*Environmental Proteomics Conference, Keystone, CO*
*January 2010*

# Approaches to Microbial Community Proteomics Data Analysis

**Angela D. Norbeck**
Senior Scientist
Pacific Northwest National Laboratory

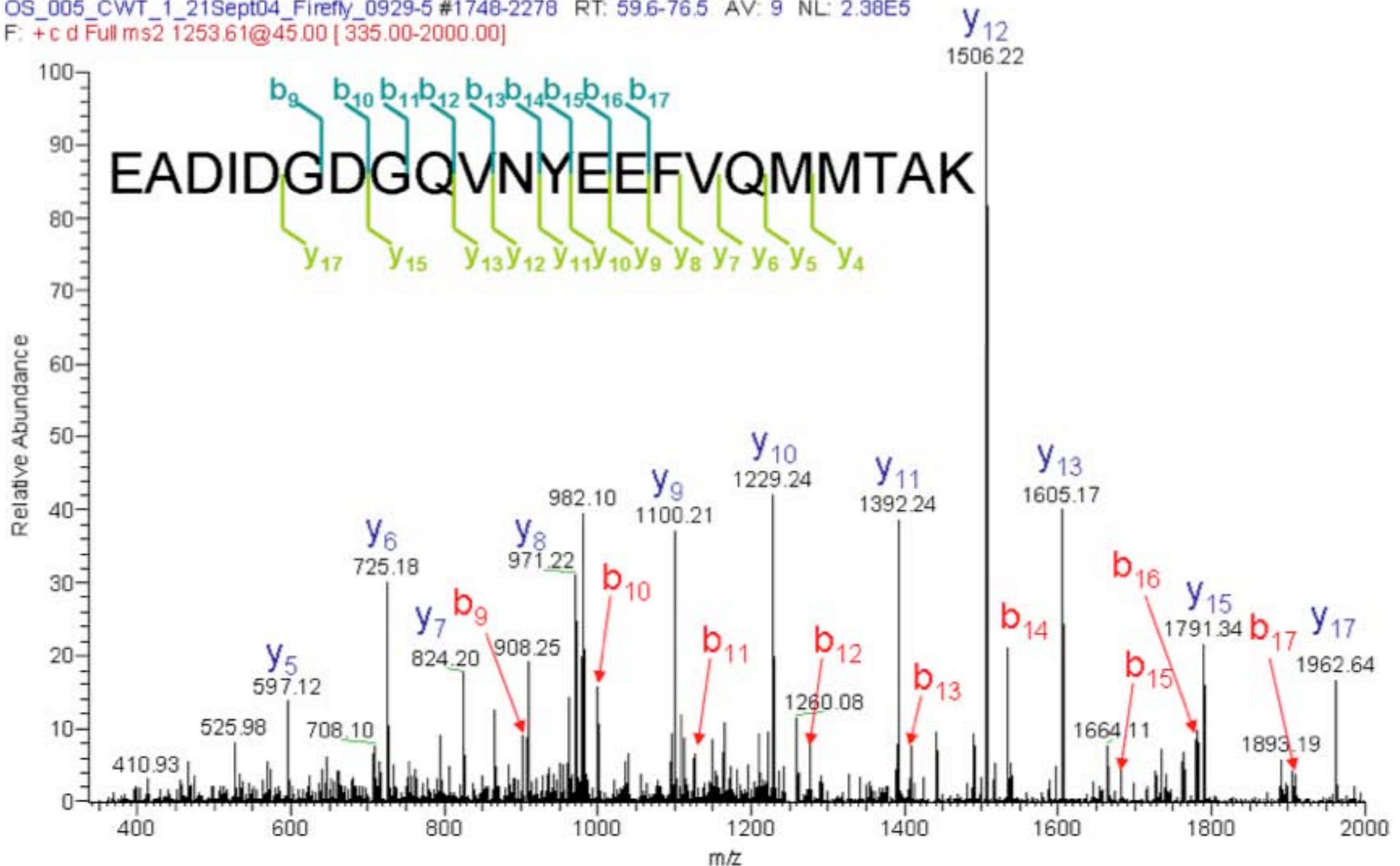**Pacific Northwest**
NATIONAL LABORATORY

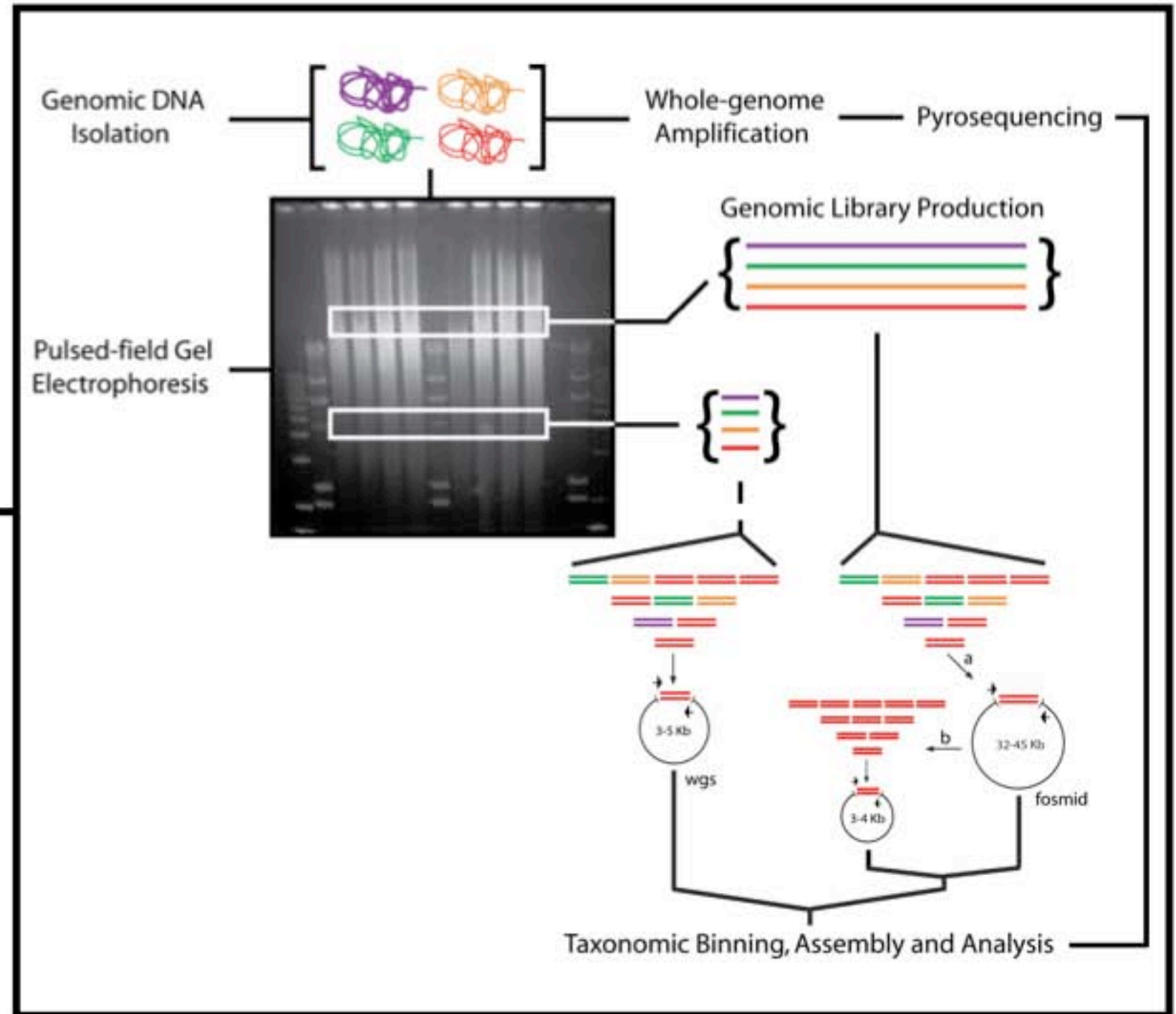*Proudly Operated by* Battelle *Since 1965*

# Challenges

- **Extracting protein from samples**
  - Sample prep method development

- **Protein file needed for searching spectra**
  - Sequences from metagenome sequencing
  - Groups of annotated organism files
  - *In-silico* derived sequences

- **Finding proteins of interest within large results set**
  - Data analysis methods

# Protein files: Matching fragmentation (MS/MS) spectra to protein sequence files



OS_005_CWT_1_21Sept04_Firefly_0929-5 #1748-2278   RT: 59.6-76.5   AV: 9   NL: 2.38E5
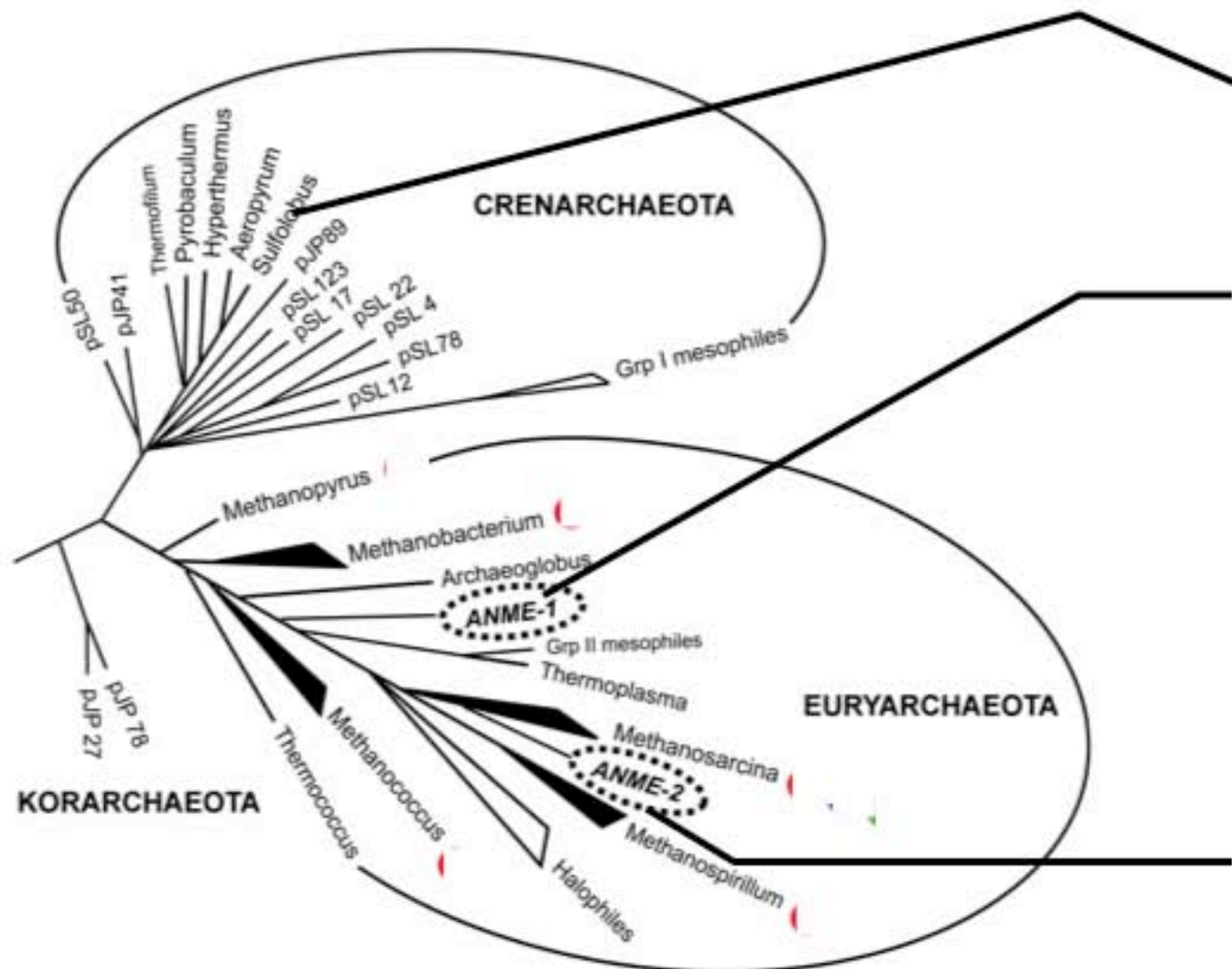F: + c d Full ms2 1253.61@45.00 [ 335.00-2000.00]

EADIDGDGQVNYEEFVQMMTAK

# Metagenomic sequences: library production



Slide from Steven Hallam, UBC

# Groups of annotated organism files

## Methane oxidizing Archaea



Figure from Hinrichs, K. et al. (1999) *Nature* 398:802-805

>gi_15605614_ref_NP_212987.1_ elongation factor Tu [Aquifex aeolicus VF5]

MAKEKFERTKEHVNVGTIGHVDHGKSTLTSAITCVL AAGLVEGGKAKCFKYEEIDKAPEEKERGITINIT

>gi_15605615_ref_NP_212988.1_ ribosomal protein S10 [Aquifex aeolicus VF5]

MEQEKIRIKLRAYDHRLLDQSVKQIIETVKRTGGVVK GPIPLPTRKRKWCVLRSPHKFDQSREHFEIREF

SRILDIIRFTPQTIEALMEISLPAGVDVEVKMRG

>gi_15605616_ref_NP_212989.1_ ribosomal protein L03 [Aquifex aeolicus VF5]

MPLGLIGEKVGMTRVLLKDGTAIPVTVIKFPVNYVVQ VKSQNTKDGYNALQIGAYEAKEKHLTKPLIGHF

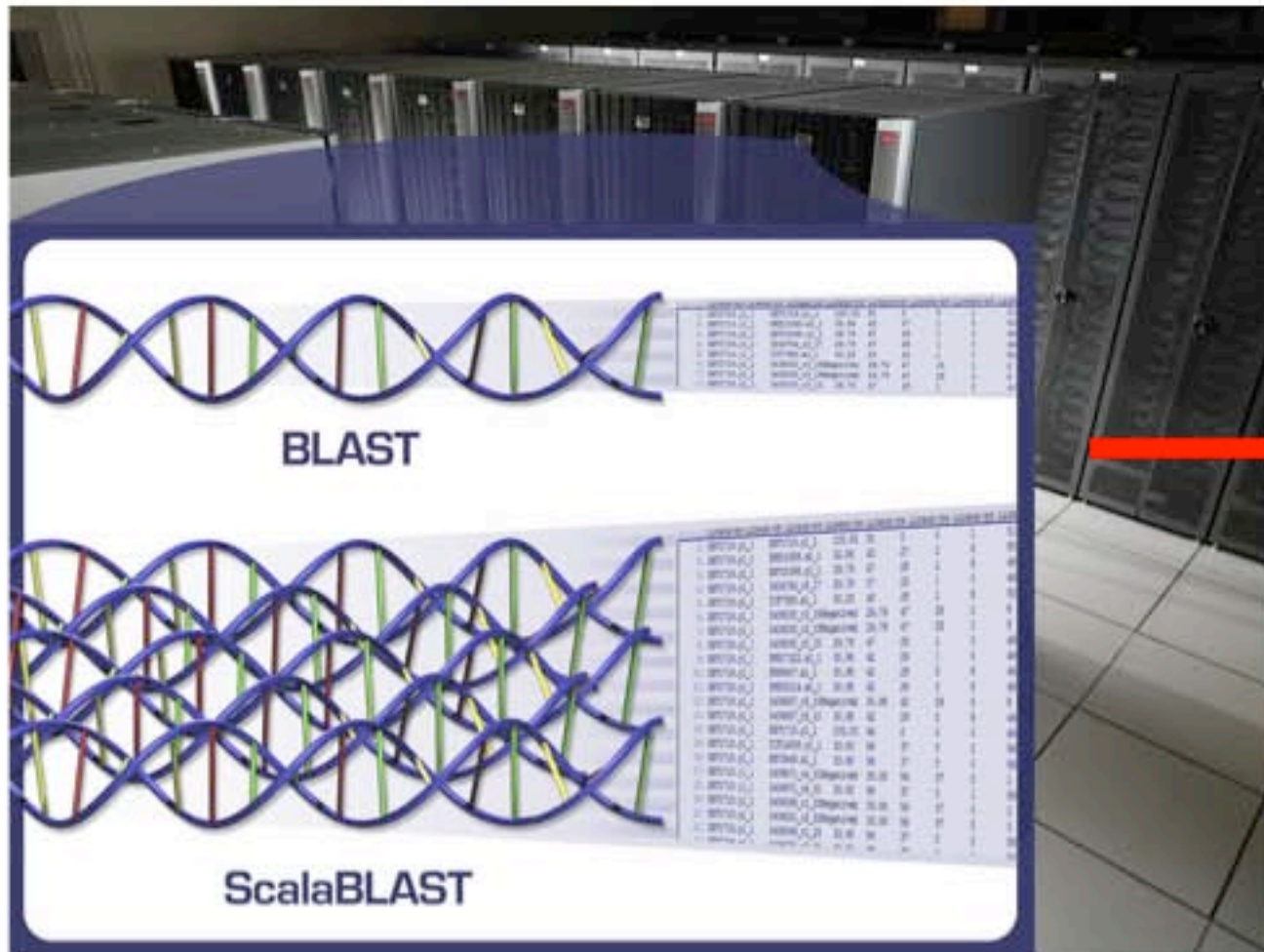K>gi_11498795_ref_NP_070024.1_ replication factor C large subunit [Archaeoglobus fulgidus DSM 4304]

DQRSWRVIERIVGEGAFNETISDEGEFLSSRIGKLK LIILDEVDNIHKKEDVGGEAALIRLIKRKPAQPL

I>gi_220904928_ref_YP_002480240.1_ NADH (or F420H2) dehydrogenase, subunit C [Desulfovibrio desulfuricans subsp. desulfuricans str. ATCC 27774]

MESLEIADRLRGFFPEEVLDVREFRGQLAVLVRSGR ILELLAYLRDVLDMRHLQALCGVDNSRRNEPGLS

QGHPLRKEYPVKIPARGHEEWEGLTALRKRAAELD ALSWQGGARHE

# *In-silico* derived sequences: generating protein lookup tables from multiple organisms having similar function

# Data analysis methods: software requirements are tough to meet

- **Globally view a community proteome using mass spectrometry and highlight identification regions**
  - Challenge: data files in GB with millions of rows
  - Data complexity: files are interrelated
  - No existing capability to detect patterns
  - Previous methods involved sorting lists and looking at top 1000 rows

- **Zoom in on regions and extract information about proteins**
  - Challenge: Redrawing of interface takes time and memory and information extraction about subset requires a new query

- **Change between different peptide and protein information**

# IMPROV – Integrated MetaPROteomics Viewer

# Galaxy View Navigation with Control Panel

# Example: Proteomics exploration of the ocean floor



1 — Dive T201 on the southern ridge of the Eel River Basin.

2 — Push core PC45 sampled below bacterial mat. Cells enriched from 6-9 cm interval associated with methane oxidation.

Slide from Steven Hallam, UBC

# Application to Eel River Basin microbial community



**Methanogenesis protein clusters found to be dominant**

# 172,000 proteins in 1700 clusters
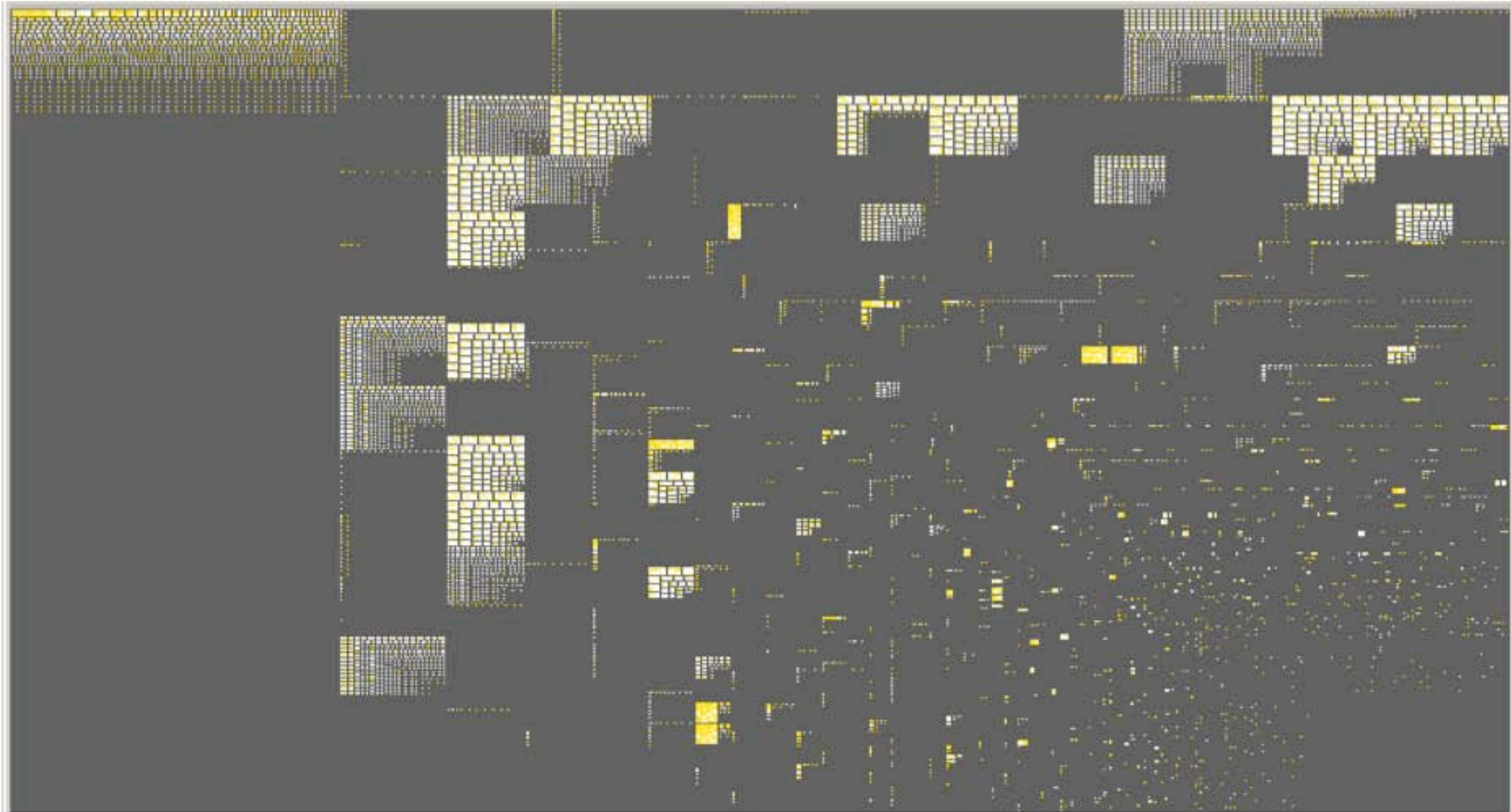
# New cellular functions discovered

# Unknown functions explored

# Example: Proteomics exploration of a Sargasso Sea microbial community

- **2008 Global Ocean Survey protein file downloaded from CAMERA**

- **Over $6 \times 10^6$ protein entries (2.2 GB file)**

- **Took longer than 1 month to search spectra**

- **Identified 16,000 clusters of proteins**

- **Peptide database is over 4 x "normal" (single organism) size**
  - Log file for database maintenance is an order of magnitude larger than usual

# Some proteins highly expressed, others not
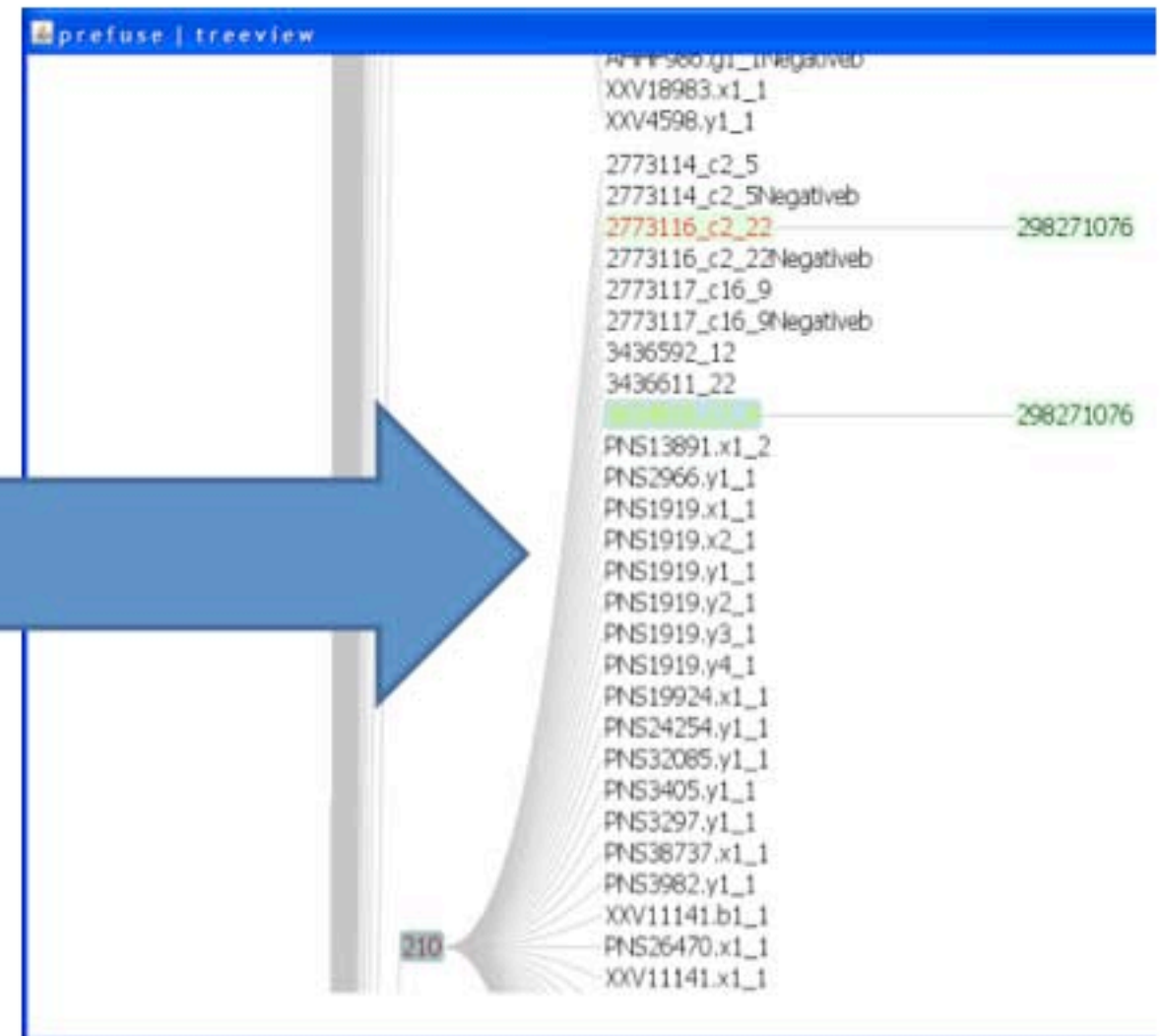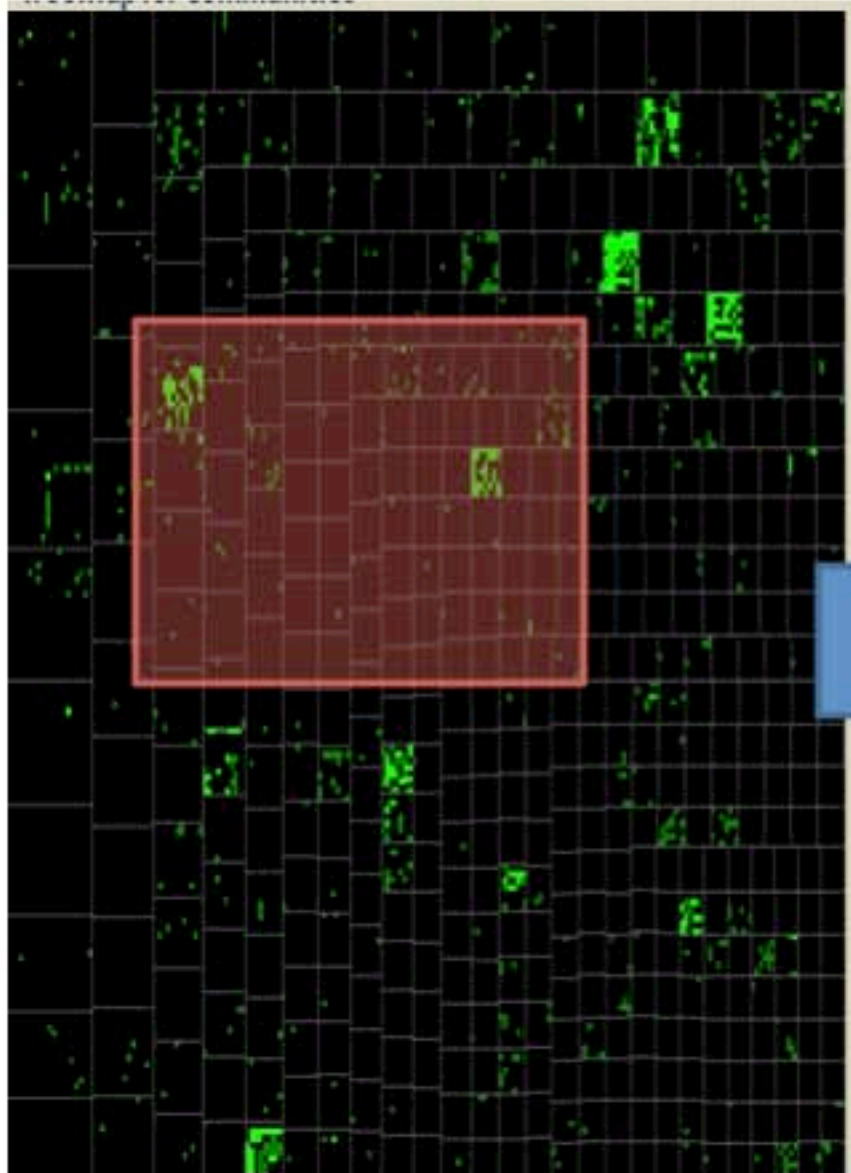
# Shows annotations

# Data mining of smaller, less abundantly expressed proteins is critical to understanding dynamics

# Additional IMPROV features

- **Multiple ways to cluster proteins (i.e. by function, cell location, pathway..)**

- **Ability to toggle between peptide/protein spectra counts and peptide/protein abundance values or modification states**

- **Text string searching**

- **Zooming functions**

- **Ability to view and analyze multiple conditions and find clusters changing dynamically**
  - Spiral graphs for time course

- **Normalization of clusters by percent coverage**
  - Reorganization of map by highest coverage

- **Link out to other databases (Kegg, pfam..)**

# Acknowledgments

**Software Development**
- Mudita Singhal
- Kelly Domico
- Getiria Onsongo

**Informatics Team at PNNL**



**Funding and Facilities**



**Collaborators**
- Dr. Steven Hallam
  University of British Columbia
- Dr. Stephen Giovannoni
  Oregon State University

**High-Throughput Proteomics Production Team at PNNL**

# Discussion Points

- **What are the main kinds of biological applications for environmental proteomics?**

- **Should data repositories adhere to standard formats?**

- **How much data can or should be shared?**

- **How should orthologs amongst different organisms be used to relate functions in a community?**