

# **DOE Joint Genome Institute**

## **Computing Requirements**

**Shane Canon, Rob Egan, David Goodstein**

**Victor M Markowitz**

**NERSC BER Requirements for 2017**

**September 11-12, 2012**



**Established in 1999, located in Walnut Creek**  
**Supported by the DOE Office of Science.**  
**Budget of ~70 Million/ yr ~ 300 Employees**

## Mission

**User facility for large scale genomics to enable bioenergy & environmental research**

## Mission Areas



**Bioenergy**



**Carbon Cycling**



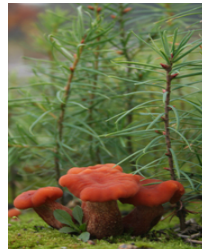
**Biogeochemistry**

## JGI Programs

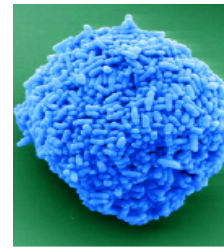
**Plants**



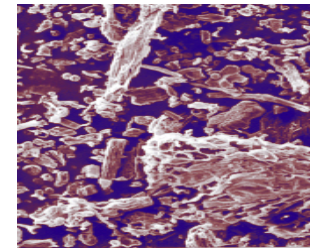
**Fungi**



**Microbes**



**Metagenomes**





## Energy and Environmental Science

**Experimental  
Data Generation**

**Sequencing**

**Data  
Interpretation**

**Assembly of  
Sequence Data**

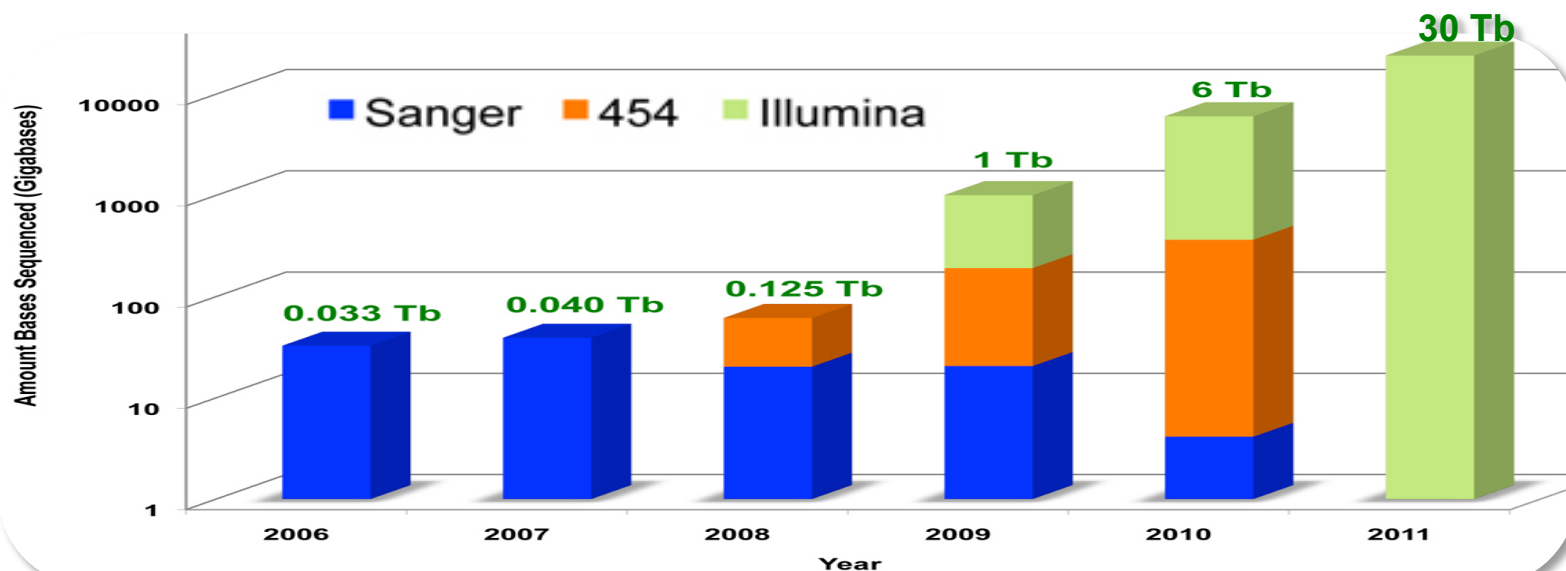
**User  
Interactions**

**Data Distribution**

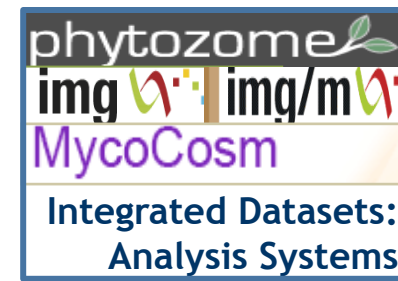
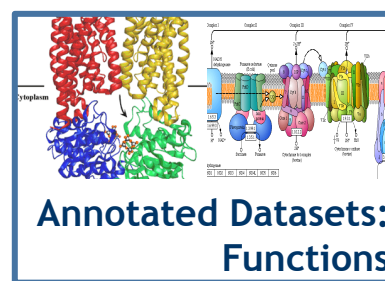
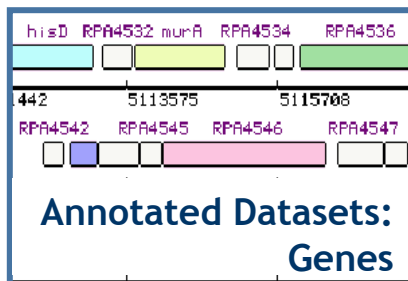
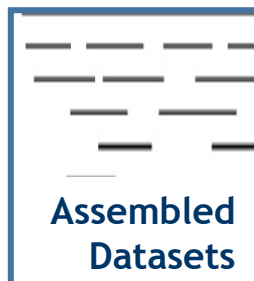
# JGI drives scientific discovery



Sequence  
data



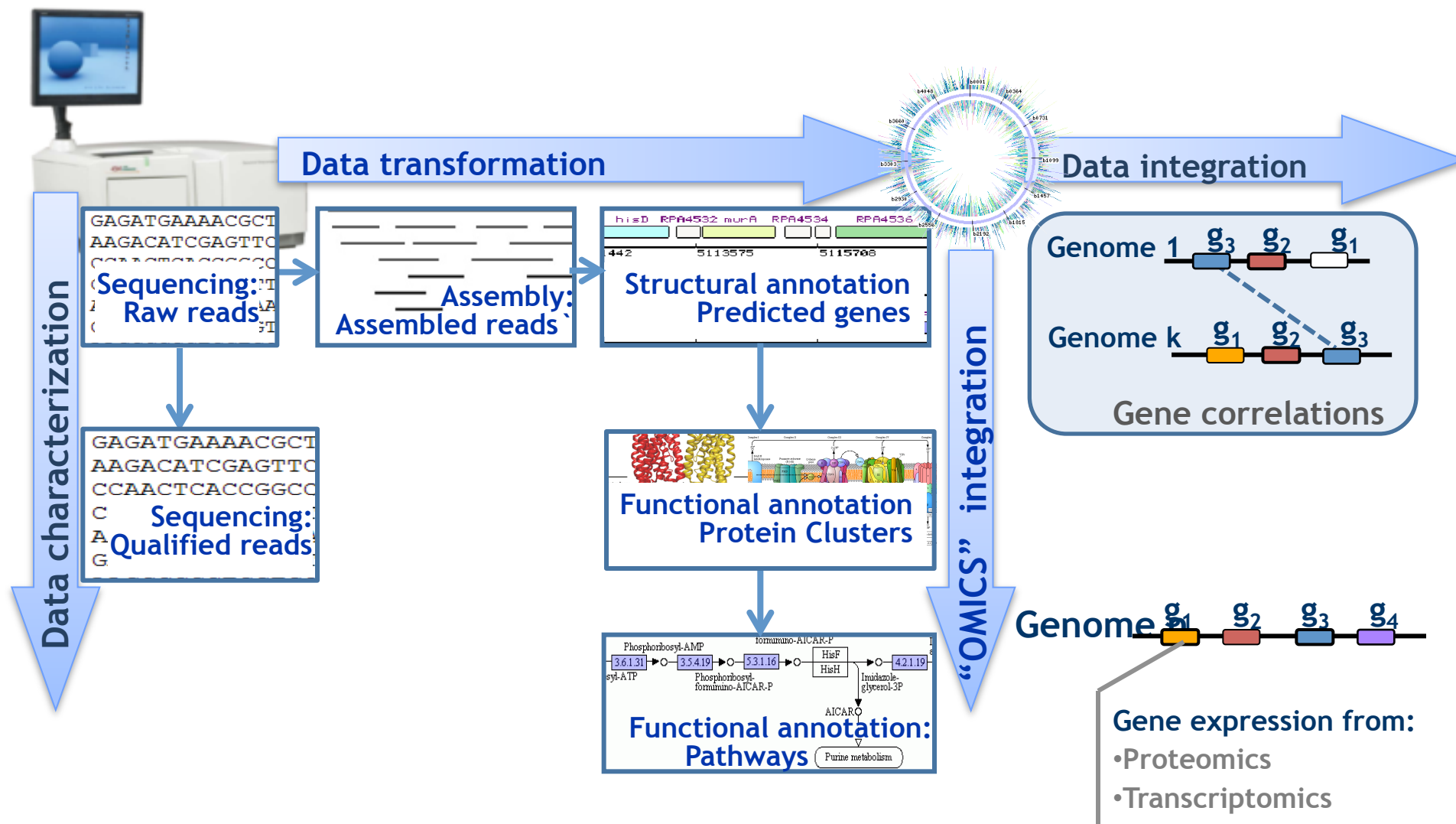
From raw to  
interpreted  
data



Science  
publications

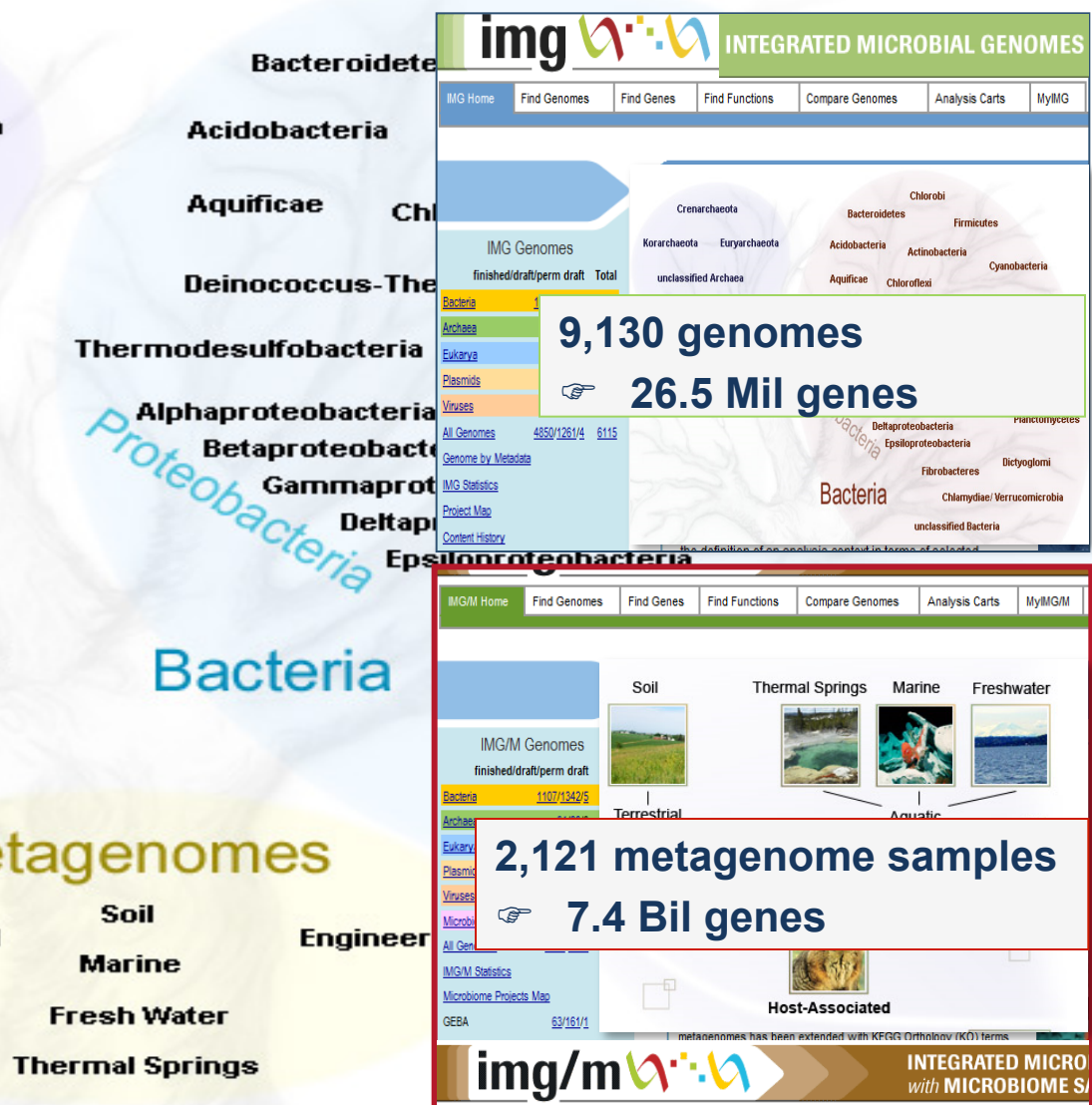
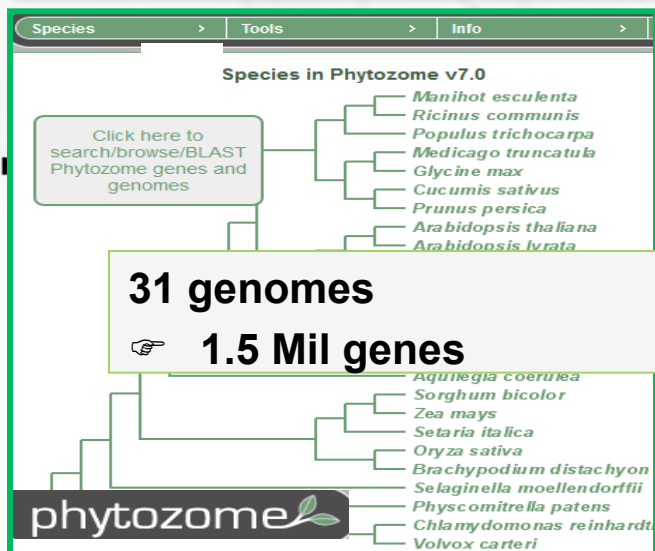
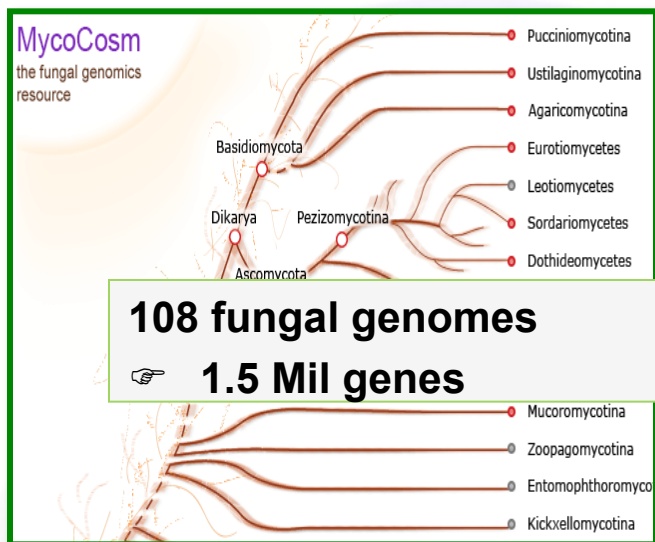






Iterative process of genomic data interpretation & refinement

## Data integrated framework for (meta) genome sequence data analysis





## Energy and Environmental Science

### Experimental Data Generation

**Sequencing**

**Sample Prep**

**Single Cell Gen**

**DNA Synthesis**

**Whole Genome  
Fxn Annotation**

### Data Interpretation

**Assembly of  
Sequence Data**

**Integrated Data  
Management**

**Functional  
Discovery and  
Annotation**

**High Performance /  
High Throughput  
Computing**

### User Interactions

**Data Distribution**

**Onsite User  
Interactions**

**User Training**

**Community  
Coordination**

**Data Analysis  
Platforms**



ESNet Bay Area  
100 Gbits/sec



CONNECTION  
TO CHICAGO

2010

2011

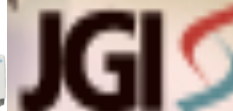
2012

2013

...

2017

Start JGI system  
deployment  
@ NERSC



Supercomputing Infrastructure

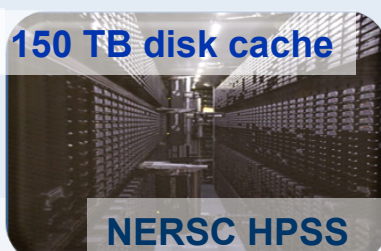
**NERSC**

153,216 cores

150 TB disk cache



NERSC Supercomputer



NERSC HPSS

Core Computing Infrastructure

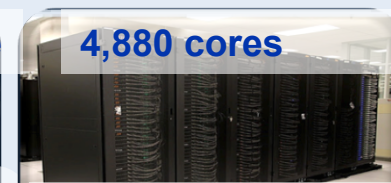
Web & DB Servers

2 PB storage

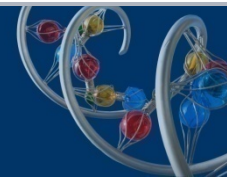
File System

4,880 cores

Compute Cluster



# Scaling challenge



## ❖ Computing needs

**Existing** sequence data processing methods, pipelines, systems **cannot keep up** with rapid increase in size & number of sequence datasets

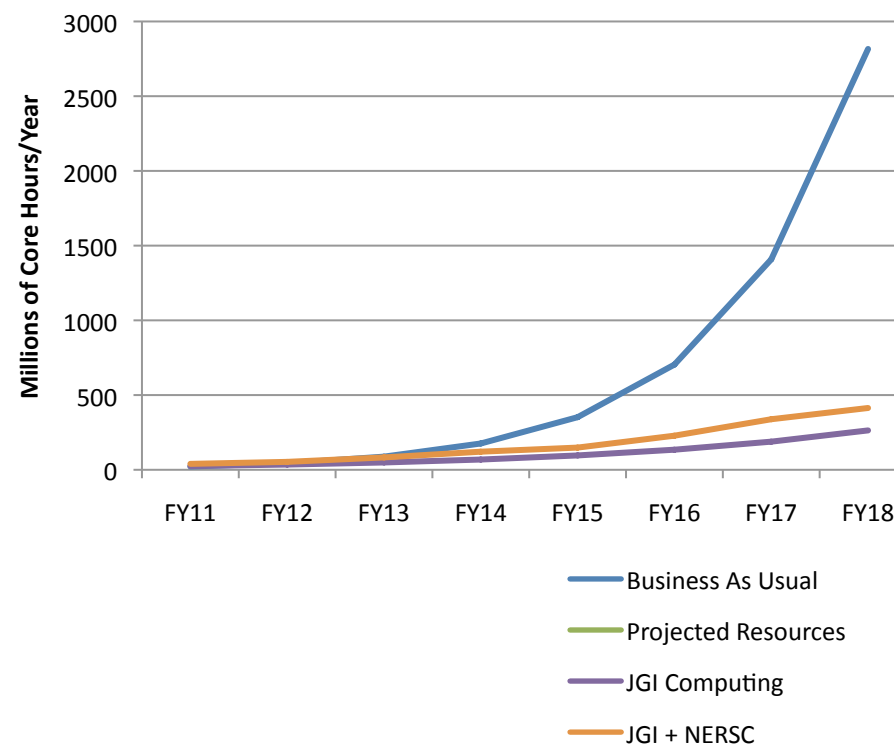
## ❖ Alternatives

### ☐ Increase supply

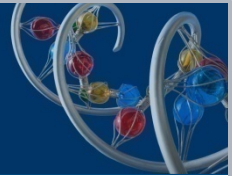
- ✓ ■ Buy more
- ✓ ■ Use other resources
- Use the cloud

### ☐ Reduce consumption

- Improve efficiency
- Reduce data processing







## ❖ Optimize

Re-engineer key methods/tools/workflows

## ❖ Innovate

Design new analysis methods

## ❖ Prioritize

Selective computations

## ❖ Limit

Sequence what you can analyze

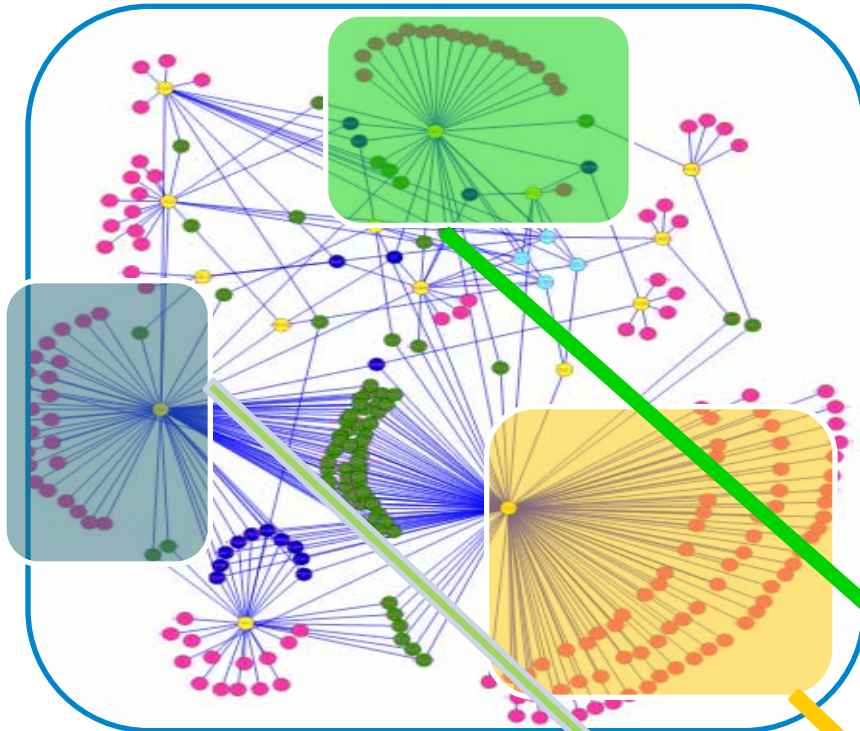
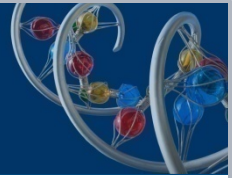
## ❖ Adapt

Ask different questions

Re-engineered workflows addressed scaling challenge for microbial genome & metagenome data processing

Advanced data management (Fastbit) techniques applied to gene cassette conservation (operon) analysis results in substantial performance boost





## IMG/M: 7 Billions Genes

Metagenome analysis limited to known protein families from sequenced organisms

☞ ignores about 70% of genes of

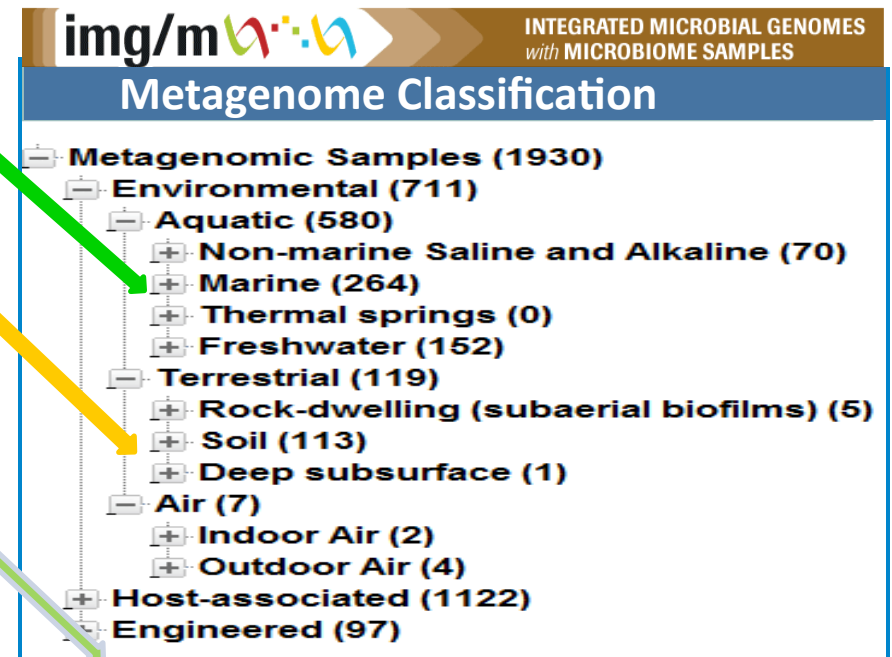
potentially new organisms that are not similar to these protein families

## Protein clustering

Allow comparing metagenome samples in terms of known & potentially new protein families

☞ requires 206 M CPU core hours annually for maintaining clusters using linear gene pledging

☞ alternative strategies that can provide 10X acceleration are being explored



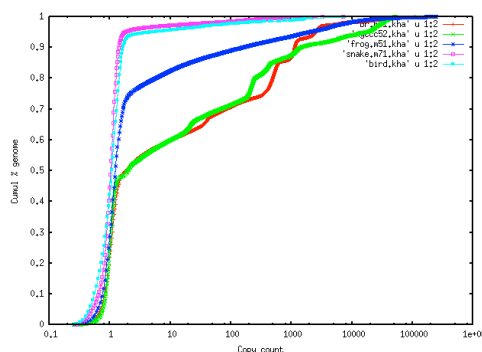


(~700Mb plant genome, with Sanger reference)

**Meraculous (Chapman):** scales up to 1000+ cores (full cluster); peak (bottleneck) memory ~10-20Gb. 2-3 days elapsed time. 486Mbp scaff. total; 74.4kb scaffN50; 7.9kb contigN50

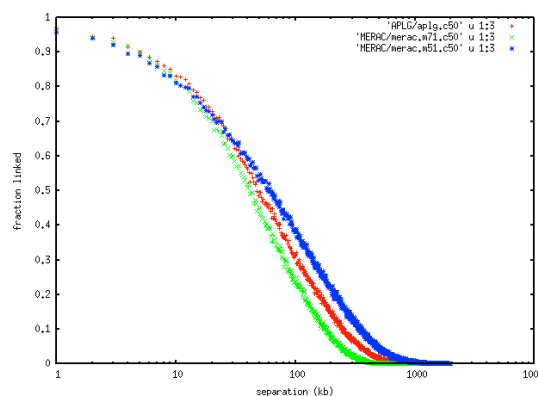
**APLG (Labutti):** 15 days elapsed time; peak memory 365Gb; not “clusterable”. 438Mbp scaff. total; 69kb scaffN50; 19.4kb contigN50

**SGA (Sunkara):** 11 days elapsed; 30Gb peak memory; (in principle) parts can be run on cluster. 425Mbp scaff total; 7.9kbp scaffN50; 6.4kbp contigN50



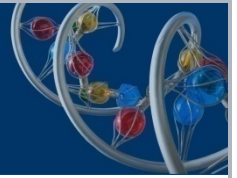
Plants are harder than (typical) vertebrates to assemble due to repeat content and polymorphism

**Correctness matters! (N50 is often our only metric, but it is a pretty poor one)**



“C50” (only measurable with some known reference sequence)  
APLG 48kbp;  
meraculous 62kbp.

2-4X more global misjoins in APLG result



- ❖ **Machines currently using: Carver and Hopper**
- ❖ **Hours used in 2012: ~12M core hours**
  - ❑ Primarily on Hopper using the taskfarmer
  - ❑ Key applications included BLAST, hmmsearch, and usearch
  - ❑ Scales 1k-12k (some larger runs too)
  - ❑ Native MPI applications included the JGI developed K-mernator and the ANL developed Kiki
- ❖ **Reading ~16 GB and ~TB of output**
- ❖ **Try to use around 16 GB per node (partition references to fit in memory)**
- ❖ **Heavy Use of TaskFarmer**
- ❖ **Typically stick to scratch for performance**



Software	HPC Scope	State	Speed on 1TBase, Reduced Metagenome
ABYSS	Contigs Only	Mature & Active	3h / 1800 cores (Reduced 787Gbase)
Ray	Scaffolds	Mature & Active	Unsuccessful
Kiki	Contigs Only	In Development	Unsuccessful
Pasha	"Pre-Graph" Only	Mature & Active	Not Attempted
Forge	Scaffolds	Not Mature, Inactive	Not Able to Compile
Contrail (Hadoop)	Scaffolds	In Development	Unsuccessful
Kmernator	Data Reduction & Contig Extension	Mature & Active	3h / 7200 cores (Full 1.115 TB raw data)



- ❖ **JGI funded resources will address steady state needs**
- ❖ **HPC allocation will be used to address other needs**
  - ☐ Problems beyond Genepool's scale
  - ☐ Development and testing new methods
  - ☐ Accelerating Time Sensitive Operations
  - ☐ Handling Demand during Peak Times
  - ☐ MPI-enabled applications benefiting from interconnect
  - ☐ Non-throughput workloads (parallel assemblers)



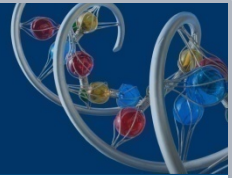
- ❖ **Core hours should keep pace with Sequence Growth -> ~75M core hours in 2017**
- ❖ **Changes in sequence technology will likely change workload characteristics**
  - ❑ NGS created an assembly challenge
  - ❑ Single molecule have longer reads with more errors
- ❖ **Increased support for data intensive workloads**



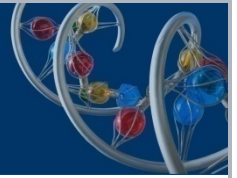


Software	Targeted Application	Speedup*	Comments
GPU HMMER (2.0)	Homology Search	20x	HMMER 3.0 is superior
GPU-BLAST	Homology Search	3-5x	Protein only
CUDASW++	Sequence Alignment	10x	
MUMerGPU	Sequence Alignment	3-10x	
SOAP3	Sequence Alignment	10x	
CuSHAW	Sequence Alignment	8-10x	
GSNP	Variant Detection	7x	
DecGPU	Read Correction	N/A	MPI based

**\* Speedup is not the appropriate metric; most calculations are against the (sometimes poorly optimized) single threaded code or a single core. Power consumption and capital costs should be rigorously evaluated against a large data set.**



- ❖ **Support for non-relational database services (NoSQL, Key-Value stores)**
- ❖ **Continue and enhanced support for Hadoop**
- ❖ **Enhanced support for High-Throughput Workloads (enable more workloads to move from Genepool to the big systems)**



- ❖ **JGI's partnership with NERSC is critical to meeting its mission requirements**
- ❖ **JGI will expand the use of its HPC allocation with a focus on new analysis, large-scale problems, and development**
- ❖ **JGI is interested in new services and capabilities around data and high-throughput computing**



# Thank You