- **DOE mission**: predict, control and design the biological components of energetic processes and environmental balance.

- Complex missions with rapidly expanding, intricately related diverse data types requires a mean to augment scientists' ability to:
  - Filter information
  - Focus attention
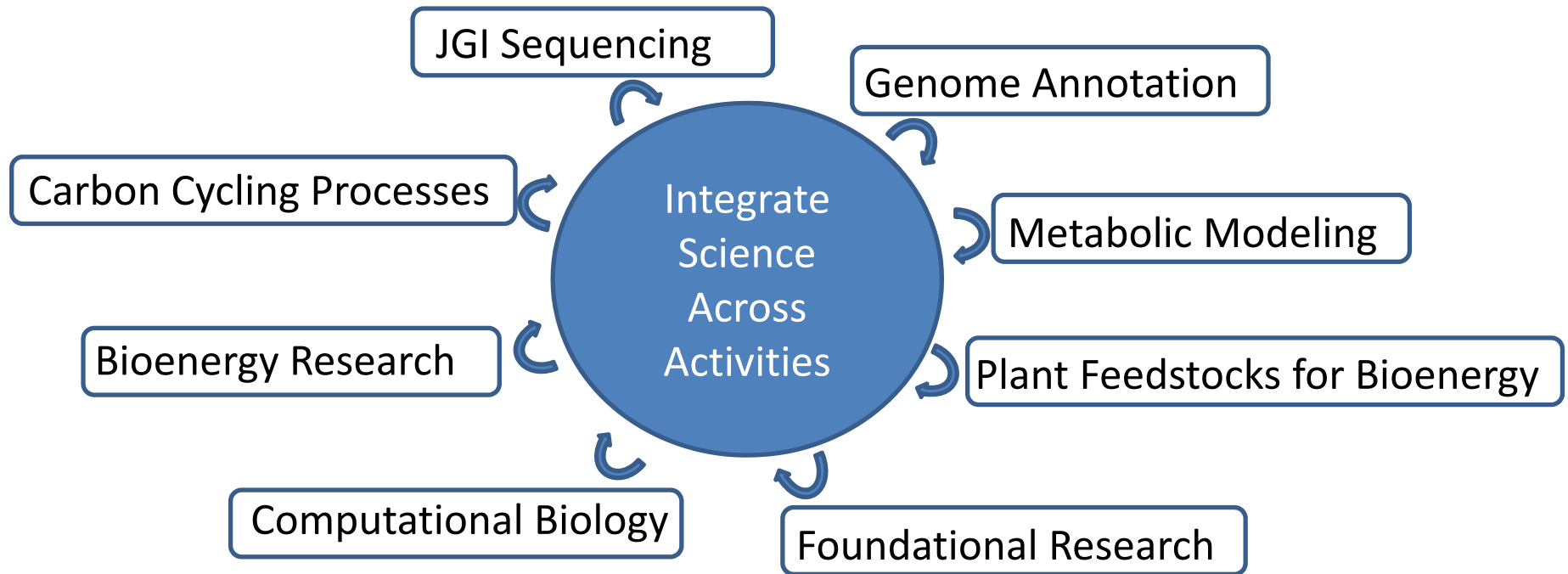  - Ask the right questions
  - Leverage other minds

JGI Sequencing

Genome Annotation

Carbon Cycling Processes

**Integrate Science Across Activities**

Metabolic Modeling

Bioenergy Research

Plant Feedstocks for Bioenergy

Computational Biology

Foundational Research
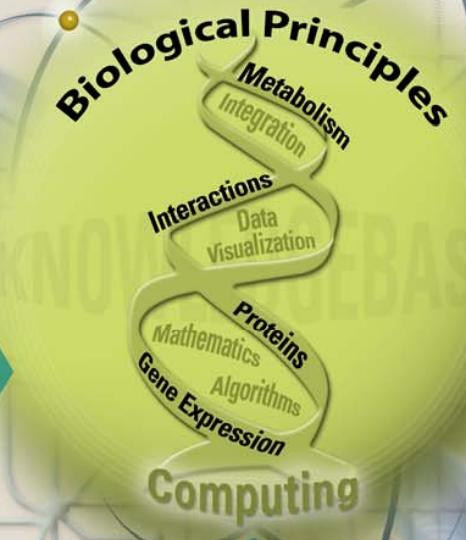
There is a tremendous wealth of data and information in the Genomic Sciences program. The Knowledgebase (Kbase) is an opportunity to integrate this data and information both within individual activities as well as to integrate together different activities.

# DOE Systems Biology Knowledgebase
## Components Enabling Its Development

Plants and Microbes for Energy and Environment

## Biological Principles

Metabolism
Integration
Interactions
Data Visualization
Proteins
Mathematics
Gene Expression
Algorithms

**Computing**

### 2010 Knowledgebase R&D Project

Year-long effort funded by American Recovery and Reinvestment ACT (ARRA).

Results from this project, which was completed in September 2010, are underpinning Knowledgebase development:

- **DOE Systems Biology Knowledgebase Implementation Plan**
- **ARRA pilot projects**

### User Community Data and Resources

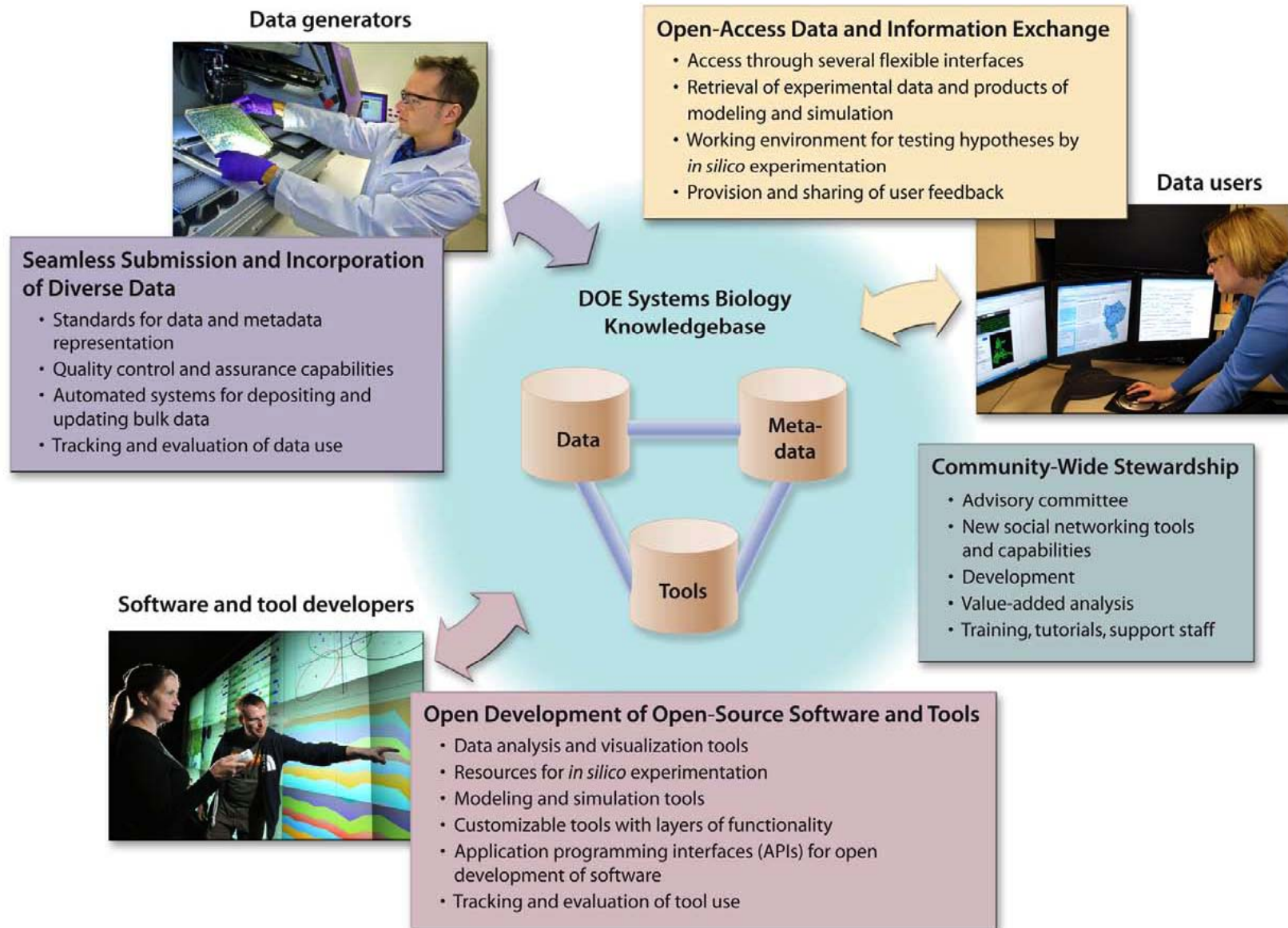The Knowledgebase will leverage and establish critical partnerships with other synergistic efforts:

- **DOE Joint Genome Institute**
- **National Center for Biotechnology Information**
- **iPlant Collaborative**
- **DOE Office of Advanced Scientific Computing Research**
- **Others**

### University-Led Projects to Develop Computational Biology and Bioinformatic Methods Enabling the Knowledgebase

In 2010, the DOE Office of Biological and Environmental Research awarded funding to 11 projects for basic research that will support Knowledgebase development in four areas:

- **Omic data integration**
- **Integrated pathway reconstruction**
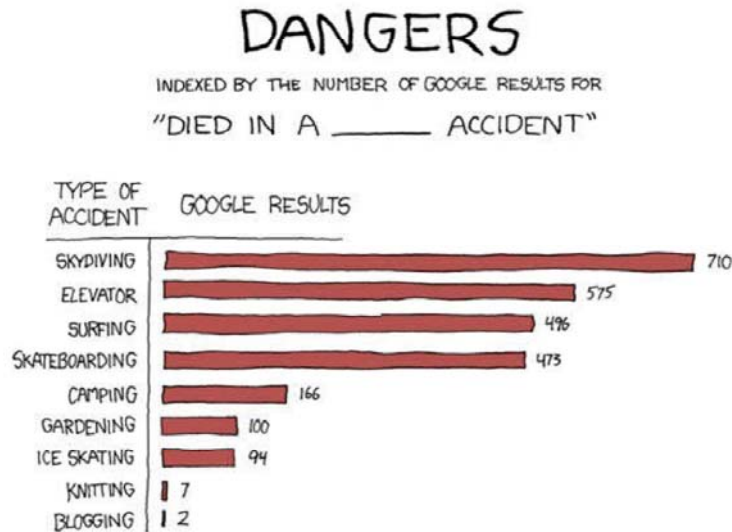- **Genomic annotation**
- **Whole cellular simulations**

Predictive Understanding

**KBASE**
predictive biology

DOE Systems Biology Knowledgebase

**Data generators**

**Open-Access Data and Information Exchange**
- Access through several flexible interfaces
- Retrieval of experimental data and products of modeling and simulation
- Working environment for testing hypotheses by *in silico* experimentation
- Provision and sharing of user feedback

**Data users**

**Seamless Submission and Incorporation of Diverse Data**
- Standards for data and metadata representation
- Quality control and assurance capabilities
- Automated systems for depositing and updating bulk data
- Tracking and evaluation of data use

**DOE Systems Biology Knowledgebase**

Data

Meta-data

Tools

**Community-Wide Stewardship**
- Advisory committee
- New social networking tools and capabilities
- Development
- Value-added analysis
- Training, tutorials, support staff

**Software and tool developers**

**Open Development of Open-Source Software and Tools**
- Data analysis and visualization tools
- Resources for *in silico* experimentation
- Modeling and simulation tools
- Customizable tools with layers of functionality
- Application programming interfaces (APIs) for open development of software
- Tracking and evaluation of tool use

**DOE Office of Science • Office of Biological and Environmental Research**

# Kbase, working to build a(n)…

- *Knowledge*base enabling *predictive* systems biology.

- Powerful modeling framework.

- **Community-driven**, extensible and scalable **open-source** software and application system.

- Infrastructure for integration and reconciliation of algorithms and data sources.

- Framework for standardization, search, and association of data.

- Resource to enable **experimental design** and **interpretation** of results.

# Data is extracted and displayed

## DANGERS

INDEXED BY THE NUMBER OF GOOGLE RESULTS FOR

"DIED IN A _____ ACCIDENT"

| TYPE OF ACCIDENT | GOOGLE RESULTS |
|---|---|
| SKYDIVING | 710 |
| ELEVATOR | 575 |
| SURFING | 496 |
| SKATEBOARDING | 473 |
| CAMPING | 166 |
| GARDENING | 100 |
| ICE SKATING | 94 |
| KNITTING | 7 |
| BLOGGING | 2 |

This is the database model

**Knowledgebases** should *learn* a "model" of the data to provide "conclusions" (hypotheses)

M-C Jenkins (http://www.scienceforseo.com)
Images from xkcd comics

**Databases** enable the rapid organization and **search** of data

# Knowledge is learning & answering

PROBABILITY BOOK IS GOOD

NUMBER OF WORDS MADE UP BY AUTHOR

"THE ELDERS, OR FRAÁS, GUARDED THE FARMLINGS (CHILDREN) WITH THEIR KRYTOSES, WHICH ARE LIKE SWORDS BUT AWESOMER..."

This is the knowledgebase model

*Utilize existing commercial software technology and leveraging DOE internet resources (ESNet) and DOE cloud computing platforms (Magellan)*

> **Kbase is a framework for data collection, integration and analysis tools to enable the simplified use and modeling of large scale genome and genome enabled information**

## By 2013 Deliver on the Initial Goals:

▪Kbase Infrastructure firmly established at 4 laboratories including high performance and cloud computing with routine 10+ Gb/s data transfer over ESNet between all Kbase sites

▪First public release includes:
  o Integration of data to reconstruct and predict metabolic and gene expression regulatory networks for up to 1,000 microbes to manipulate microbial function.

  o Integration of phenotypic and experimental data for bioenergy plants to predict metabolic and regulatory genotypes enabling manipulation of biomass properties.

**DOE Office of Science • Office of Biological and Environmental Research**

# HOW ARE WE GOING TO DO IT?

# Enable DOE Mission Science

**Communities**

**Plants**

**Microbes**

KBase Workflow Model

Driving data towards dynamic models of function

Functional Inference in Genomes and Metagenomes

Model-based Analysis & Eng

Metabolic, Regulatory, and Community Network Inference

Comparative Analysis of Predictions

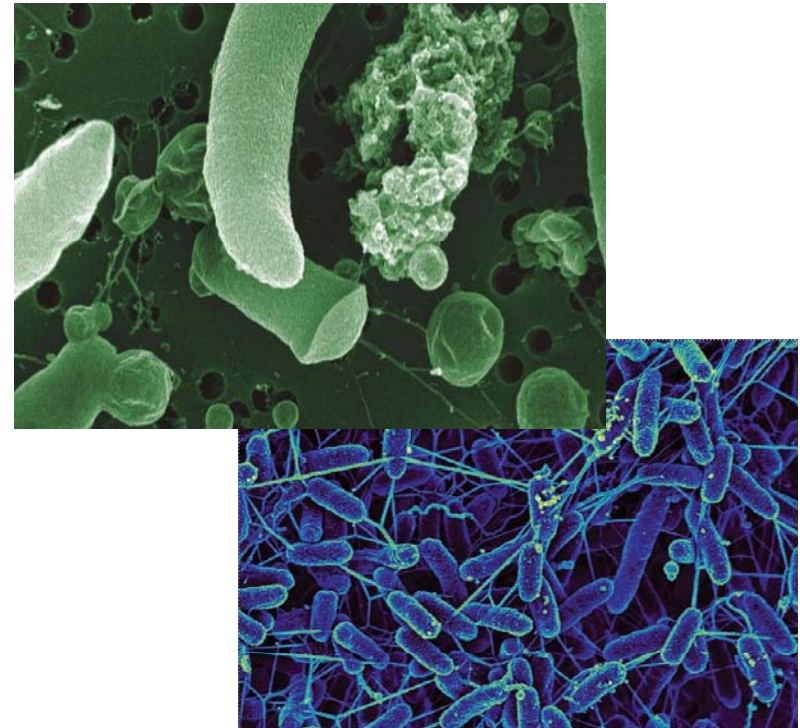| Inference of gene structure and annotation by homology | Direct measurement and Guilt-by-Association Function Inference | Inference of networks and suggestions for hole-filling | Behavioral Prediction And Design |

**KBase Workflow Model**

Driving data towards dynamic models of function

← Functional Inference in Genomes and Metagenomes ← Model-based Analysis & Eng →

Comparative Analysis of Predictions

← Metabolic, Regulatory, and Community Network Inference → ←

| Inference of gene structure and annotation by homology | Direct measurement and Guilt-by-Association Function Inference | Inference of networks and suggestions for hole-filling | Behavioral Prediction And Design |

Measures of Confidence and Quality

Clearinghouse of formal Predictions/Hypotheses

User Communities

Inference of gene structure and annotation by homology

Direct measurement and Guilt-by-Association Function Inference

Inference of networks and suggestions for hole-filling

Measures of Confidence and Quality

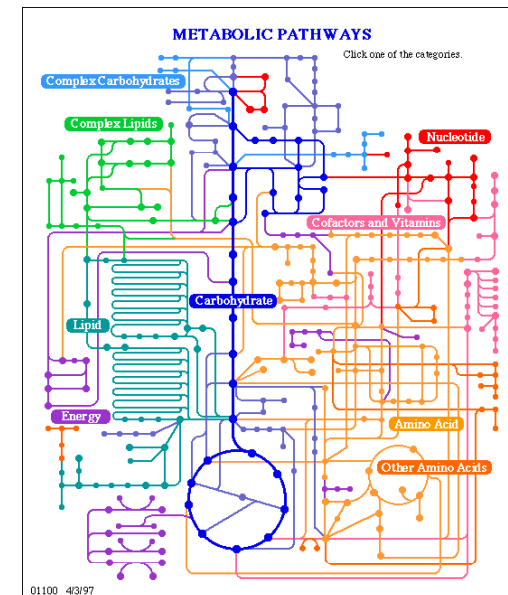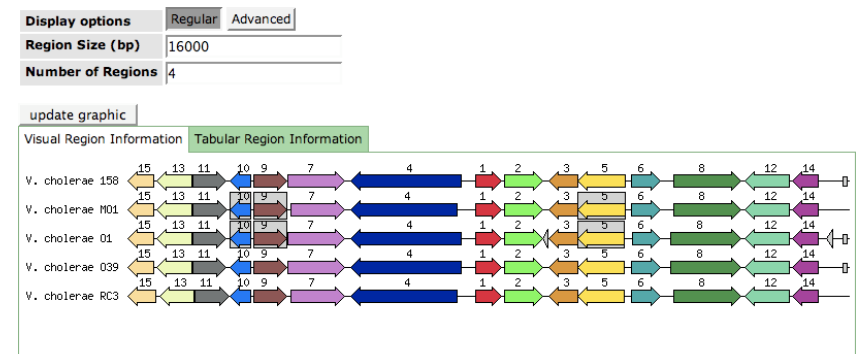Clearinghouse of Predictions/Hypotheses

User Communities

## Our overall goals are to:

- Reconstruct and predict metabolic and gene expression regulatory networks to manipulate microbial function

- Vastly increase the capability of the scientific community to communicate and utilize their existing data

- Enable the planning of effective experiments and maximizing understanding of microbial system function
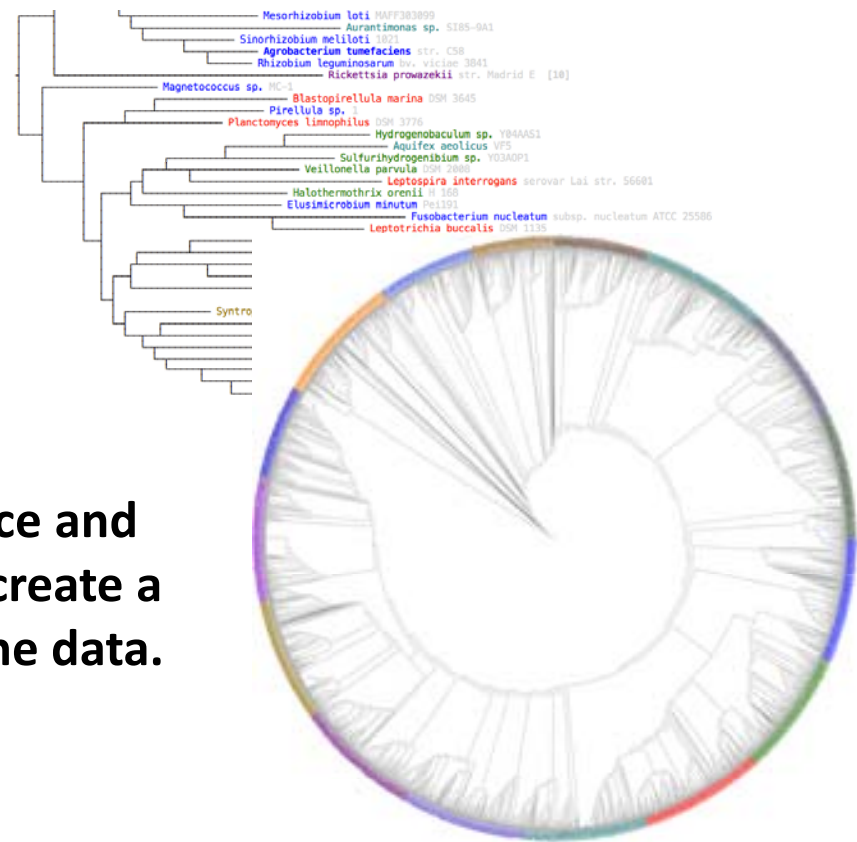
## We propose to do this by:

- Annotating genomes and assigning confidence

- Reconstructing metabolism and optimizing for function

- Reconstructing regulation and assessing agreement with expression data

- Integrating and standardizing -omics data from multiple data sources

- Constructing models of microbial organisms and interlinking models with data

**Within 13 months, we will be able to demonstrate use of the following:**

- Data integration and data model
- Next generation organism pages
- Phylogenetic tree services
- Next generation gene pages
- Metabolic modeling
- Regulatory/Transcriptional networks

**A microbiologist with a genome sequence and phenotypic growth data will be able to create a metabolic model fully reconciled with the data.**

**KBASE**
*predictive biology*
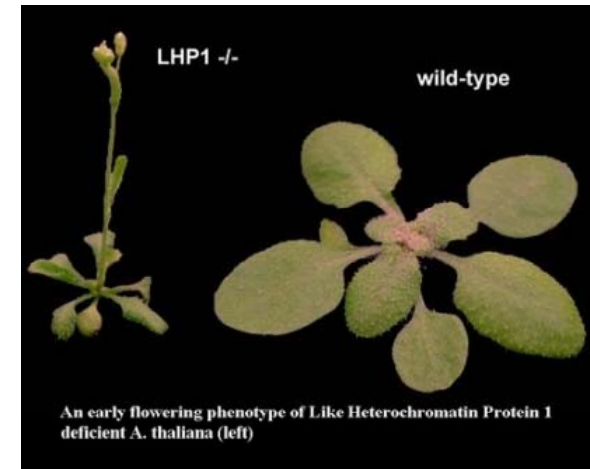DOE Systems Biology Knowledgebase

## Our overall goals:

*In order to extract knowledge from the wealth of high-throughput data in plant biology we need the ability to meaningfully integrate data.*



LHP1 -/-          wild-type

An early flowering phenotype of Like Heterochromatin Protein 1 deficient A. thaliana (left)

**Therefore we aim to:**

- Deploy Kbase capability that will allow for interactive, data-driven analysis and exploration across multiple -omics experiments.

- Provide researchers in plant sciences access to comprehensive datasets from high-throughput experiments together with relevant analytical tools and resources.

- Provide a platform for researchers to analyze their own experimental data, and have these results incorporated into the data exploration framework.

## We will accomplish this by using:

- Advance storage and indexing strategies for fast but persistent retrieval of large-scale genomics data

- Massively parallel processing of raw data using cloud compute resources

- Interactive charting libraries

- Using controlled vocabularies and ontologies for describing and storing genomic data
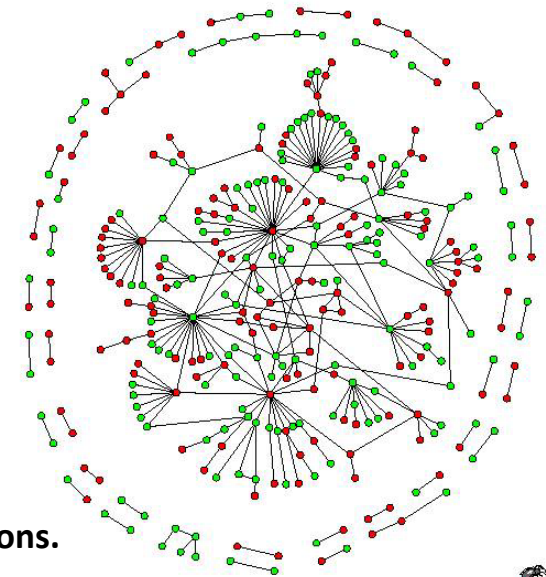
**Within 13 months, we will have the following capabilities:**

**Genotyping Workflow**
- Create a workflow for rapidly converting sequencing reads into genotypes
- Demonstrate more than 100-fold speedup over serial version by leveraging capabilities of KBase cloud
- Workflows will be developed as part of the KBase CyberInfrastructure

**Data exploration: Linking of gene targets from phenotype and genotype studies with co-expression, protein-protein interaction, and metabolic models**
- Allow users to narrow candidate gene lists based on these integrated data types
- Recommender system based on "guilt-by-association" principle: identify other genes, expression datasets, etc. associated with the user-selected group of genes
- Project genetic variations and network edges onto metabolic pathways
- Visualize co-expression network and node Interactions among subnetworks correlated to phenotype of interest

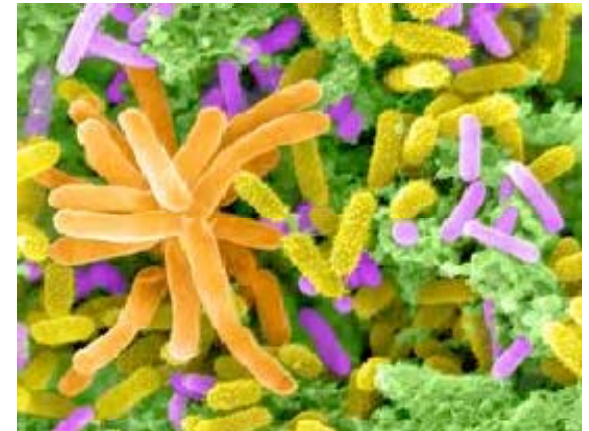**Plants work is cross-cutting with microbial and  communities data and predictions.**

**There has been an explosion of metagenomics data:**

- Systems biology is driven by the ever-increasing wealth of data
- Metagenomics is >90% of the data
- Computation needs to be smarter

**Our overall goal is to build a Kbase metagenomic platform that provides:**

- Scalable, flexible analyses, link physiological and metadata sets to metagenomic sequences

- Data QC and GSC compliant data and standards for data collection

- Enable modeling of metabolic processes within a community

- Predict microbial growth in isolation and in a community

**Within 13 months, we will have the following capabilities:**



**Metagenomic Experimental Design Wizard**

The foundation laid will enable researchers to perform *in silico* experimentation and hypothesis testing

**Bioprospecting**

- Find communities with similar alpha diversity

- Find communities in similar biomes

- Locate novel proteins (unknowns)

- Suggest functions (based on metadata) that might be encoded in "abundant unknowns"
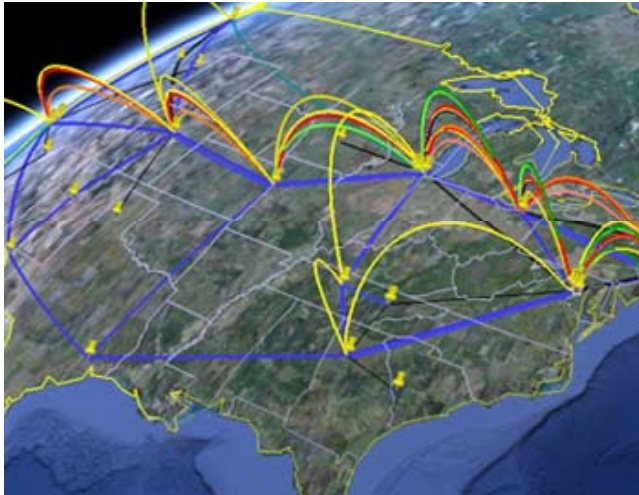
- Identify optimal candidates for screening

Communities work is cross-cutting with microbial and plant data and predictions.

**KBase is providing robust and scalable infrastructure to support new science by offering…**

- High speed data transfers will enable better response times
- Access to back end storage systems for large data sets
- Remote compute services for HPC, cluster and cloud based
- Workflow support
- Web UX access to data and computing
- Persistent and transient data management capabilities
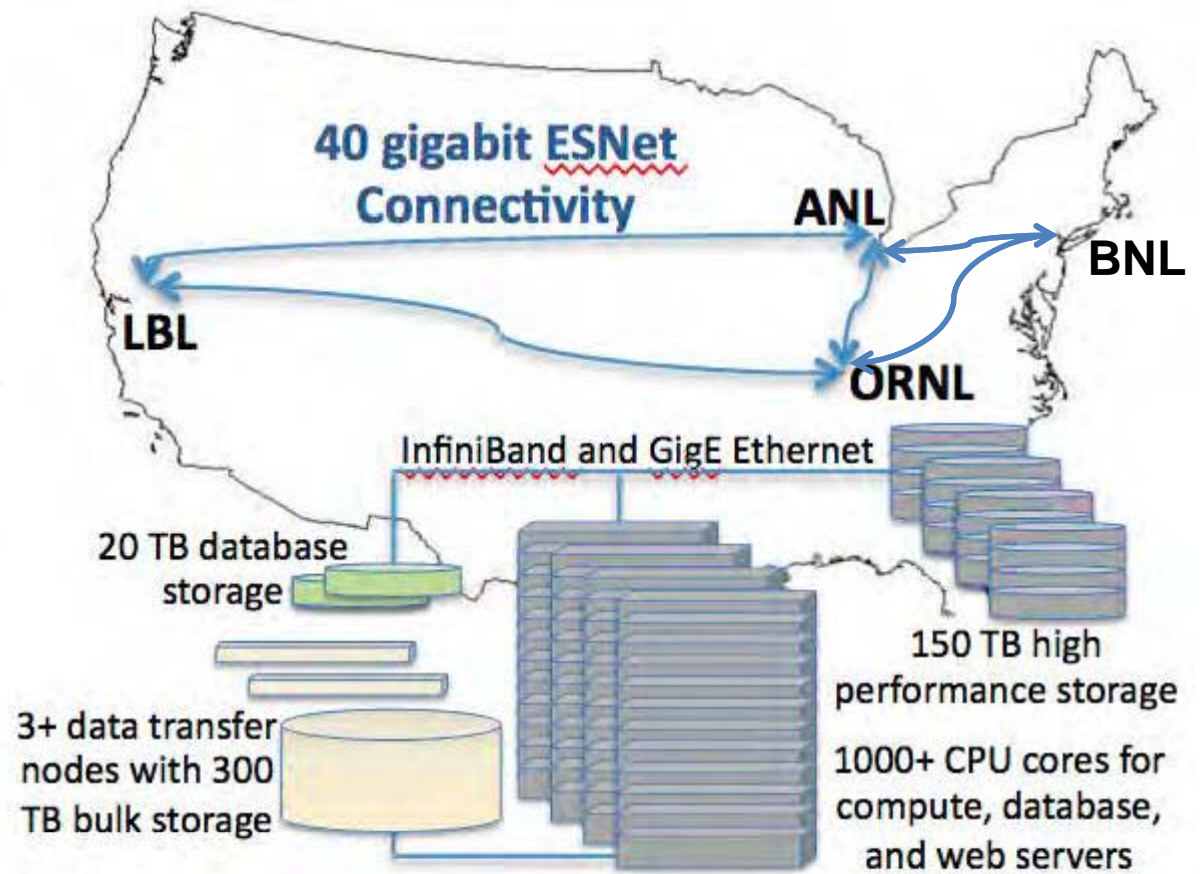- Support for users, teams, projects and cross talk

**DOE Office of Science • Office of Biological and Environmental Research**

Kbase leverages ESNet for 10+ Gb/s data transfer between all nodes



- ESnet backbone ( ESnet4) is anational 10 Gbps optical circuit infrastructure

- ESnet shares its optical network with **Internet2**

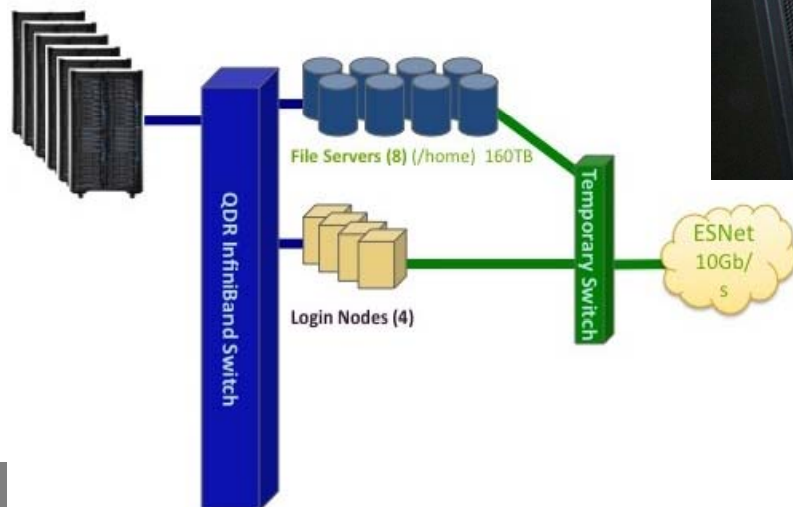- ESnet's IP network functions as a Tier 1 internet service provider

![KBASE predictive biology - DOE Systems Biology Knowledgebase]

**Argonne Magellan Hardware**

**Built on DOE Magellan Cloud**
**Magellan Goals:**

Establish a nationwide scientific mid-range distributed computing and data analysis testbed.

http://magellan.alcf.anl.gov/
http://science.energy.gov/ascr



Current config at 252nodes
Expected to reach ~700 nodes in FY12

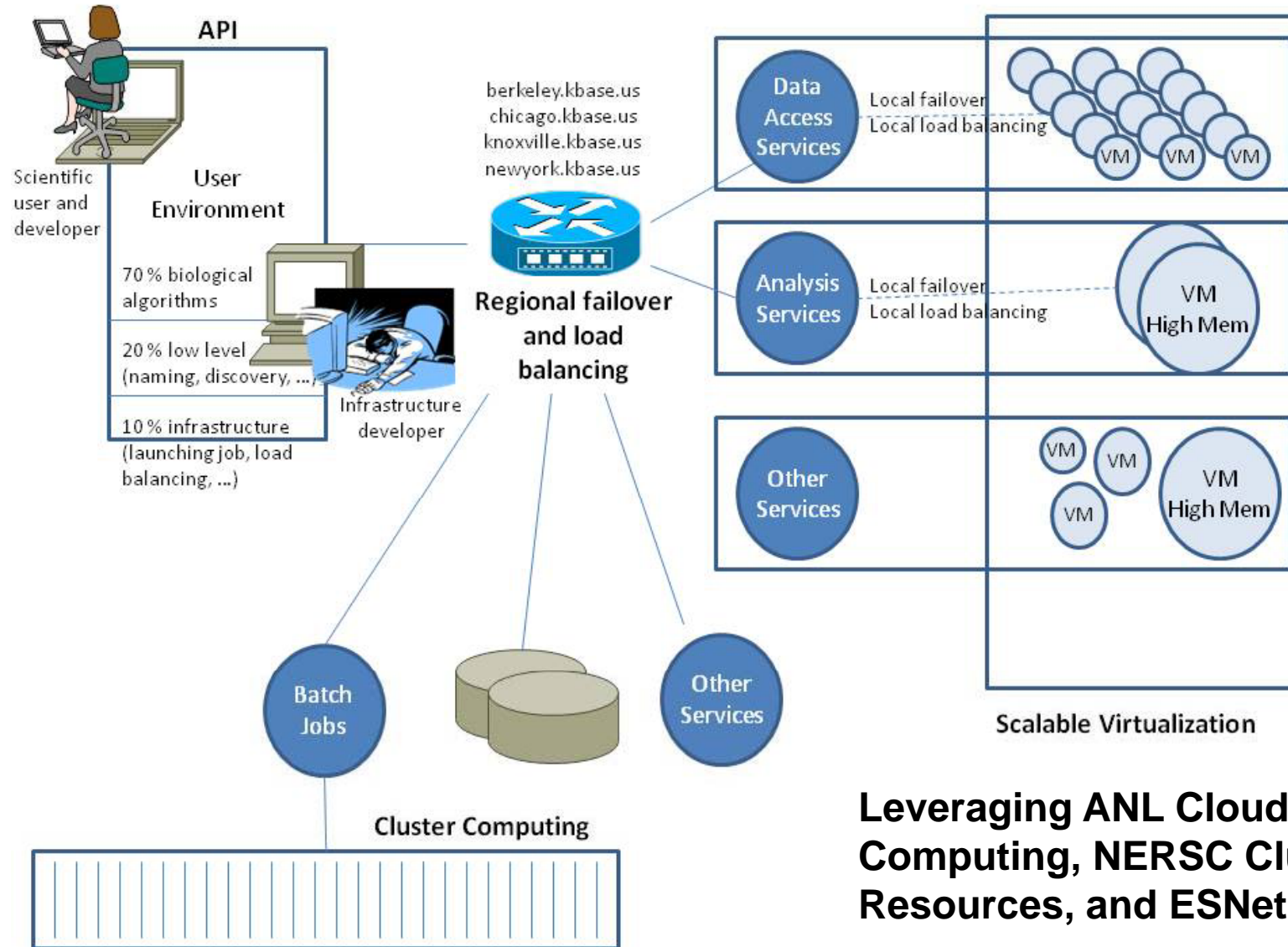Smaller system housed at ORNL with around ~70 nodes

**Leveraging ANL Cloud Computing, NERSC Cluster Resources, and ESNet**

> *KBase will integrate initial f core services and enable KBase Unified Application Programming Interface (API)*

> The *KBase Unified API is the primary programming target for KBase tools and applications*

> *KBase Unified API provides*

- Access to core functions and data
- Low-level systems services
- User interface and graphics services and
- Integration services



Kbase Core Services: Microbes Online, Meta Microbes Online, RegFam, SEED, RAST, Model SEED, MG-RAST, Phytozome, External Cores, Future KBase Cores

➤ Uniform view of KBase core services

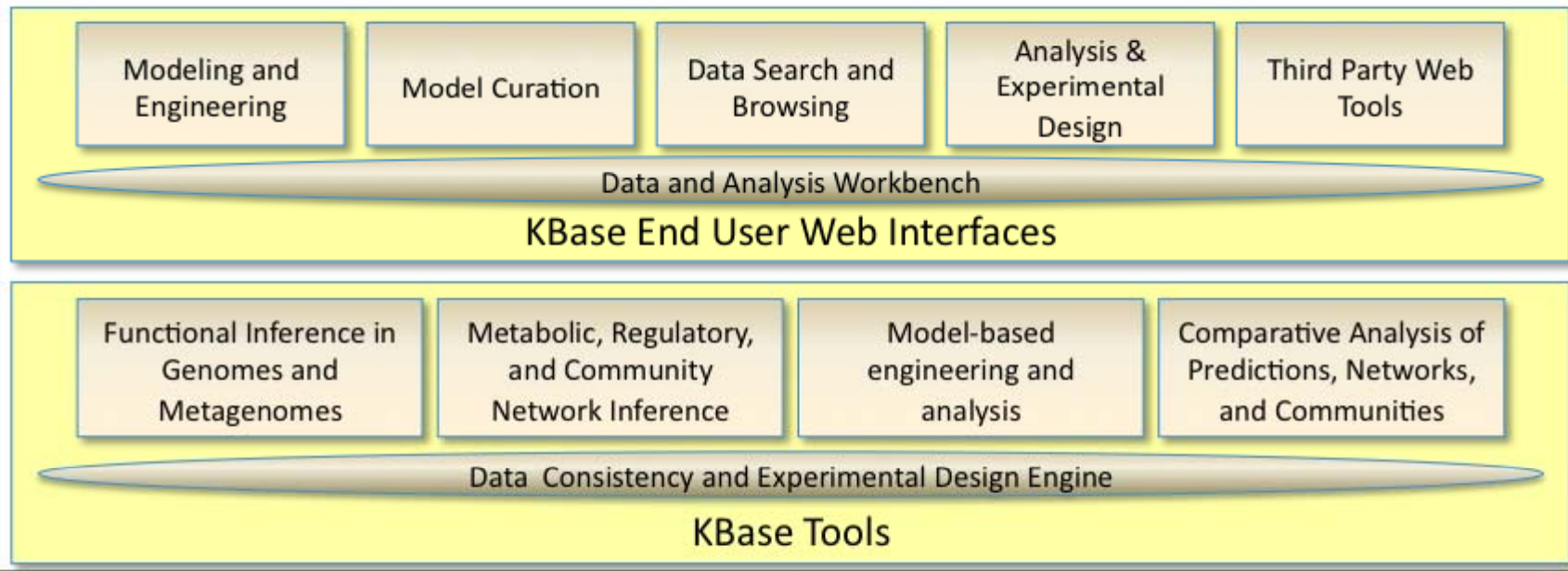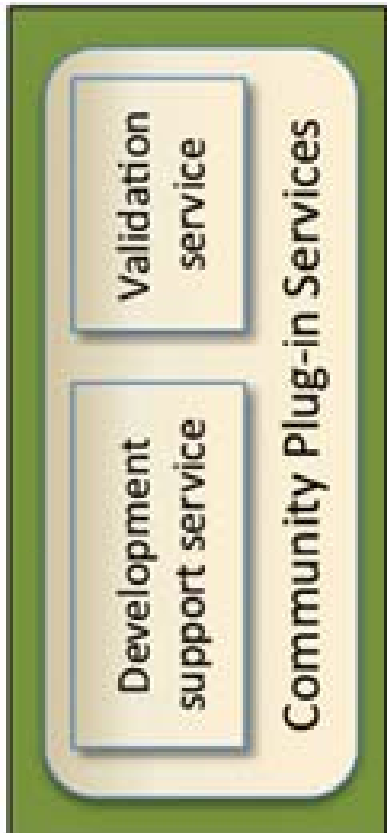➤ Access to data management, security and computational services in the KBase infrastructure

➤ Rich set of standard graphics and UX components for building web applications with single look and feel

KBASE
predictive biology

DOE Systems Biology Knowledgebase

➢Most KBase code development will be layered on top of the KBase Unified API

➢Supports a layered architecture consisting of collections of KBase Tools and End User Web Interfaces

➢Each layer extends the API by providing additional support for high-level functions



| Modeling and Engineering | Model Curation | Data Search and Browsing | Analysis & Experimental Design | Third Party Web Tools |

Data and Analysis Workbench

**KBase End User Web Interfaces**

| Functional Inference in Genomes and Metagenomes | Metabolic, Regulatory, and Community Network Inference | Model-based engineering and analysis | Comparative Analysis of Predictions, Networks, and Communities |

Data Consistency and Experimental Design Engine

**KBase Tools**

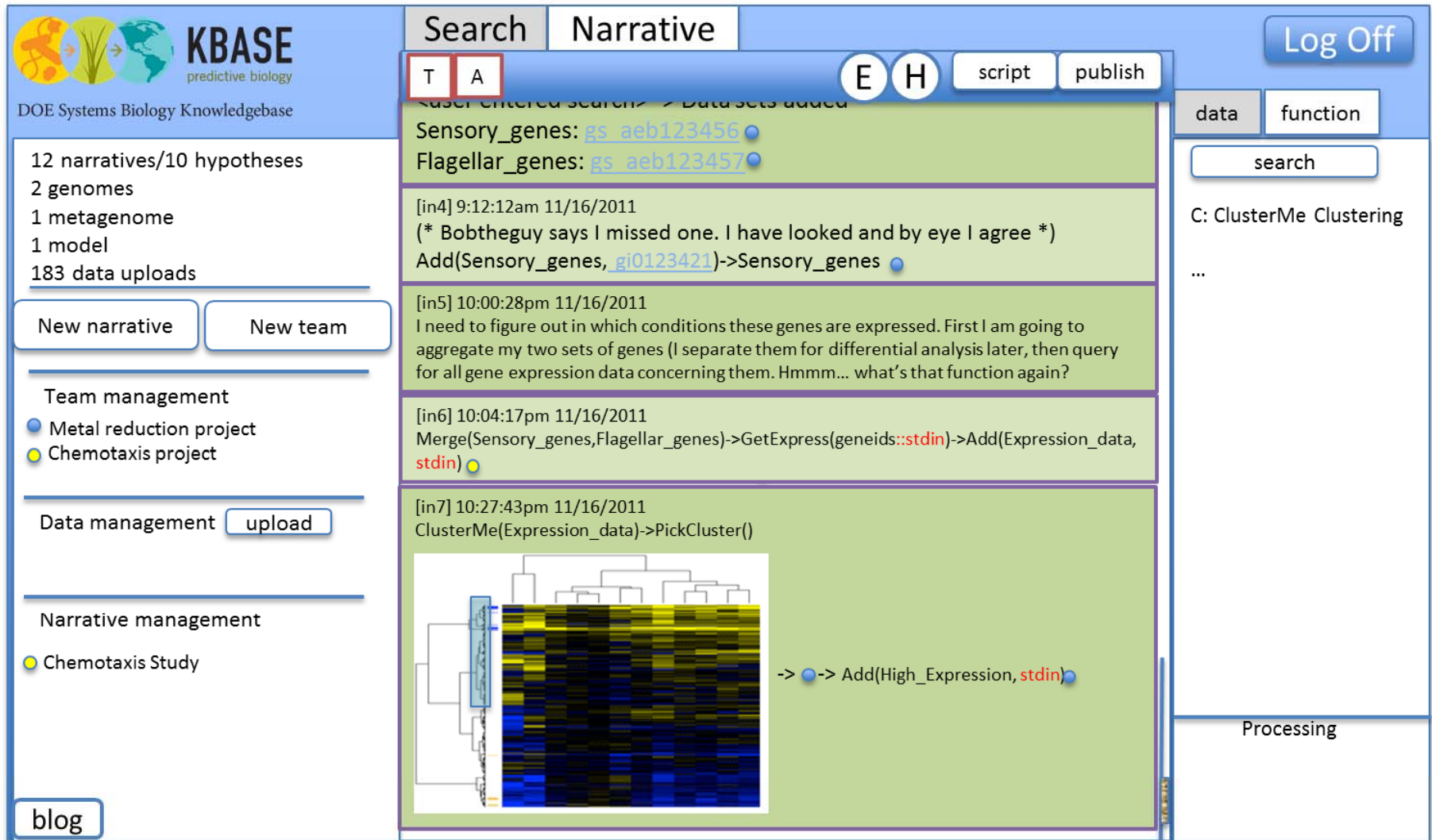**DOE Office of Science • Office of Biological and Environmental Research**

## Two ways to work with KBase

➢Extend from below with additional core services that *extend the KBase API*

➢Extend from above by plugging in applications that *use the KBase API*

Concept: Kbase User Experience

Narrative Graphs Provide Measures of Activity and Influence

Function and Data contents can be tracked to assess "use"

Scenario A

A.1

Scenario B

Narrative Idea

Narrative Code ...sioning and ...nching

Narrative Query Update

Narrative Data Change

Project Linkage and Citation

# KBase will allow users to…

- Keep pace with the rate of data generation.

- Use interrelated data of different types.

- Develop analytical tools.

- Perform rapid comparisons across multiple genomes.

- Allow free and open access to KBase data, models and simulations.

- Grow partners, developers, data generators, or users.

- Coordinate with other researchers to target what experimental work is needed, and to share results.

- Learn how to use KBase features of models and move into predictive systems biology.

# BUILDING THE KBASE COMMUNITY

# Who are the Kbase core developers

## The KBase Team:

## Collaborators:

- Cold Spring Harbor Laboratory
- University of California, Davis
- Hope College in Michigan
- University of Illinois at Urbana-Champaign
- Yale University
- The Joint Genome Institute
- InterBRC Knowledgebase
- Several University Knowledgebase Projects

→**New Partners and Collaborators from Users and Stakeholders**

**DOE Office of Science • Office of Biological and Environmental Research**

**KBASE** predictive biology
DOE Systems Biology Knowledgebase

| LBNL | ANL |
|------|-----|
| BNL | ORNL |

**Team lead
* Team co-lead

*Steering Committee Member*

**DOE Program Manager**

**Board of Directors**

**Governance Board**

**Executive Team/Management**

| Sergei Maslov, ASO | Rick Stevens, CTO | ⭐Adam Arkin, CEO, CSO | Bob Cottingham, COO, Planning and Partnerships |

**Project Manager ACOO**

**Science Domains Crosscut**

| Science Team Lead/ CSO: **Adam Arkin | Science Team Co-Lead/ASO: *Sergei Maslov |

| Microbes | Plants | Communities |
|----------|--------|-------------|
| **Paramvir Dehal | ** Doreen Ware | **Folker Meyer |
| Steven Brenner | Sergei Maslov | Narayan Desai |
| Ben Bowen | Jer-Ming Chia | Andreas Wilke |
| Pavel Novichkov | Michael Schatz | Jared Wilkening |
| John-Mark Chandonia | Mark Gerstein | Elizabeth Glass |
| Keith Keller | | |
| | *David Weston | **Dylan Chivian |
| *Chris Henry | Priya Ranjan | John Bates |
| Scott Devoid | Guruprasad Kora | |
| FangFang Xia | | |
| | Chris Henry | |
| | Elizabeth Glass | |
| | David Goodstein | |
| | Pam Ronald | |

**Infrastructure and Services (CI)**

| CI Team Lead/CTO: ** Rick Stevens | *Tom Brettin |

| | |
|---|---|
| Narayan Desai | Shane Canon |
| Andreas Wilke | Dantong Yu |
| Dan Olson | Shinjae Yoo |
| Ross Overbeek | |
| Robert Olson | Daniel Quest |
| Terry Disz | Michael Galloway |
| | Miriam Land |

Other Roles:
JGI POC: Rick Stevens
BRC-KB POC: Adam Arkin
Science Team POC for Outreach: Paramvir Dehal

**Outreach and Public Relations**

**Outreach**

**Elizabeth Glass
*Brian Davison
Jennifer Salazar
Bob Cottingham
Paramvir Dehal
Dylan Chivian

**Public Relations**

**Jennifer Salazar

Outreach Reps for Other Teams:
Plants, Communities: Elizabeth Glass
Microbes: Brian Davison
CI, Management: Jennifer Salazar

## Formal project management under 413-lite

The Kbase core team is creating the initial infrastructure to:

- Allow effective and easy data federation and centralization.

- Lower the bar to adding and accessing new algorithms and analyses

- Enable nimble use of of advanced DOE network and compute resources

- Create a community framework for development and biological knowledge sharing

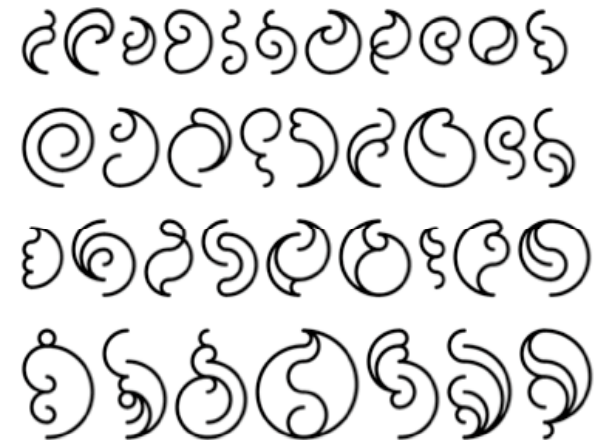- Enable *their own research and yours!*

**KBASE:**

A. Professional Computational Biologists
B. Data generators and basic analysts
C. Knowledge Seekers
D. Knowledge Generators

instances of "minimum inventory/maximum diversity" systems, a term coined by Peter Pearce in his book, Structure in Nature Is a Strategy for Design (MIT Press, 1978).

***Therefore we aim to:***

- Create a powerful framework for programmatic access to data and functions of Kbase. (Users A,B)
  - Ultimately provide stubs for use in PERL, PYTHON, R, MATLAB, Galaxy, etc.

- Create a set of packaged "Widgets" that make placement and recognizable display of Kbase "functions" on web pages (or within perhaps other apps), easy and identifiable.  (Users B)

- Create a "simplified" portal for search and aggregation of data for data consumers and Knowledge Seekers. (Users C,D)

- *Create a innovative platform for knowledge creation, evolution and sharing.*

**KBASE**
*predictive biology*
DOE Systems Biology Knowledgebase

## We will build a diverse community of users by…

- Creating a unique a powerful resource for people to share new analyses and data and the use of these to produce new science.

- Advertising and informing via:
  - Online communication
  - Conference workshops
  - Press releases

- Interacting with the community and fostering relationships

- Collecting feedback from the community to ensure we are meeting their needs

- Measuring our impact

- Participating in standards groups



**Community science drives development!**

(Planning → Requirements → Analysis and Design → Deployment → Testing → Evaluation)

**Visit us at kbase.science.energy.gov**

**Contact us at outreach@kbase.us**

# WHEN WILL WE DO IT?

**KBASE**
*predictive biology*

DOE Systems Biology Knowledgebase

- Feb 2012 – Development release (internal target)
  - debug release engineering, prototype deployments, initial data models and data loads, non-unified API, performance testing, architecture refinement
- May 2012 – Alpha release (internal target, limited invited testing)
  - draft tutorials, v0.0 database loads, draft API (performance and ubiquity unified prototypes), draft UI library, domain workflow drafts, cloud and cluster services
- Aug 2012 – Beta Release (early adopter beta testing)
  - workflow function complete, API refinement, v1.0 database loads, prototype plugin interfaces, prototype galaxy support, performance debugging
- Nov 2012 – Production Release Candidate (public beta testing)
  - draft website, draft documentation, full functionality API, draft UI v1.0, database loads v2.0, significant number of external beta test users
- Feb 2013 – KBase Production Release
  - public website, unified API, initial production UI, database loads v3.0 (microbes, community, plant databases), production demonstration workflows, replication and fail over services,

# Outreach Timeline?

Website, Online Resources — Nov 2011

Workshop at DOE Annual Meeting — Feb 2012

You are here!

Session at JGI Annual Meeting — March 2012

Ongoing workshops and sessions at conferences — …

Draft Documentation and Tutorials

Feb 2013

Beta Release of Demos

Microbes   Plants   Communities

**For more information, and to follow our progress, please visit**
**http://kbase.science.energy.gov/**
**http://outreach.kbase.us**

# Thank you!

**Visit us at kbase.science.energy.gov**
**Contact us at outreach@kbase.us**