**Summary of the 2008 BER S&O Review**

– **New JGI Management and Organizational Structure**

– **Changes in JGI Informatics**

**2009 JGI Five Year Strategic plan**

**Dec 3-5, 2008, Report Issued March 3, 2009**

- **Science**
- **Management**
- **Operations**
- **Informatics**

# Recent JGI Publication Metrics

**2009**

**Total Peer-Reviewed Publications (Science/Nature/PNAS)**   **81 (18)**

**2009 citations of JGI-Authored Papers published 2005-Present**   **18,919**

*Sorgum Genome*   *Nature* **2009**

**Two algal (Micromonas) Genome**   *Science 2009*

**Manuscripts in various stages of review**

*Genomic Encyclopedia of Bacteria and Archea*   *(under review Nature)*
*Brachypodium Genome*   *(submitted to Science)*
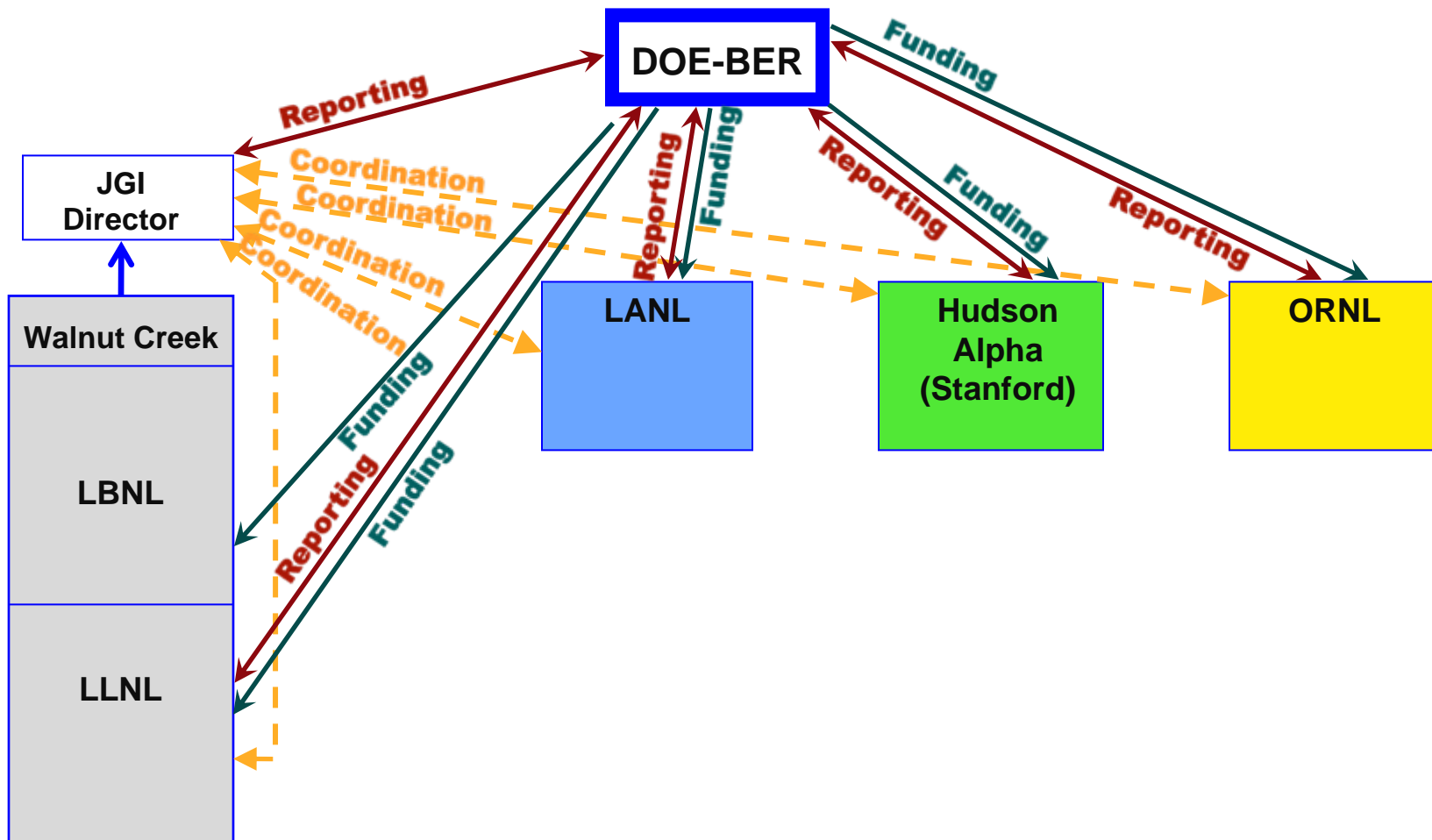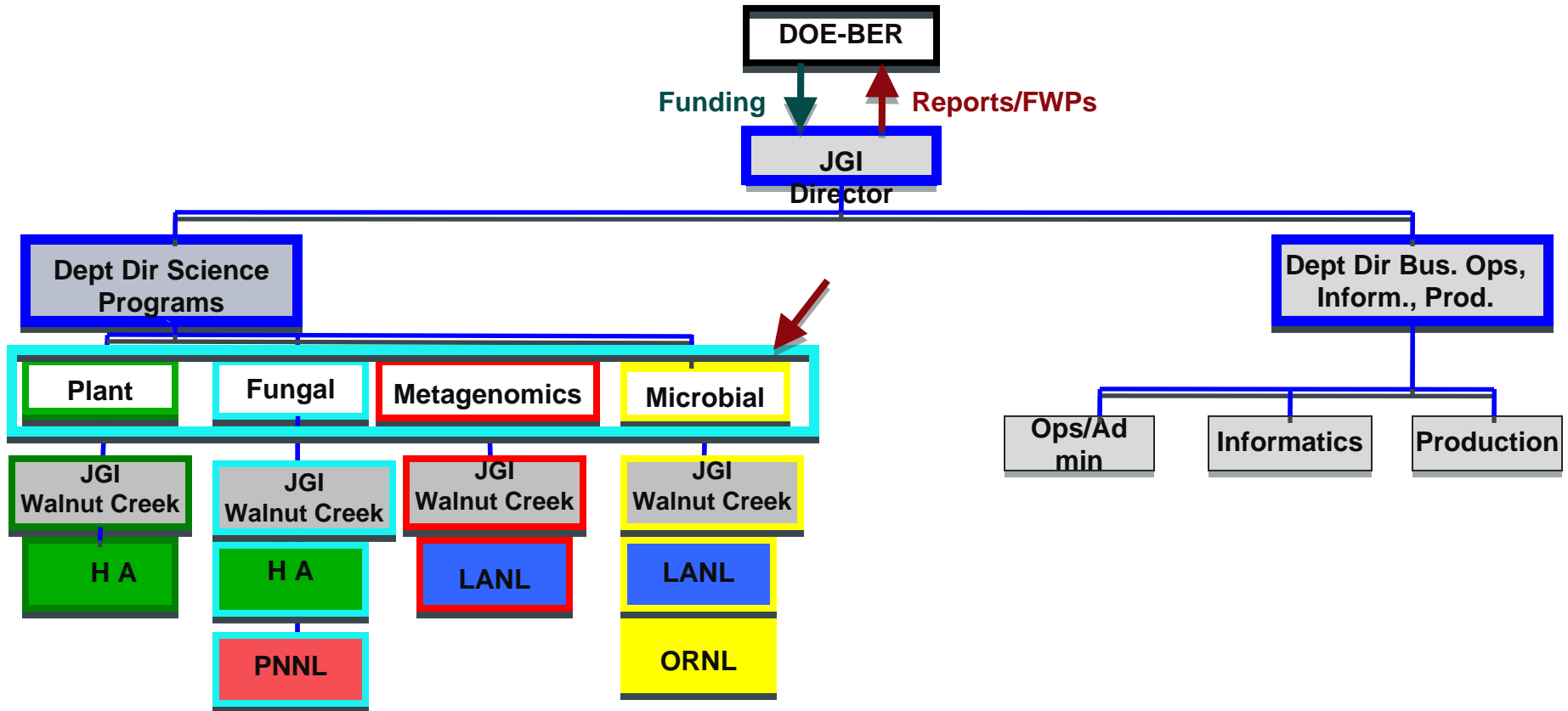*Soybean Genome*   *(soon to be submitted)*

**Dec 3-5, 2008, Report Issued March 3, 2009**

- Science
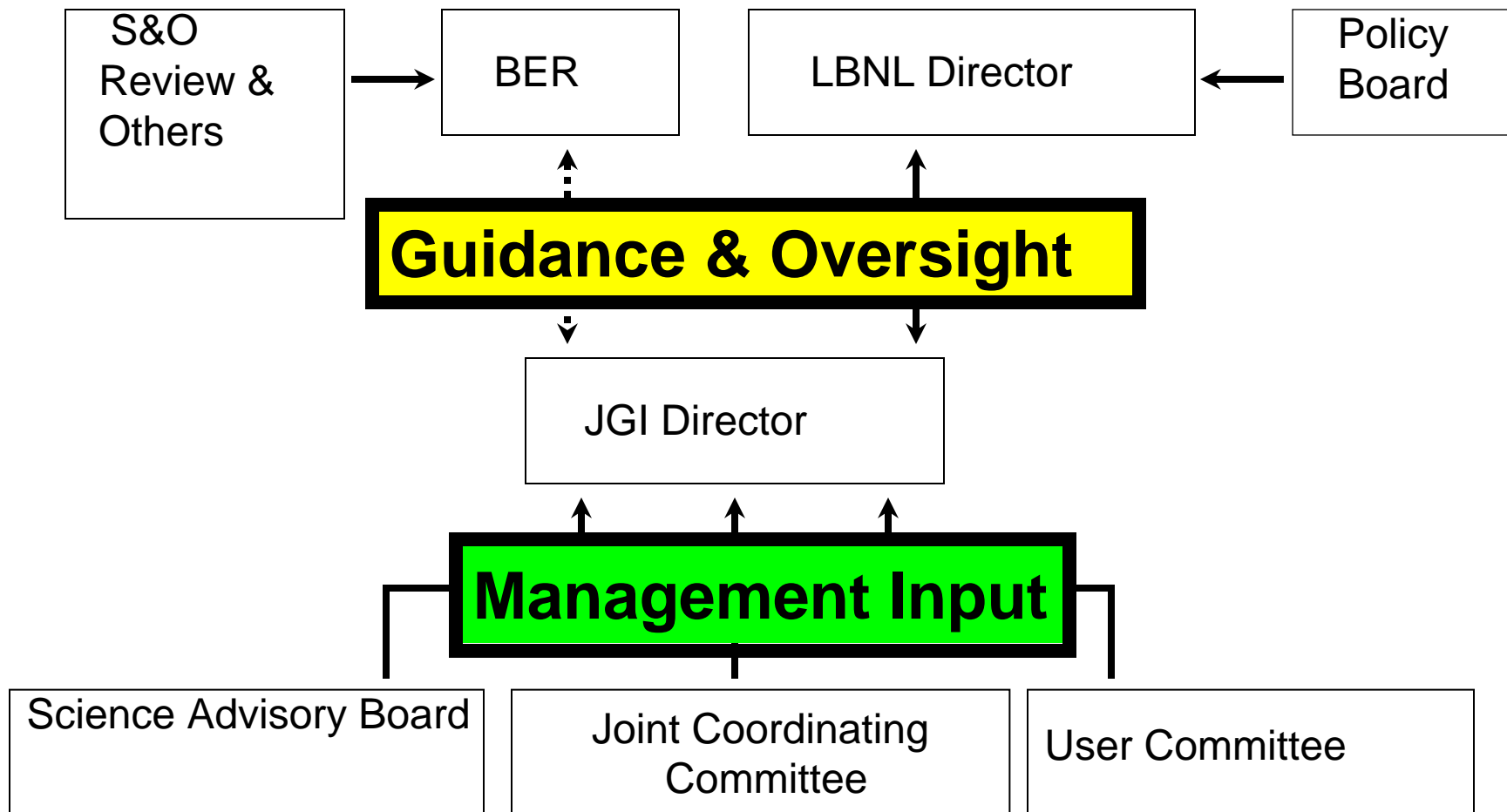- **Management**
- **Operations**
- Informatics

- **"Establish a centralized management and organizational structure"**

- **"Establish funding flow from LBNL to partners to align authority and responsibility"**

# JGI Operations and Management

- **Coordination of JGI Scientific Goals Unlinked to Funding and Reporting**

# JGI Operations and Management



- JGI Director Responsible and Accountable for JGI Productivity
- Organization Broken into Programs
- **Programs Coordinate the JGI Scientific Goals Linked to Funding**

**JGI Science Program Management Model Implemented**

**Science Program Leads Named and Functioning**

**Implemented Program Based Financial Planning & Reporting**

*(Science based programmatic funding decisions have already been made)*

**Dec 3-5, 2008, Report Issued March 3, 2009**

- Science
- Management
- Operations
- **Informatics**

**"Remedy Senior Leadership Deficiencies"**

**Concern about the Informatic Planning and
Hardware Infrastructure Faced with
Dramatic Increases in Data Generation**

## JGI Informatics Department Head
### Svilen Tzonev (*Illumina / Solexa , Sr. Director of Software Development and Bioinformatics*)

*JGI Chief Informatic Officer/ JGI Associate Director (New position)*
*Jill Mesirov  Part Time Consultant to the JGI Director (CIO / Associate Director Broad Harvard MIT))*

JGI Informatics Project Management Program Head
(New position)
Evi Dube  (Previously Computing Division Leader (LLNL))

**JGI Informatics Department Head**
**Svilen Tzonev (*Illumina / Solexa , Sr. Director of Software Development and Bioinformatics*)**

*JGI Chief Informatic Officer/ JGI Associate Director*
*(New position)*
*Jill Mesirov  Part Time Consultant to the JGI*
*(CIO / Associate Director Broad Harvard MIT)*

**JGI Informatics Project Management Program Head**
*(New position)*
**Evi Dube  (Previously Computing Division Leader (LLNL))**

**JGI Informatics Department Head**
Svilen Tzonev (*Illumina / Solexa , Sr. Director of Software Development and Bioinformatics*)

*JGI Chief Informatic Officer/ JGI Associate Director (New position)*
*Jill Mesirov  Part Time Consultant to the JGI Director (CIO / Associate Director Broad Harvard MIT))*

**JGI Informatics Project Manager
*(New position)*
Evi Dube  (Previously Computing Division Leader (LLNL))**

# Informatics

## Informatics Advisory Committee
### *(New advisory group)*

**First meeting scheduled Dec 09**

# I. ARRA Funding Targeted for IT Upgrades:  $11.1M

# Funding Received-August 7th

Add Infrastructure for Management, Storage and Movement of Genome-Scale Datasets

*Increase Storage (100% increase)*
*Upgrade JGI's Central Computing Cluster (200% increase)*
*Upgrade ESNet Access*

**II.** Explore the DOE National Lab's high performance computing capabilities to help JGI's data intense challenges

### Projects already **underway** or being developed

**Using ScalaBLAST for Metagenomics**

– Porting ScalaBLAST to JGI clusters; help parallelize short read assembler

**Working with JGI on designing computing nodes**

– Scientific engagement with JGI on computational challenges

– IMG update using large BLAST runs focusing on isolate genomes and potentially some of the new Titanium metagenome datasets

**Two large-scale metagenomics problems identified**

- Goal 1: evaluate metagenome assembly by using simulated dataset
- Goal 2 : evaluate gene calling through a large blastx search against NR using joined Illumina pair-end reads.

– Argonne team identified to assist JGI team

– Assisted IMG in solving contentious parallel computing bug

– Working with JGI on designing computing nodes

– Possible Site for Remote Computing

# Informatics

**Summary:**
**Recruiting leadership to strategically as well as tactically chart the JGI's informatic directions.**

**Additional informatic personnel and hardware.**

**Engaging the DOE National Labs to adapt their high performance computing capabilities to address JGI needs**

- **23 Committee Recommendations**
  - Informatics (10 Recommendations)
  - Operations & Safety (7 Recommendations)
  - Science (4 Recommendations)
  - Management (2 Management)

- **JGI Corrective Action Status**
  - All Recommendations are either fully implemented or are in the process of be implemented

- **Corrective Action Milestones-Reviewed by JGI Senior Managers and Communicated with BER**

# Revolutionary Technological Changes

**Mission:**

User Facility for Large Scale Genomics and Analysis
to Enable Bioenergy and Environmental  Research


**Unique Capabilities**

In carrying out genomics of focus, scale, and complexity
to help users solve important but hard DOE relevant problems

# Plants, Microbes, Metagenomes User Facility

ACHIEVEMENTS OF THE
NATIONAL PLANT GENOME INITIATIVE
AND NEW HORIZONS IN
PLANT BIOLOGY

NATIONAL RESEARCH COUNCIL
OF THE NATIONAL ACADEMIES

"It is critical that JGI continue to serve a broad remit for sequencing and resequencing of plant genomes…"

# Plant Genomes

| | |
|---|---|
| **Populus trichocarpa** | **black cottonwood** |
| **Chlamydomonas reinhardtii** | **chlorophyte alga** |
| **Physcomitrella patens** | **bryophyte moss** |
| **Sorghum bicolor** | **sorghum** |
| **Glycine max** | **soybean** |
| *Panicum virgatum* | *switchgrass* |
| *Miscanthus x giganteus* | *miscanthus* |
| *Volvox carteri* | *multicellular green alga* |
| *Arabidopsis lyrata* | *Lyrate rockcress* |
| *Selaginella moellendorffii* | *lycophyte fern* |
| *Brachypodium distachyon* | *false purple brome* |
| *Mimulus guttatus* | *monkeyflower* |
| *Manihot esculenta* | *cassava* |
| *Prunus persica* | *peach* |
| *Setaria italica* | *foxtail millet* |
| *Capsella rubella* | *Pink Shepherd's purse* |
| *Aquilegia coerulea* | *Goldsmith columbine* |
| *Gossypium raimondii* | *diploid cotton* |
| *Eucalyptus grandis* | *rose gum* |
| *Phaseolus vulgaris* | *Common bean* |
| *Zea mays (Mo17 inbred)* | *maize* |

# Flagship Plant Genomes
## (High Priority Organisms)

- **Produce higher quality accurate and well annotated genomes**
- **Add additional information**



*Sorghum Bicolor*

**Foxtail millet**

*Chlamydomonas*

**Soybean**

*Brachypodium*

**Poplar**

**Switchgrass**

*Miscanthus*

# Problem with plant genomes:
# Though we have the sequence we cannot assign function to most of it

# Discovery of functional elements through comparative genomics



Vol 450|8 November 2007|doi:10.1038/nature06340

nature

ARTICLES

## Discovery of functional elements in 12 *Drosophila* genomes using evolutionary signatures

Alexander Stark[1,2]*, Michael F. Lin[1,2]*, Pouya Kheradpour[2]*, Jakob S. Pedersen[3,4]*, Leopold Parts[5,6], Joseph W. Carlson[7], Madeline A. Crosby[8], Matthew D. Rasmussen[2], Sushmita Roy[9], Ameya N. Deoras[2], J. Graham Ruby[10,11], Julius Brennecke[12], Harvard FlyBase curators†, Berkeley Drosophila Genome Project†, Emily Hodges[12], Angie S. Hinrichs[4], Anat Caspi[13], Benedict Paten[4,5,14], Seung-Won Park[15], Mira V. Han[16], Morgan L. Maeder[17], Benjamin J. Polansky[17], Bryanne E. Robson[17], Stein Aerts[18,19], Jacques van Helden[20], Bassem Hassan[18,19], Donald G. Gilbert[21], Deborah A. Eastman[17], Michael Rice[22], Michael Weir[23], Matthew W. Hahn[16], Yongkyu Park[15], Colin N. Dewey[24], Lior Pachter[25,26], W. James Kent[4], David Haussler[4], Eric C. Lai[27], David P. Bartel[10,11], Gregory J. Hannon[12], Thomas C. Kaufman[21], Michael B. Eisen[28,29], Andrew G. Clark[30], Douglas Smith[31], Susan E. Celniker[7], William M. Gelbart[8,32] & Manolis Kellis[1,2]

# JGI Pilot Project to Sequence and Analyze 12 Arabidopsis Relatives
## (Project Leaders: Joe Ecker, Detlef Weigel, Dan Rokhsar)



The 12 Arabidopsis Relatives Project will:

• Identify functional elements in these genomes based on evolutionary constraint

• Provide a road map for similar projects in plant groups with much larger genomes
  e.g., grasses

**Produce the genomic tools necessary to understand cell wall construction, feedstock production and carbon sequestration and sunlight energy harvesting pathways**

# JGI Microbial Program:

**Bacteria and Archaea 939 Complete Genomes** (NCBI, August 2009)

JGI
31%

Other
47%

TIGR
10%

Genoscope
3%

JCVI
4%

Sanger
5%

**JGI has produced ~30% of the world's complete archaeal and bacterial genomes**

bacterial and archaeal phyla with cultured isolate
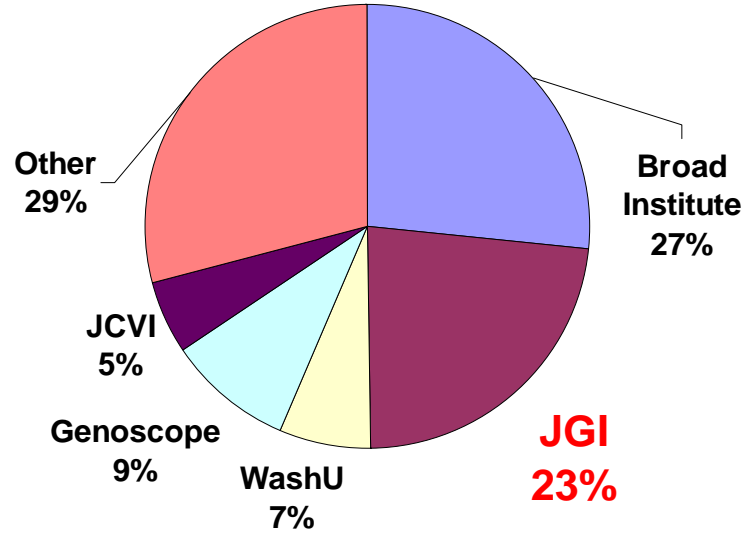
bacterial and archaeal phyla with cultured isolate

# Pilot Genomic Encyclopedia of Bacteria and Archaea
## GEBA (Sequencing and analysis of 64 genomes)

H. Utahensis

New Sequenced Phyla GEBA

Halobacteria
Methanomicrobia
Thermoplasmata
Methanobacteria
Thermococci
Chloroflexi
OP10
Th
Gemmatimonadetes
Fibrobacteres
Elusimicrobia
Planctomycetes

**(Final stages of Nature review)**

▪ **"Significantly improves analysis of metagenomic data"**

▪ **"Markedly improves the identification of new gene families and members of gene families"**

## *including those of DOE relevance*

Bacteroidetes
OP5
Spirochetes
Cyanobacteria
Deferribacte
Fusobacteria
Firmicutes

*H. Utahensis*

**Great Salt Lake**

*Halorhabdus utahensis*

27% NaCl

6x saltier than ocean

A cellulase gene cluster

# Optimized Cellulases

Ionic Liquid Detergents for Cellulose Dissolution



before          after

DOE Bioenegy Center (JBEI)
of Synthesized H Utahensis
Cellulases

High Salt
Cellulase Activity

bacterial and archaeal phyla with **cultured isolate**

The true extent of the coverage….
bacterial and archaeal diversity including major uncultured phyla

# Our ultimate goal….

**bacterial and archaeal diversity including major uncultured phyla**

**Technologies for accessing the genomes of hard to culture organisms**

# JGI Fungal Program

Fungal genome projects worldwide

# *Vol.1.* **Plant feedstock health**

*Part 1.* Mycorrhizal Symbionts (Basidiomycota)
*Part 2.* Plant Pathogens (Dothideomycota)

# *Vol.2.* **Lignocellulose degradation**

White and brown rot (Basidiomycota)

# *Vol.3.* **Fermentation**

# *Vol.4.* **Industrial organisms**

*Part 1. Trichoderma spp.*

*Part 2. Aspegillus nigri*

**Genomic Encyclopedia of Fungi (GEF)**

THE NEW SCIENCE OF
**METAGENOMICS**

Revealing the Secrets of Our Microbial Planet

NATIONAL RESEARCH COUNCIL
OF THE NATIONAL ACADEMIES

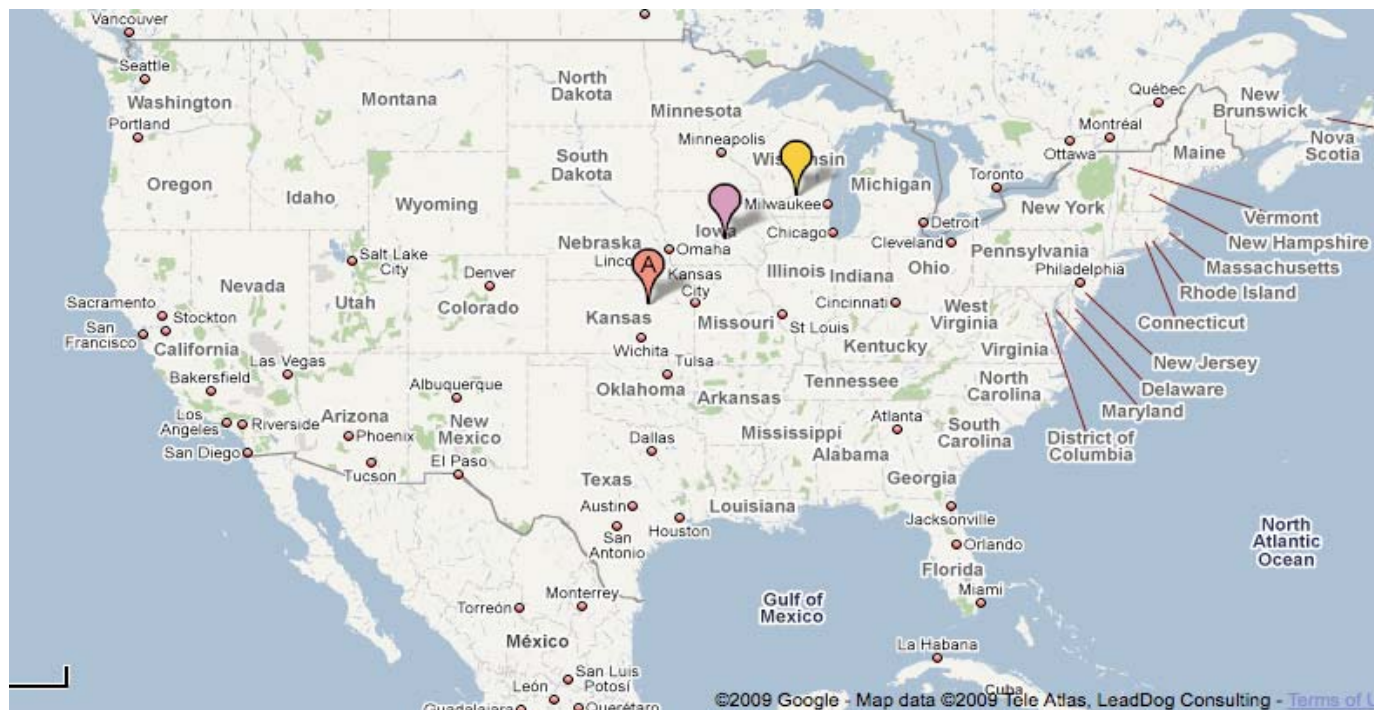recommended the establishment of a small number of large-scale projects to study particular habitats

# Soil

Represents a habitat that is important due to its role in carbon and nitrogen cycling and biomass production. It is a particularly "genomicaly" challenging environment due to its complexity.

## Multiple samplings of 3 Different Prairie Sites

- **Midwest prairie represents largest expanse of world's most fertile soils**
- **Sequesters the most carbon of any soil system in the U.S.**
- **Produces large amounts of biomass annually: Key for biofuels, carbon sequestration**

Iowa never-tilled prairie

Iowa >100 yr tilled (corn)

- **To improve soil management, carbon sequestration,**

- **Through genetic understanding maybe able to potentially manage traits such as green house gas fluxes, and carbon stability**

Established Switchgrass (Wisconsin)

62 | **Termite hindgut**, 62 Mbp

100,000 **Great Prairie metagenomic project** , ~100 Gbp

1,000,000 **Potential Future Projects:**
Terabase Projects

*Terror* terabase projects

- **Sequence the 12 Great Soil Orders (the major soil types)**
  - Covers the extremes in soil diversity, e.g. permafrost to tropical desert to alkali flats to taiga forest to wetlands to discover gene adaptation to environmental extremes.

*Sequence the State Soil of each State*

*These official state soils share the same distinction as state flowers and birds.*

# JGI User Programs

# Opportunities and Challenges

Pre Incubation

Post Incubation

**Transmission Electron Microscopy**
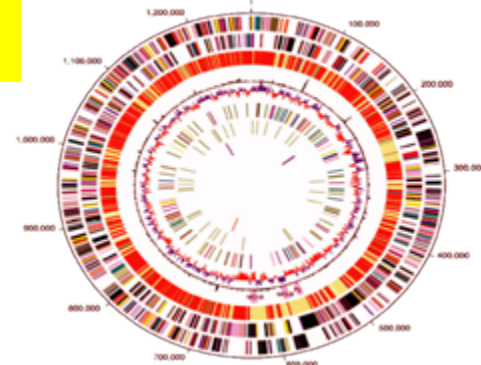
**(Termite project total size: ~62 MB)**

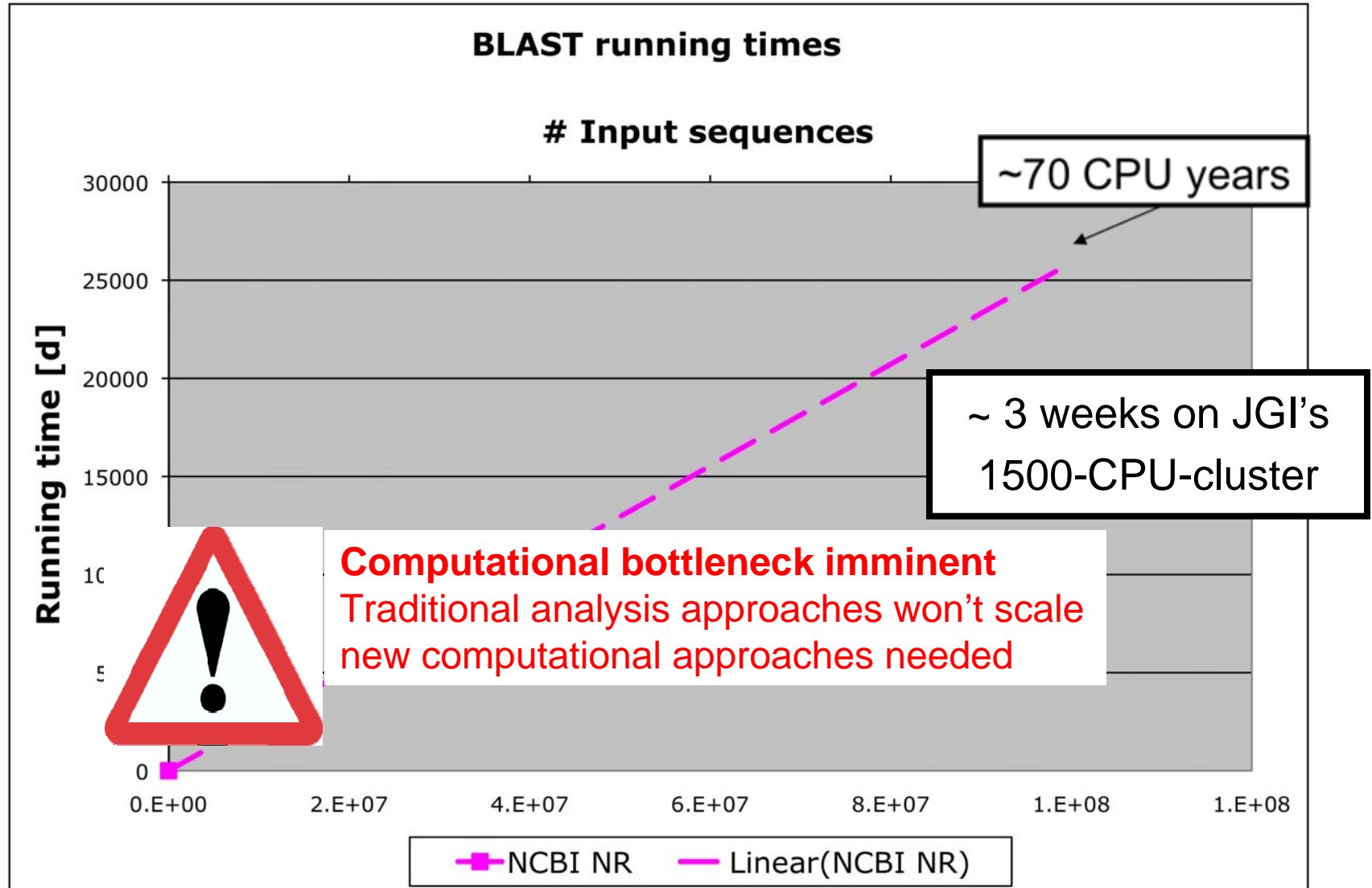**Rumen project total size: ~18 GB**

**Full Length Lignocellulolytic Genes**

## JGI-Sponsored Workshop on HPC for Next Gen Sequencing

## Focus on bringing National Lab computational infrastructure and expertise to bear on problems associated with the volume and nature of next gen sequencer out-put

# It is not about the machines