# MyEMSL & EMSLHub:
## Creating A flexible framework for scientific data sharing, discovery, and collaboration
PNNL-SA-96303

David Cowley
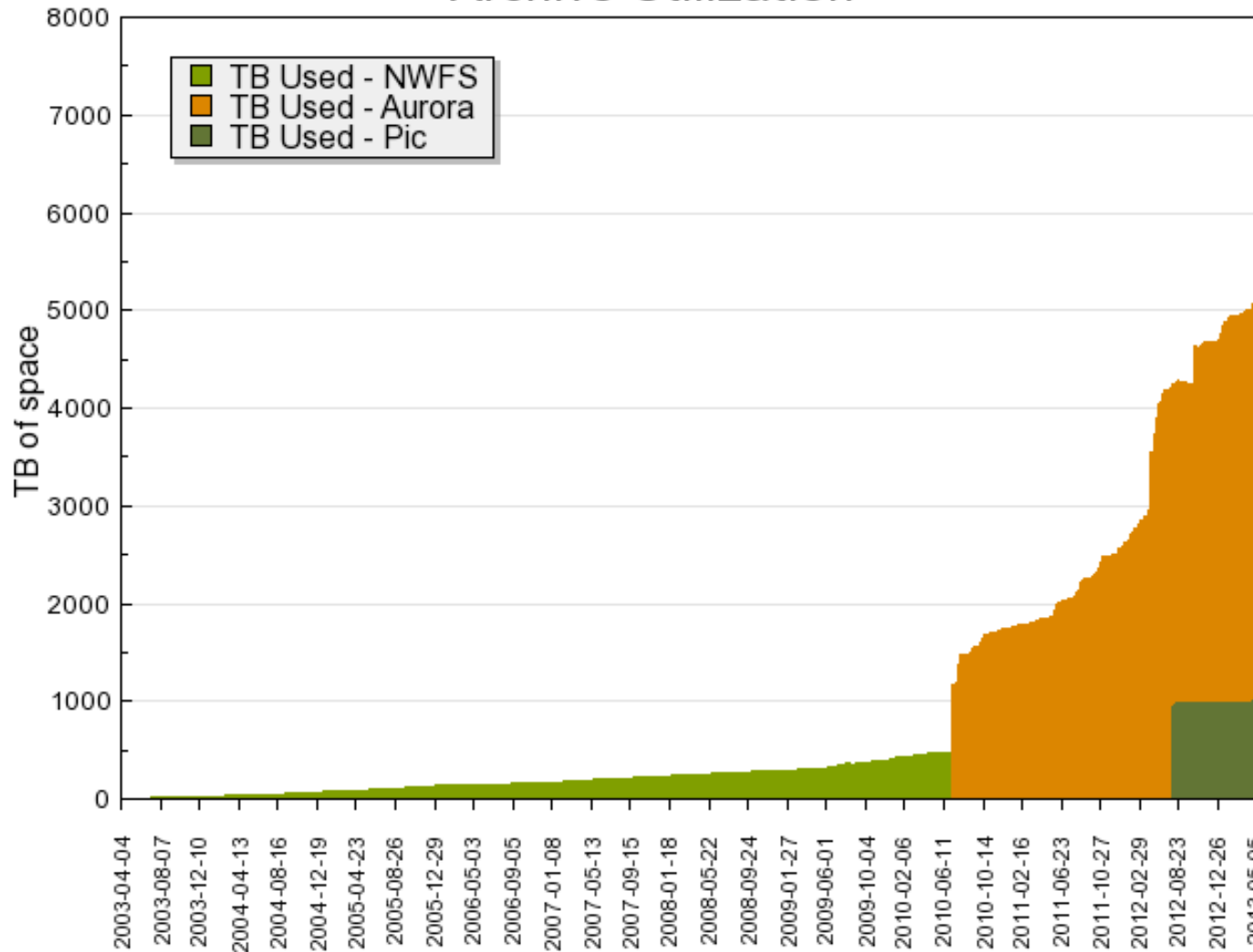
**Pacific Northwest**
NATIONAL LABORATORY

U.S. DEPARTMENT OF
**ENERGY**

*Proudly Operated by* **Battelle** *Since 1965*

Archive Utilization

# EMSL Estimated Aggregate Data Production/Transfer Rates

- **At present:**
  - 6 TB/day produced in EMSL
  - 5 TB/day transferred within EMSL
  - 200 GB/month transferred out of EMSL

- **2-5 years from now:**
  - 20 TB/day produced
  - 40 TB/day transferred within
  - 600 TB/month transferred
  - 5 TB/month transferred into EMSL

- **5+ years from now:**
  - 100 TB/day produced
  - 200 TB/day transferred within
  - 3 PB/month transferred out
  - 50 TB/month transferred in

# New Directives: Expand Public Access to the Results of Federally Funded Research

- *"The Obama Administration is committed to the proposition that* citizens deserve easy access to the results of scientific research their tax dollars have paid for. *That's why, in a policy memorandum released today, OSTP Director John Holdren has* directed Federal agencies with more than $100M in R&D expenditures to develop plans to make the published results of federally funded research freely available to the public within one year of publication and requiring researchers to better account for and manage the digital data resulting from federally funded scientific research."*

- -- Posted by Michael Stebbins on whitehouse.gov February 22, 2013 at 12:04 PM EDT

# Continuum of Data Sharing

**Nobody**
- Disks in Drawers
- Thumb Drives

**Colleagues/ Workgroup**
- Shared Filesystems
- Desktop Shares
- FTP

**Institution**
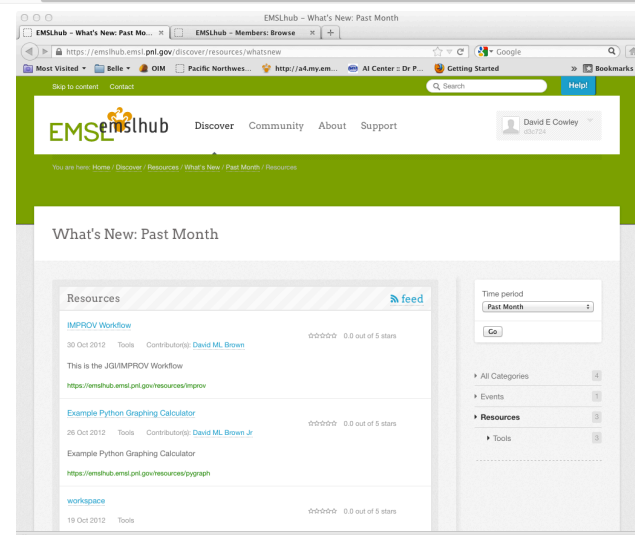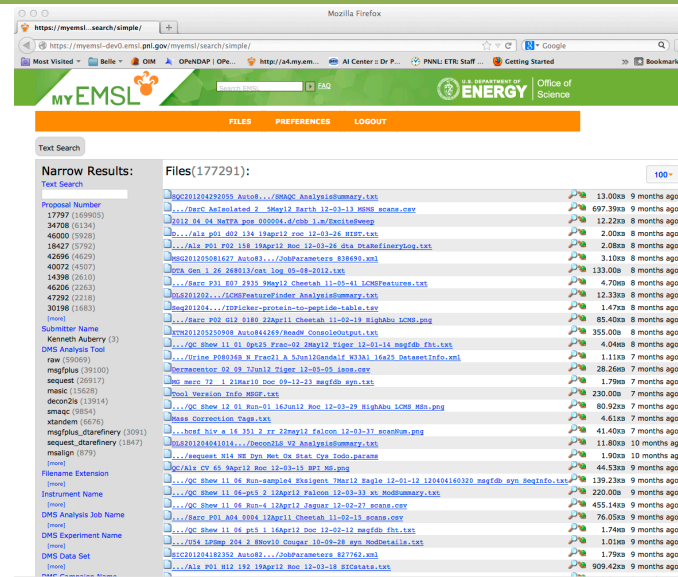- Archive Systems
- Hard drives via FedEx

**Community/ World**
- Cloud
- Semantic Web

Pacific Northwest
NATIONAL LABORATORY

U.S. DEPARTMENT OF
ENERGY

*Proudly Operated by* **Battelle** *Since 1965*

- **Aurora** is EMSL's current archive system
- It has lots of parts:
  - Multi-petabyte robotic tape library (in CSF)
  - 1 Petabyte Disk Storage (in EMSL)
  - Many servers
  - Many services
- It's mostly been used as a big ol' pile of files
  - Organization (if any) was strictly DIY
  - Hard to find anything or even know it's there!
- Access from outside EMSL is extremely limited

# MyEMSL: Data Management for a diverse user facility

- Provide a system to capture, store and share data & metadata from EMSL's instruments and computers

- Create a <u>flexible</u> data framework that supports multiple scientific fields & approaches

- Allow scientific collaboration around data via EMSLHub (aka HubZero):
  - ◆ Workflow-based computation
  - ◆ Data integration and analysis

- Leverage high quality open source software
  - ◆ Form open source development groups
  - ◆ Easily shared with community
  - ◆ Highly efficient support and maintenance model

- Long term, create a framework to federate data sources across institutions

# Supporting Scientific Collaboration with MyEMSL

- Provide search and easy access to data & metadata for authorized users

- Facilitate and enforce data release policies

- Find enthusiastic scientists who know they have data needs and partner with them

- Provide data interaction framework to facilitate "Bring Your Own":
  - Data catalogs & definitions
  - Workflow based computation
  - Data integration and analysis
  - Web and mobile apps

# Why MyEMSL?

- **MyEMSL** aims to solve a number of problems with the old archive approach:
  - Make data easy to find by good use of metadata
  - Make files easy to get via the web
  - Set up good access control
  - Make it easy to share data with collaborators
- Leverage existing HPSS archive for bulk data storage
- Augment legacy archive with other parts:
  - Metadata Database
  - Automated uploader for instruments
  - Web portal @ https://my.emsl.pnl.gov
  - Links to EMSL's User/Resource allocation systems & EMSLHub

- We have a big pile of files going back 15 years, but nobody knows what all is in there!

- The key to discovering and understanding data is to have good **metadata!**

- Metadata tells you what is known about the data

- Can be trivial:
  - File name
  - File owner
  - Time and date stamps

- Can be much more interesting:
  - Instrument settings & environmental parameters
  - Sample information
  - Registration data for multimodal analyses

# What we need to do about metadata

- The more we can get of it, the more powerful MyEMSL becomes

- We have defined a little essential metadata that all data must have

- MyEMSL must be **flexible and non-restrictive** about metadata and **enable metadata to evolve** over time

- We don't know whatever other metadata people may need – maybe they don't know yet

- Science partners will need to tell us what metadata will be useful to them and MyEMSL will need to support it

- How can we get it?
  - Our processes (i.e. autouploader) can capture or generate a little of it
  - Humans may need to enter it
  - Code can be written to extract/generate it, but this will require funding and people time!

# What exists today



- Web and service interfaces for moving data and metadata into and out of MyEMSL

- Authentication mechanisms

- Search and "shopping cart" interfaces

- EMSLHub collaboration site with pilot workflows

- Proteomics pipeline adapted to MyEMSL interfaces

- Increasing numbers of scientific instruments uploading data

# MyEMSL Conceptual Diagram

# Data, Instruments in EMSL context



**EMSL**

EMSLHub

EMSL Users

MyEMSL

EMSLfs

EUS

Archive (Aurora)

Internet

EMSL Network

IDLNet

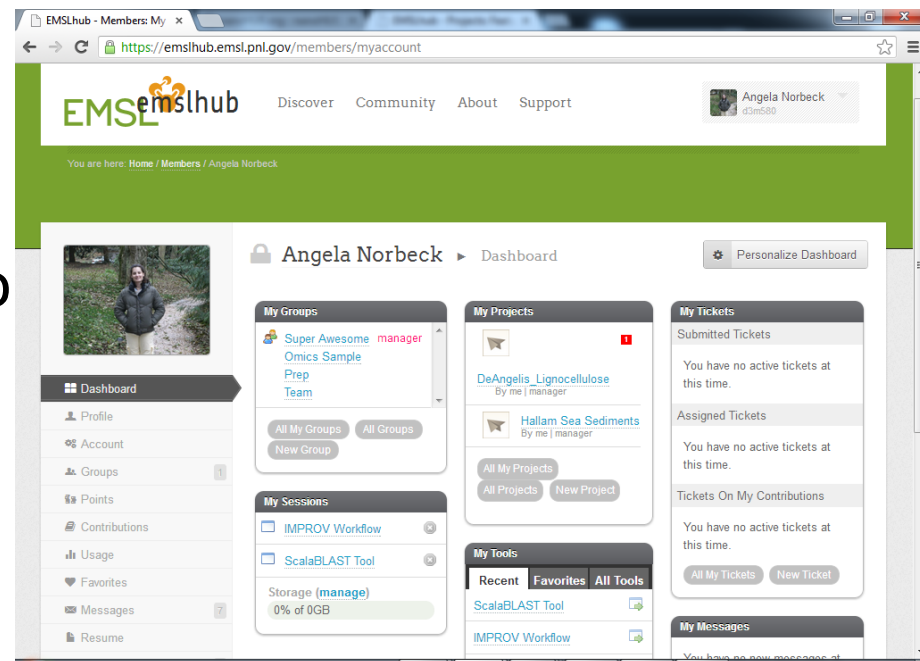# http://my.emsl.pnl.gov

# Example: Public data in MyEMSL web pages

# Establishing MyEMSL Collaborations

- We are establishing key collaborations with other institutions to:
  - Define tools and interfaces for data exchange
  - Leverage existing tools and frameworks
  - Establish open source consortia
- Partner organizations and activities:
  - DOE Joint Genome Institute (JGI)
  - DOE Systems Biology Knowledgebase (KBASE)
  - National Center for Microscopy and Imaging Research (NCMIR):
  - Semantic Physical Sciences (SPS) workshops

# EMSLHub



- **EMSLHub** is a collaborative web site at [https://emslhub.emsl.pnl.gov](https://emslhub.emsl.pnl.gov)

- Purpose: provide a system to allow scientific collaboration around *data* with a rich toolset

- Early collaborative efforts:
  - JGI – EMSL Data Integration
  - Proteomics project coordination & data dissemination
  - Integration of NWChem & NMR
  - Scientific Workflows

# Tools Development in EMSLHub

- HubZero framework allows programs to be developed, shared, used within as "tools" in EMSLHub

- Anyone with access can submit tools

- Approved tools can be shared with groups

- First examples are:
  - NMR analysis pipeline
  - ScalaBLAST job submission
  - IMPROV in EMSLHub

# Example: Metaproteomics pipeline implementation in EMSLHub



Download Sequence Data from JGI → Obtain mass spec proteomics data from EMSL → EMSLhub submits parallel jobs to compute resources → Perform Sequence Comparisons & Clustering → Clustered protein sequences viewed in IMPROV

# Example: Metaproteomics pipeline implementation in EMSLHub

Download Sequence Data from JGI → Obtain mass spec proteomics data from EMSL → EMSLhub submits parallel jobs to compute resources → Perform Sequence Comparisons & Clustering → Clustered protein sequences viewed in IMPROV
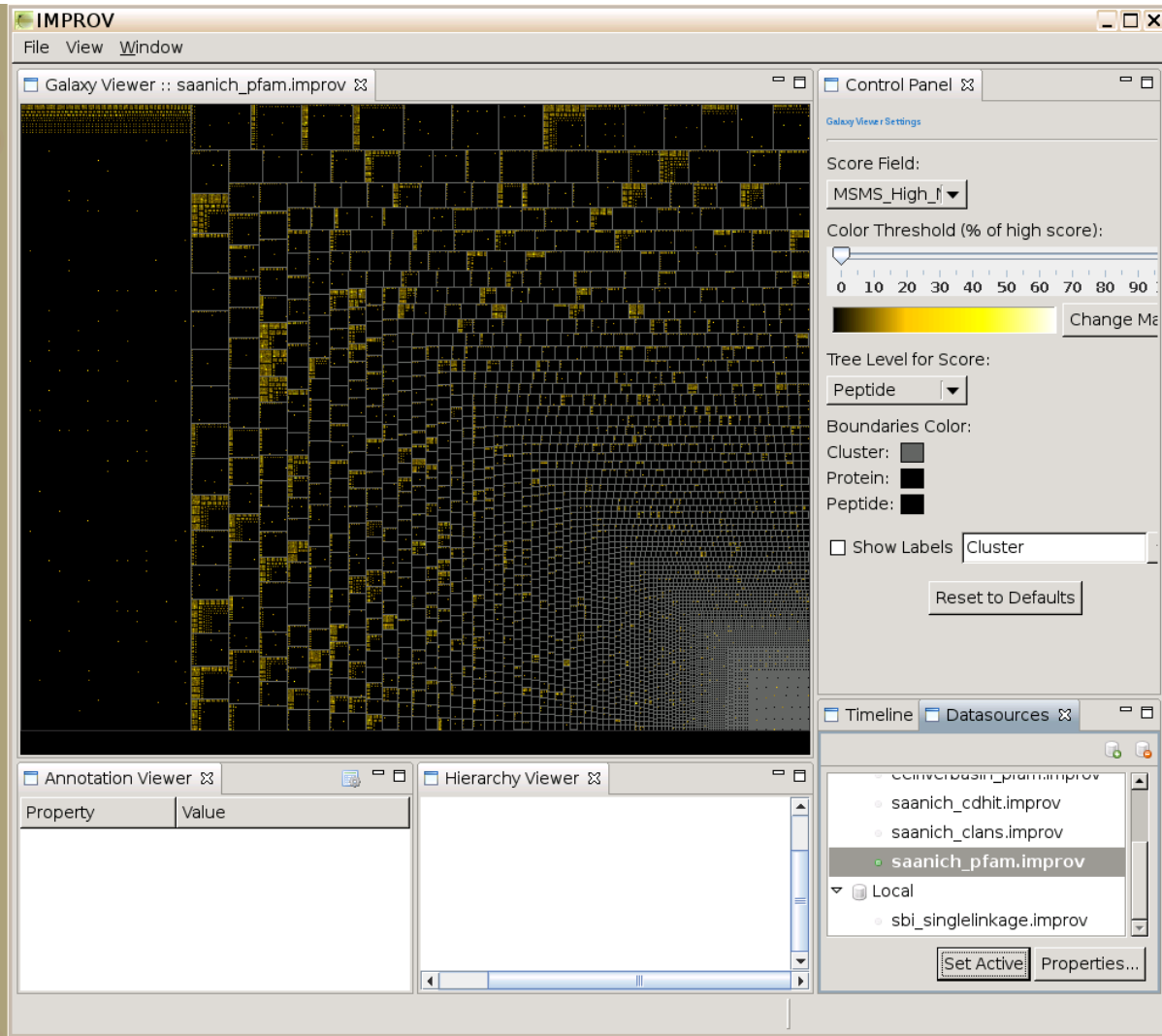
# Example: Metaproteomics pipeline implementation in EMSLHub

# Example: Metaproteomics pipeline implementation in EMSLHub



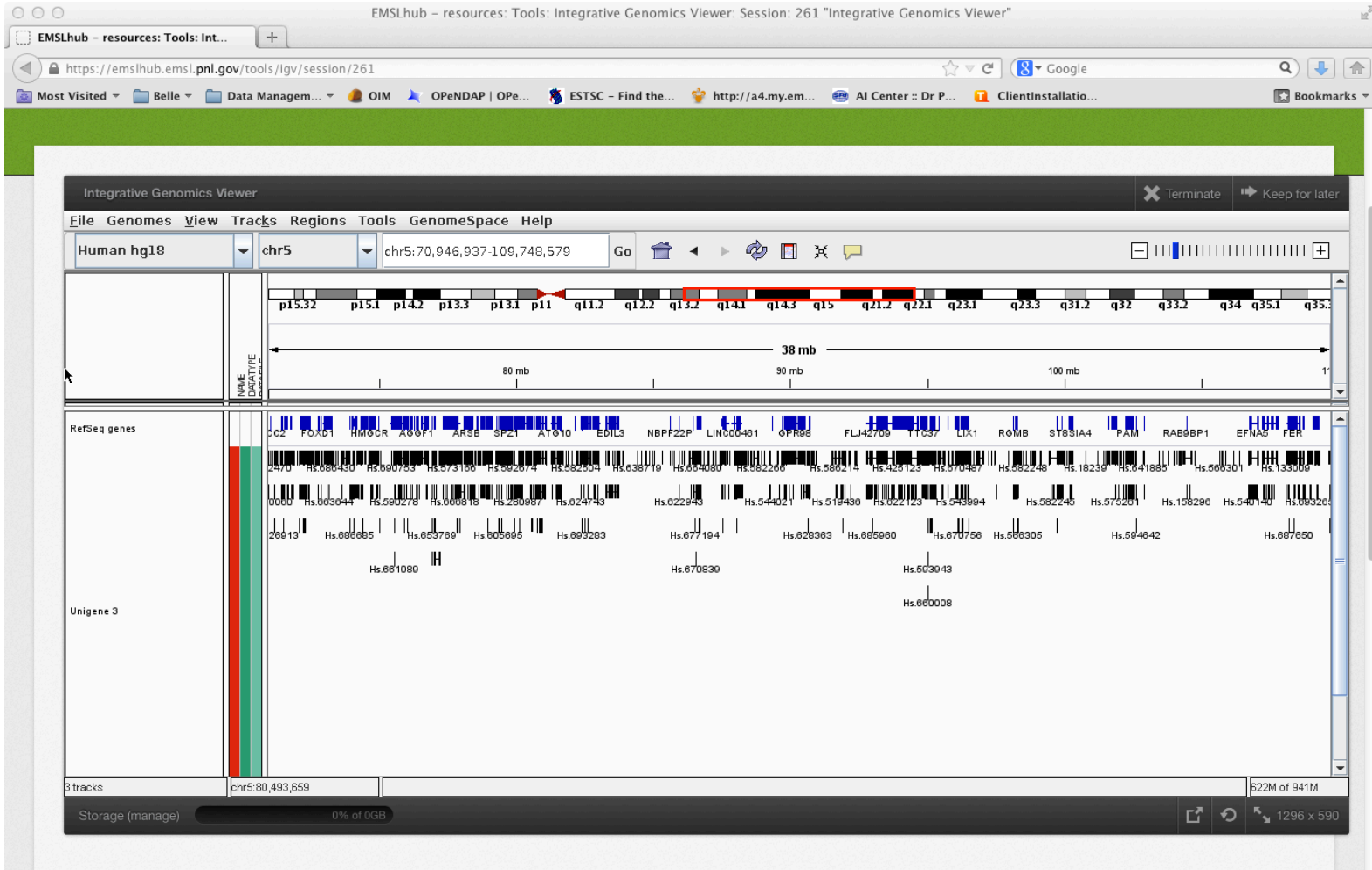Download Sequence Data from JGI → Obtain mass spec proteomics data from EMSL → EMSLhub submits parallel jobs to compute resources → Perform Sequence Comparisons & Clustering → Clustered protein sequences viewed in IMPROV

# Example: Metaproteomics pipeline implementation in EMSLHub

Pacific Northwest
NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY

Proudly Operated by *Battelle* Since 1965

# Supporting the need for new approaches to data & collaboration

- Data is being generated both experimentally and computationally at an very rapid pace

- Islands of data are growing at institutions all over the world

- The leading edge of scientific discovery in the 21st century will be attained by collaboration and combining data in new ways

- MyEMSL & EMSLHub will help us

  - Cultivate the exchange of ideas

  - Leverage shared, searchable data

  - Foster collaborations, sharing & discovery



Pacific Northwest
NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY

Proudly Operated by Battelle Since 1965

# Team Members

- Project Management
  - William Shelton
  - David Cowley
- Core Project Team
  - Kevin Glass
  - David Brown
  - Brock Erwin
  - Kevin Fox
  - Nate Trimble
- EMSLHub
  - David Brown

- Brock Erwin
- Silvia Hoisie
- Metagenomics partners
  - Angela Norbeck
  - Ken Auberry
  - David Brown
- NMR partners
  - Karl Mueller
  - David Brown
  - Herman Cho
  - Bert DeJong
  - Brock Erwin

- Nancy Washton
- Microscopy Partners
  - Nigel Browning
  - James Evans
  - James Bouwer

# Thank You!

www.**emsl**.pnl.gov

# Metadata Sample