



JGI Data Capabilities and Infrastructure

BERAC

April 21, 2023



Mission

91

The mission of the JGI is to provide the global research community with access to the most advanced integrative genome science capabilities in support of the DOE's research mission.

2,243

Primary Users leveraging JGI data generation capabilities in FY22

15,219

Secondary Users that engaged with JGI science gateways

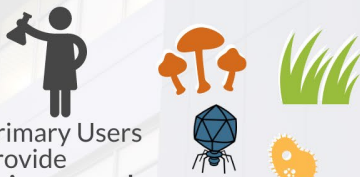
14PB

JGI Data Repository size as of December 2020

2,862

Publications

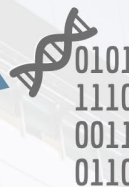
Data Generation and Reuse



Primary Users provide **unique samples** from fungi, plants, algae, bacteria, archaea, and communities as part of their studies



Samples become data



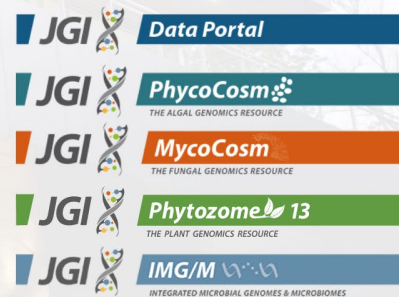
FY22 Downloads

4.1M Files

1.2PB Data



Primary and Secondary Users leverage data through JGI Flagship Science Gateways



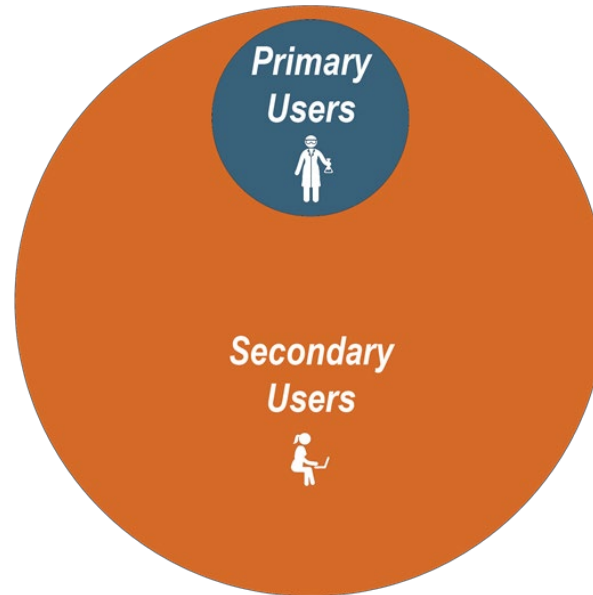
Primary and Secondary Users

Primary Users

are associated with one or more JGI **User Program proposals**

Secondary Users

build on the work of JGI personnel and primary users through **direct downstream use** of JGI data, systems, and tools.



Example Outcomes

- Publications
- Patents
- Software Adaptations
- New Technologies
- Marketable Products
- Methods & Standards
- Start-ups
- Grant Funding

JGI Contributes Data to Integrated Projects

MODELS

Improvement of watershed models to include chemical and biological processes



DATA



Data assembly, integration, and storage

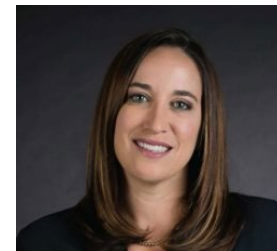


DISTRIBUTED SCIENCE APPROACH

Geochemistry, hydrology, metabolites, metagenomes, and metatranscriptomes

GROW

Genome Resolved Open Watersheds



Kelly Wrighton



Mikayla Borton

U.S. DOE. 2019. Open Watershed Science by Design: Leveraging Distributed Research Networks to Understand Watershed Systems Workshop Report, DOE/SC-0200, U.S. Department of Energy Office of Science.

JGI Public Resources for Data, Metadata, & Analysis



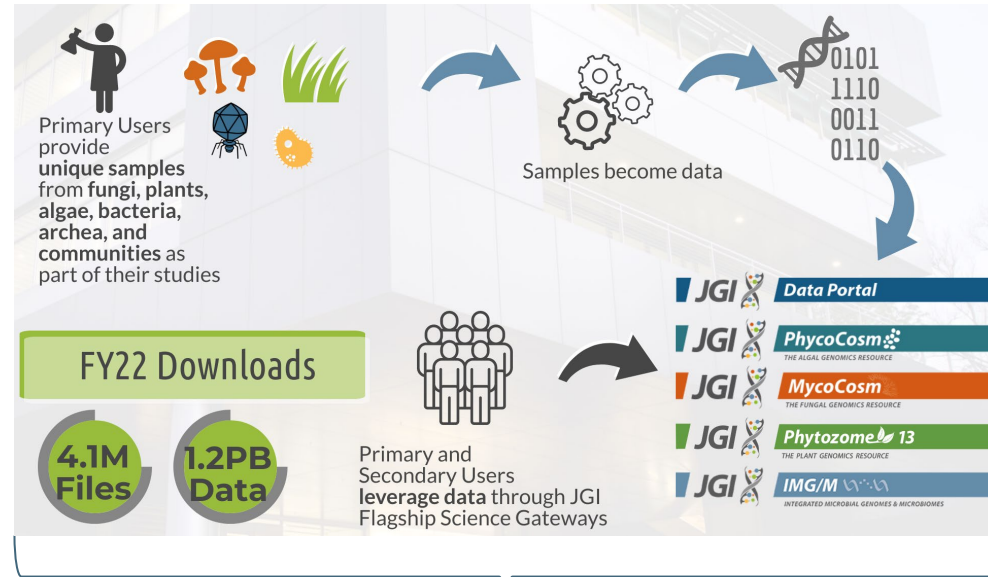
JGI Systems



External Systems

Centralized Data Access and Movement across Distributed Resources

- The JGI Archive and Metadata Organizer (**JAMO**) deployed in 2013
- Holds the **metadata and locations** for data produced by JGI
- **Powers data distribution** across JGI storage systems (file system and archives)
- Makes **centralized search** possible
- Supports **reusability research**



The JGI Data Portal – Access to all Public JGI Data



The screenshot shows the JGI Data Portal homepage. At the top, there is a navigation bar with 'Data Portal' and links for 'JGI HOME', 'GENOME PORTAL', and 'CLASSIC DOWNLOAD'. A yellow banner below the navigation bar contains a warning: 'Data Portal will be unavailable during system maintenance scheduled for Sunday, April 11, 2021 from 3:00-4:00a ET'. The main header features a green and blue abstract background with the text 'Top-quality genomic data, open to all researchers' and 'Explore and download invertebrate genomes and metagenomes'. Below this is a search bar with the placeholder text 'Search by genome, metagenome, project, or ID'. The page is divided into two main columns. The left column has a 'Search' section with a description of the data and a link to 'More about search'. The right column has a 'Download' section with a description of how to download data and a link to 'View our API docs'.

Our data

The U.S. Department of Energy (DOE) Joint Genome Institute (JGI) is a DOE Office of Science User Facility located at Lawrence Berkeley National Laboratory (Berkeley Lab). The JGI takes great pride in producing high-quality genomic and metagenomic data outputs for our users and the community. We ensure consistent quality by taking the following measures:

- Starting with top-quality samples
- Conducting ongoing quality control
- Drawing on accumulated knowledge
- Producing deeper metagenome sequences
- Developing new tools
- [Learn more](#)

New releases
New genomes will be released in Spring 2021! See the [full list of new genomes](#).

JGI in the news
Get links to [recently published studies](#) that incorporate JGI-sequenced data.

Upcoming events
[Register for upcoming JGI webinars](#) on a variety of topics.

[Contact Us](#) [Site Us](#) [Accessibility/Section 508](#)
[Disclaimer](#) [Credits](#)

© 1997-2020 The Regents of the University of California.
Genome Portal version 6.18.44 content:107000885_jgiportal-web-1 Release Train 0.9 Mar-2020 15:05:00.023 PST Current Data 04-Mar-2020 08:57:24 PST

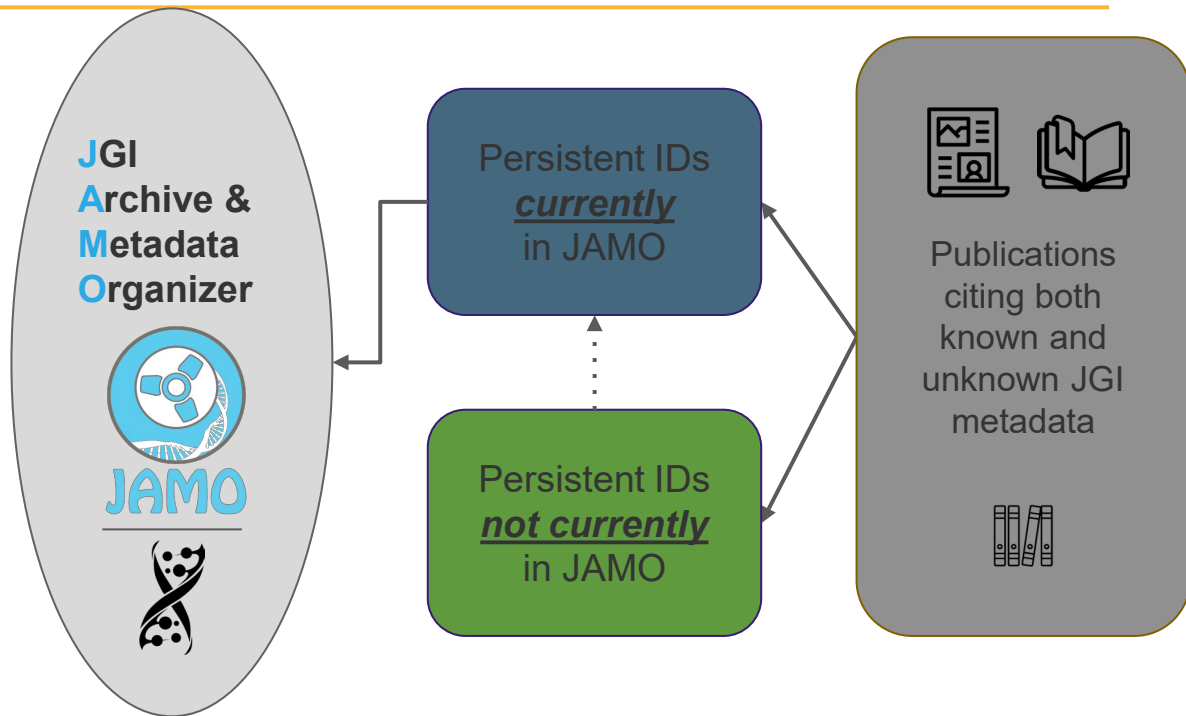


GUI: <https://data.jgi.doe.gov>
API: <https://files.jgi.doe.gov/apidoc>

Identifying Downstream Data Use

The Data Citation Explorer (DCE) Connects Data Use to User Projects

1. We know that many publications exist which cite JGI data
2. These publications cite both cataloged and uncataloged metadata
3. The DCE finds new metadata + citing publications and stores them in JAMO



Expanding Our Reach with Dimensions

How much does Dimensions expand JGI's search capabilities?

Using a sample
captured 118%
possible with

Learn more about JGI's Impact Analysis
work on 4/27 at the SC User Facility
Community of Practice talk – Neil Byers
will present!

207
papers

>98% validity for both sources

- 20,872 Citations

Data Citation Explorer – Augmenting the Data Portal



- Tie literature to data and display to data users
- Increase JGI understanding of institutional and individual impact
- Provide model for identifying data citations at other DOE repositories

Refine selections

0 files selected

Show Most Relevant Results Show All Results total 1 results shown

Everything grigoriav

Genome Acremonium chrysogenum ATCC 11550

Genome Aaosphaeria arxii CBS 175.79 v1.0 64 36 GB Select

Genome Abortiporus biennis CIRM-BRFM 1778 v1.0 135 182 GB Select

Genome metadata Acremonium chrysogenum ATCC 11550

PI name Hsi-Ling Liao

PI email sunny.liaoduff.edu

Work Completion Date 2019-12-23

Data Utilization Policy UNREST data_unrestricted

Proposal acceptance date 2018-08-16

Award DOI 10.46936/fics.proj.2017.50003/00006234

Links for more info

Genome details

Citation information

Citation: Acremonium chrysogenum ATCC 11550

Find the publications related to this dataset

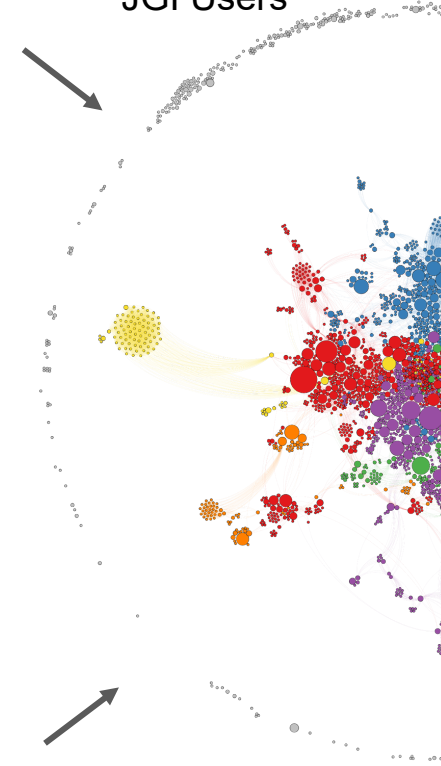
<< Return to the results page

Condon RL, Leng Y, Wu D, Bushley KE, Olin RA, Ollaly R, Martin J, Schwabatz W, Grimwood J, MohdZahudin N, Xue C, Wang R, Manning VA, Dillon B, Tu ZJ, Stephenson BJ, Salamon A, Sun H, Lowry S, LaBetti K, Han J, Copeland A, Lindqvist E, Barry K, Schmeitz J, Baker SE, Cluffetti LM, Grigoriav IV, Zhong S, Turgeon BG. Comparative genome structure of *Acremonium chrysogenum* ATCC 11550. *PLoS Genet.* 2013;9(1):e1003233. doi: 10.1371/journal.pgen.1003233. Epub 2013 Jan 24. PMID: 23357949; PMCID: PMC3554632. Primary Publication

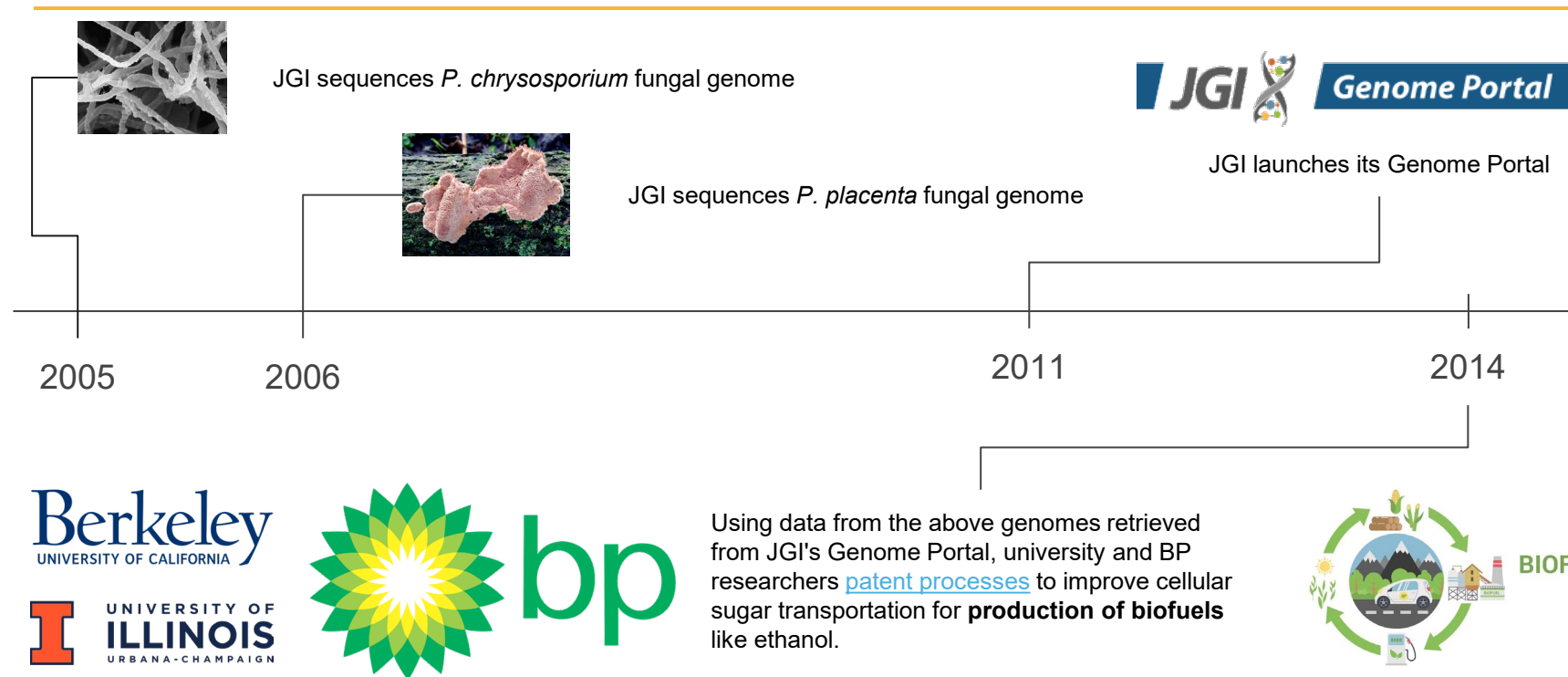
Piggeler S, Hoff B, Kick U. Asexual cephalosporin C producer *Acremonium chrysogenum* ATCC 11550 carries a functional mating type locus. *Appl Environ Microbiol.* 2008 Oct;74(19):6096-16. doi: 10.1128/AEM.01189-08. Epub 2008 Aug 8. PMID: 18689517; PMCID: PMC2565971.

Diez B, Mellado E, Fouces R, Rodriguez M, Barredo JL. Recombinant *Acremonium chrysogenum* ATCC 11550 strains for the industrial production of cephalosporin. *Microbiologia.* 1996 Sep;12(3):359-70. PMID: 8697416.

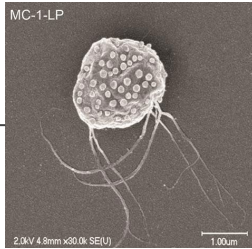
JGI Users



How does 'impact' begin?



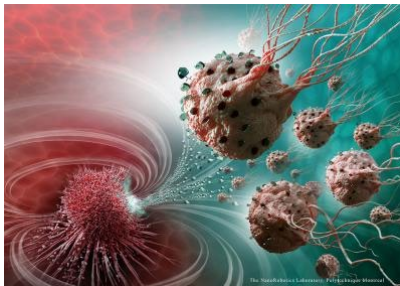
How does 'impact' begin?



JGI sequences the genome of *Magnetococcus* MC-1, a bacterium with special mobility traits that thrives in low-oxygen marine environments

2007

2016



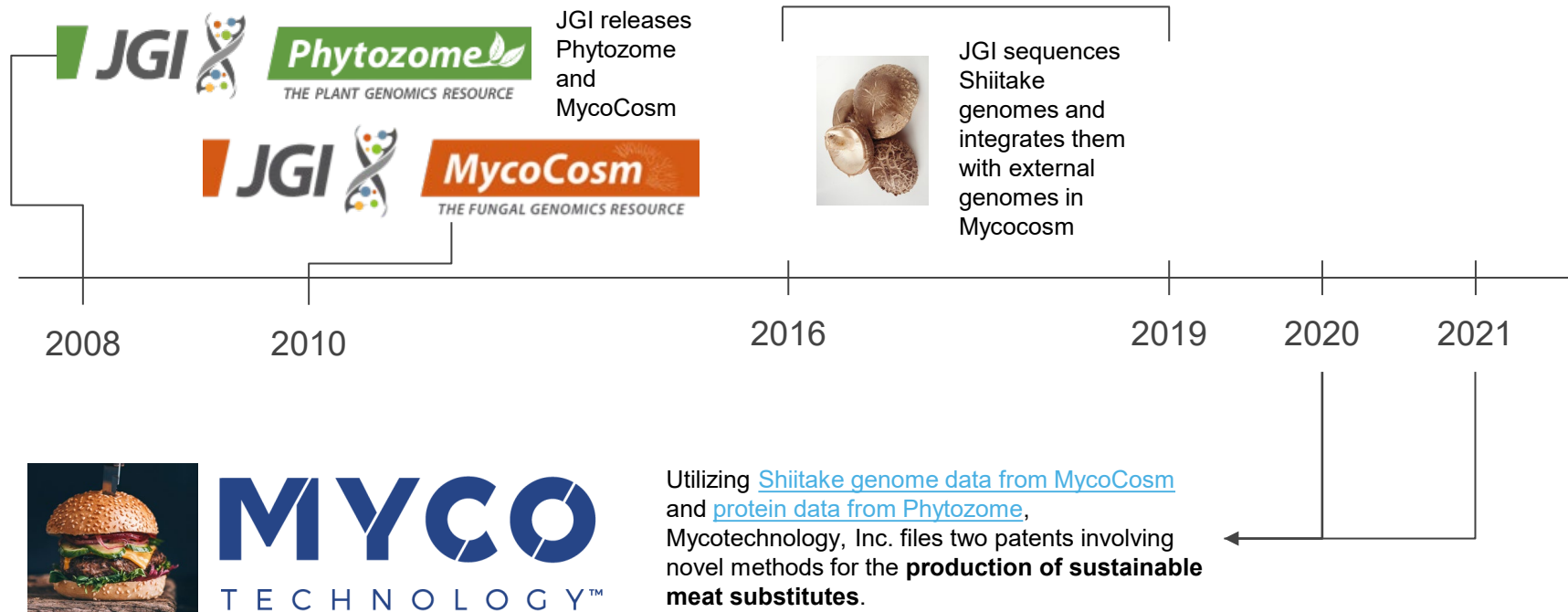
Aided by available genomic information, [researchers determine](#) that this bacteria is a very effective **medication delivery tool** for tumors in hard-to-reach areas of the brain.



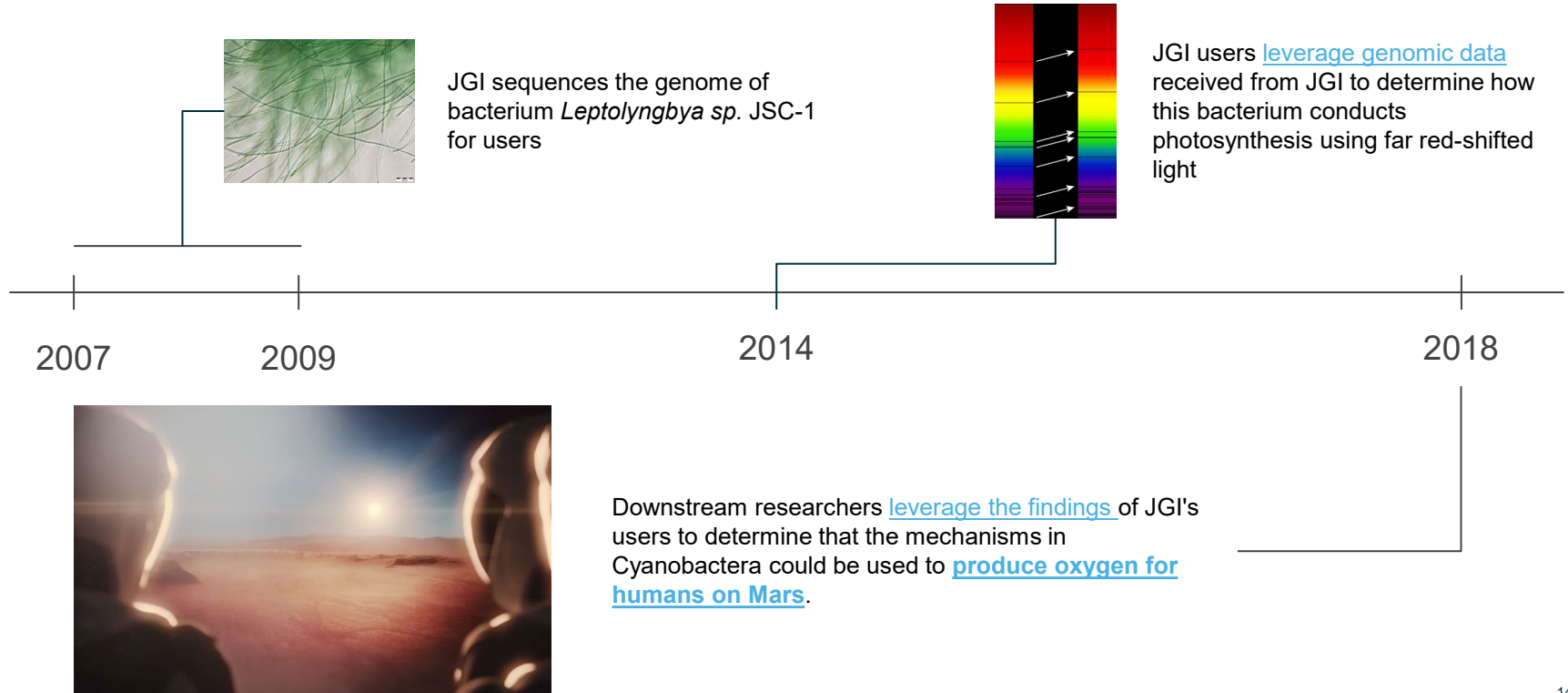
**POLYTECHNIQUE
MONTRÉAL**

UNIVERSITÉ
D'INGÉNIERIE

How does 'impact' begin?



How does 'impact' begin?

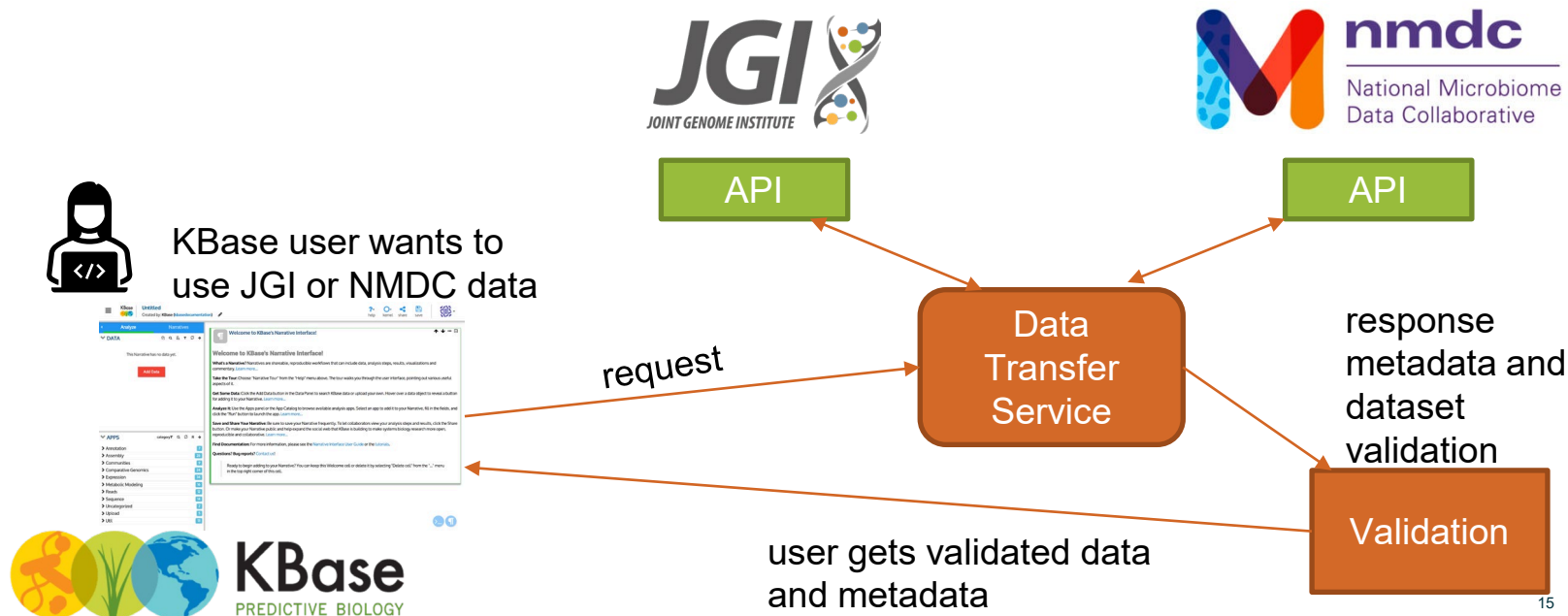


- **Validated Data Transfer**

- Make sure data moves between resources with citation and contextual information

- **Clear citation guidelines in Credit metadata**

- Consistent citations will aid in identifying data use in publications



Multiomics Data

Information about the type of samples being submitted

Alternative Names

Project, study, or sample set names that are also associated with this submission or other names / identifiers

GOLD Study ID

Provide the GOLD study IDs associated with samples for this study.

NCBI BioProject Accession

Provide the NCBI BioProject Accession Number associated with the listed NCBI BioProject Title.

Data types *

Check all -omics data types associated with samples collected for this study.

Other Non-DOE

- Metagenome
- Metatranscriptome
- Metaproteome
- Metabolome
- Natural Organic Matter (FT-ICR MS)

* indicates required field

[← GO TO PREVIOUS STEP](#)

Data types *

Check all -omics data types associated with samples collected for this study.

Joint Genome Institute (JGI)

- Metagenome
- Metatranscriptome
- Metabolome

JGI Proposal ID/Study ID *

This is the 6 digit ID assigned to your JGI Proposal and is required when completing metadata for samples to be sent to JGI for sequencing.

Environmental Molecular Science Laboratory (EMSL)

- Metaproteome
- Metabolome
- Natural Organic Matter (FT-ICR MS)

EMSL Proposal / Study Number *

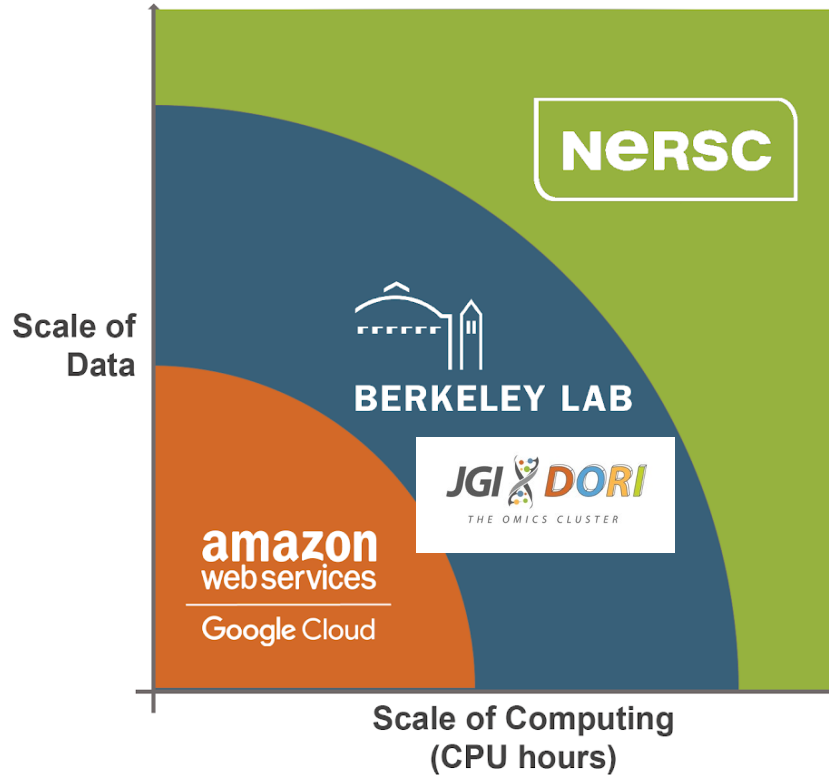
EMSL Study Number is required when processing was done at EMSL

Other Non-DOE

- Metagenome
- Metatranscriptome
- Metaproteome
- Metabolome
- Natural Organic Matter (FT-ICR MS)



JGI's Computing Infrastructure Spectrum



Binning JGI Compute Infrastructure Requirements



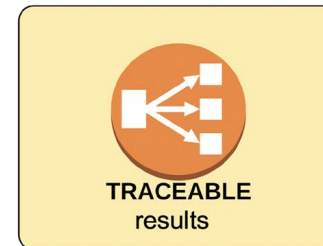
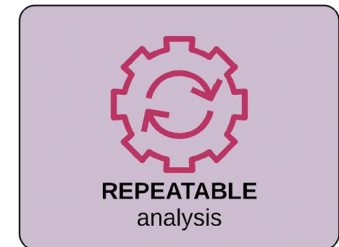
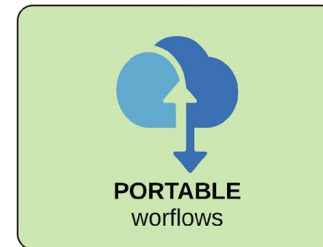
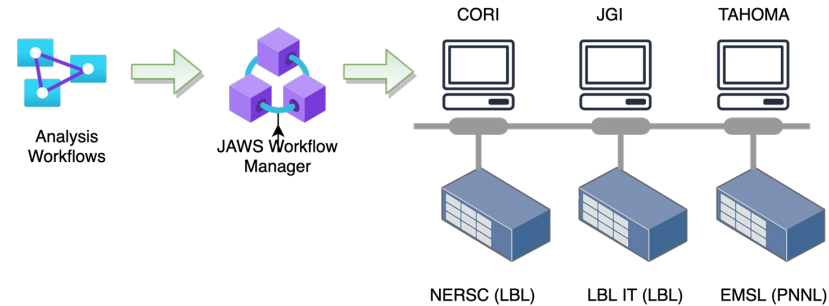
Prototyping,
exploratory
analysis, small-scale
production work



Large-scale data or
compute needs
(>100,000 CPU hours)

Our Solution to Unify Workflow Execution Across JGI

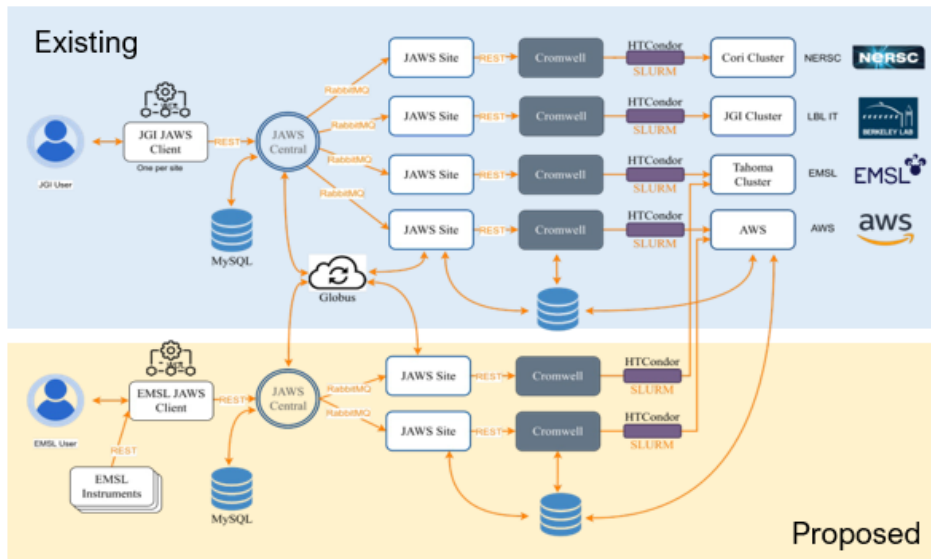
- Developed a workflow manager called **JGI Analysis Workflow Service (JAWS)** to run complex computational workflows with support for distributed computation across multiple HPC enabled sites.
- Provides a **user-friendly common interface** to seamlessly route jobs and data across multiple sites.
- Improves the **reusability** and **robustness** of bioinformatics workflows in evolving and/or diverse high-performance computing (HPC) and cloud environments.



- Consolidated workflow efforts between JAWS and NMDC team
- JAWS provides the overall framework for submitting, running and collecting results
- JAWS team will provide operational support for core JAWS service (initially for NERSC workflows)
- NMDC will be able to customize their cluster configurations and work with JAWS team on new features to support NMDC requirements
- JAWS team supporting NMDC workflow team

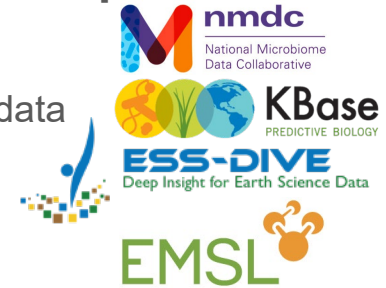


Accelerators



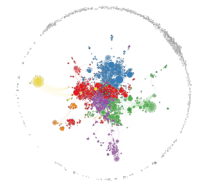
- **Build stronger connections between data resources and computing platforms**

- Submission portal to harmonize sample, environment, and process metadata
- Sample registration for all BER samples
- Collaborative User-centered Design



- **More analysis of the Data Citation Explorer results**

- Explore differences in data citations (e.g. methods vs background citation)
- Expose the Data Citation Explorer as a public-facing resource
- Add support for other identifiers (e.g. IGSNs or unique IDs from other fields)



- **Shared software and hardware infrastructure**

- Expand the number of sites JAWS supports
- JAMO infrastructure available to partners



Thank you

Staff who make this happen



Steve Chan
Advanced
Analysis GL

- Dani Cassol
- Jeff Froula
- Ed Kirton
- Angie Kollmer
- Ramani Kothadia
- Mario Melara
- Seung-Jin Sul
- Stephan Trong
- Nick Tyler



Data Portal



Steven Wilson
Advanced
Analysis GL

- Tatyana Smirnova
- Shijie Yao
- Alexandre Poliakov



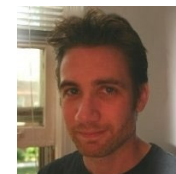
Georg Rath
Systems Infrastructure
Team Supervisor

- Harper Mann
- Matt Dunford
- Brian Yumae
- John White
- Karen Fernsler
- Wei Feinstein

The Data Citation Explorer



Neil Byers
Impact Analyst Data
Scientist



Chuck Parker
Computer Systems
Engineer



Chris Beecroft
Staff Data Scientist

- Ed Lee