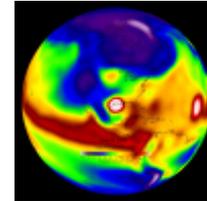
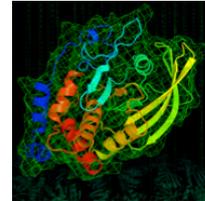
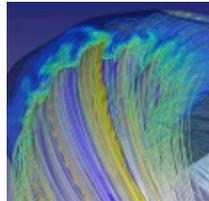
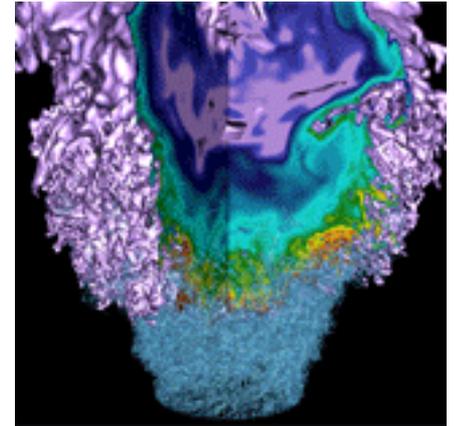
Federated Identity



- **Authentication: verifying the identity of the user**
 - Are you who you say you are?
- **Local authentication**
 - Identity tied to local user credentials from institution that hosts the gateway
- **Federated authentication**
 - Identity verification delegated to external party (such as user's home institution, or verified 3rd party)
 - External party returns an assertion to the gateway with ID information.

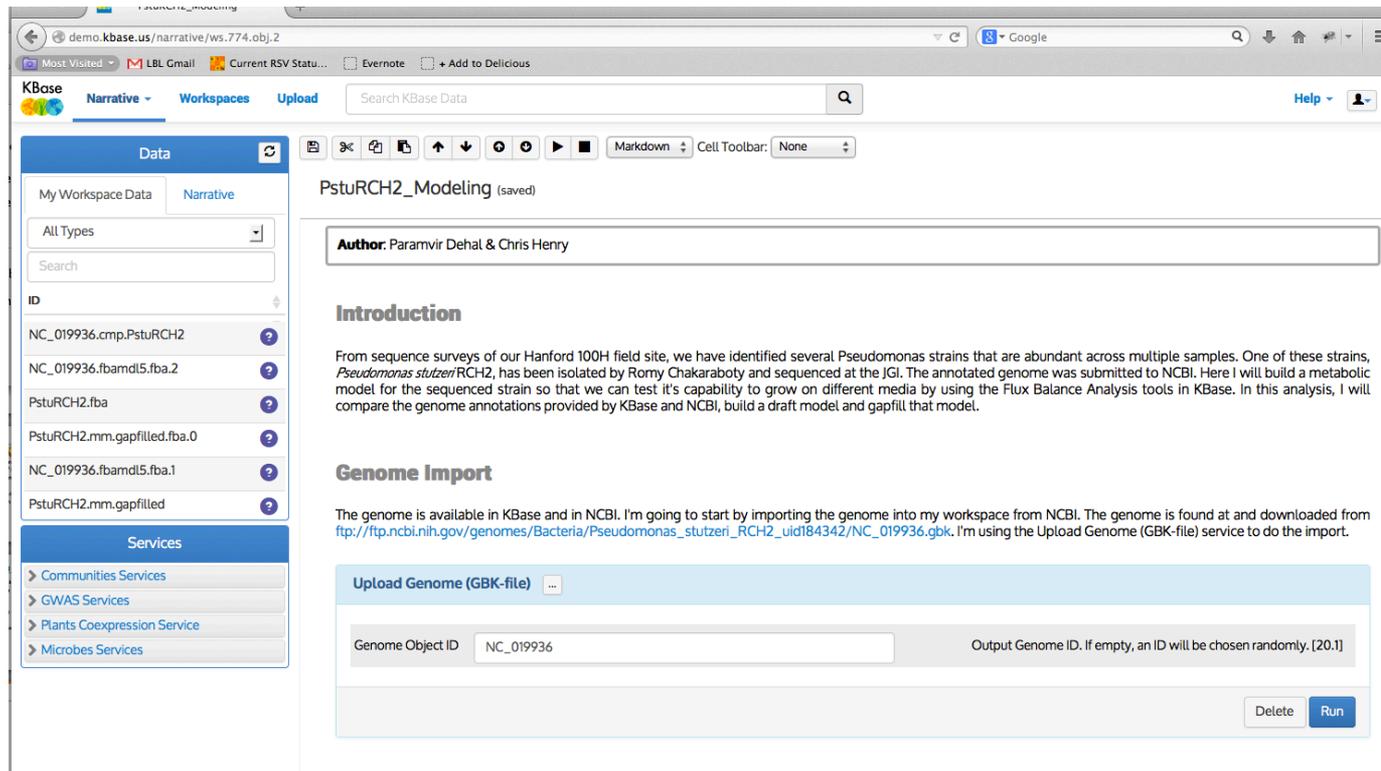
- **One identity to rule them all** 
- **Allow users to log-in using common credentials from their home institution / 3rd party provider**
- **Assumes trust relationship between science gateway and identity provider**
- **Map federated ID to local account**
- **Examples:**
 - Globus Nexus (OAuth), CILogon (Shibboleth)
- **Enables access across users from multiple institutions**

Social Data



- **Generate a discussion around science results**
- **Combine Science Gateways with social interaction tools**
 - Disqus
 - Forums integrated and cross-linked with datasets
 - Social Media integration
- **Referencing and sharing results and user generated datasets**
- **Users can upload datasets**

- Allows you to reference data and run models from other users in system



The screenshot shows a web browser window displaying the KBase Narrative interface. The URL is `demo.kbase.us/narrative/ws.774.obj.2`. The interface includes a navigation bar with 'Narrative', 'Workspaces', and 'Upload' tabs, and a search bar for 'Search KBase Data'. A left sidebar contains a 'Data' panel with 'My Workspace Data' and 'Narrative' tabs, a search box, and a list of data objects including `NC_019936.cmp.PstuRCH2`, `NC_019936.fbamd15.fba.2`, `PstuRCH2.fba`, `PstuRCH2.mm.gapfilled.fba.0`, `NC_019936.fbamd15.fba.1`, and `PstuRCH2.mm.gapfilled`. Below the data list is a 'Services' section with links for 'Communities Services', 'GWAS Services', 'Plants Coexpression Service', and 'Microbes Services'. The main content area shows a narrative titled 'PstuRCH2_Modeling (saved)'. It includes an 'Author' field with 'Paramvir Dehal & Chris Henry', an 'Introduction' section with text about *Pseudomonas stutzeri* RCH2, and a 'Genome Import' section with text about importing a genome from NCBI. At the bottom of the 'Genome Import' section is an 'Upload Genome (GBK-file)' form with a 'Genome Object ID' field containing 'NC_019936', an 'Output Genome ID' field, and 'Delete' and 'Run' buttons.

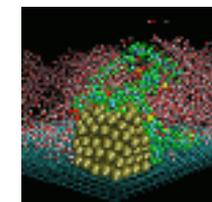
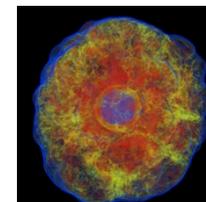
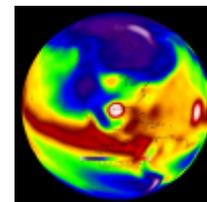
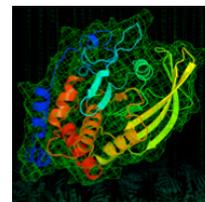
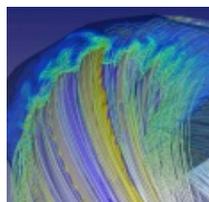
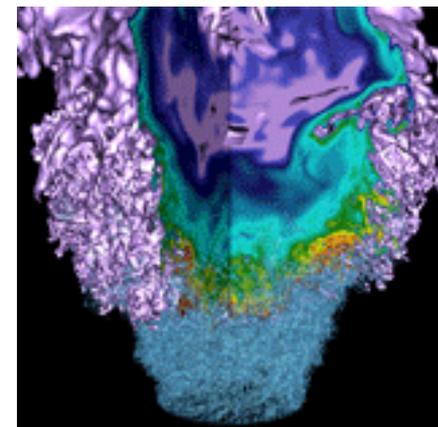
Data Upload and Provenance



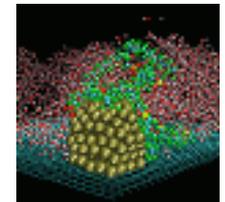
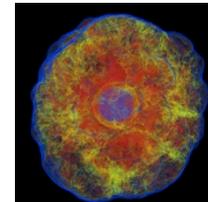
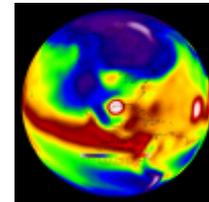
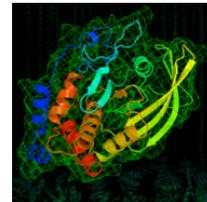
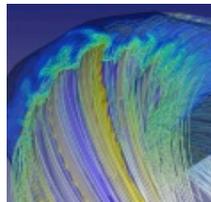
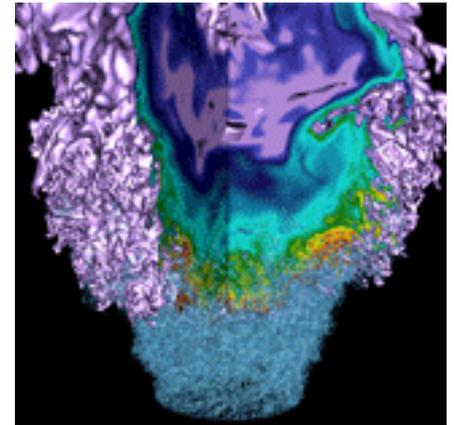
- **Users upload data to portals**
- **But this introduces challenges wrt data management**
 - Vetting/Validating the data
 - Managing Provenance
 - Data ownership and licensing
 - Public/Private/Protected data

- **Probably not a complete list**
 - Intended to spur further discussion on other interesting patterns
- **NERSC Science Gateways**
 - <http://portal.nersc.gov>
- **Contact:**
 - Shreyas Cholia: scholia@lbl.gov
- **Questions? Comments?**

Fin



Extra



- **Authorization: Access privileges for a given identity**
 - What is the user allowed to do?
- **Data portals see some combination of the following sharing modes:**
 - Public – data are publicly viewable to all users
 - Protected – data are only visible to members within a designated group or collaboration
 - Private – data are private to the individual user