

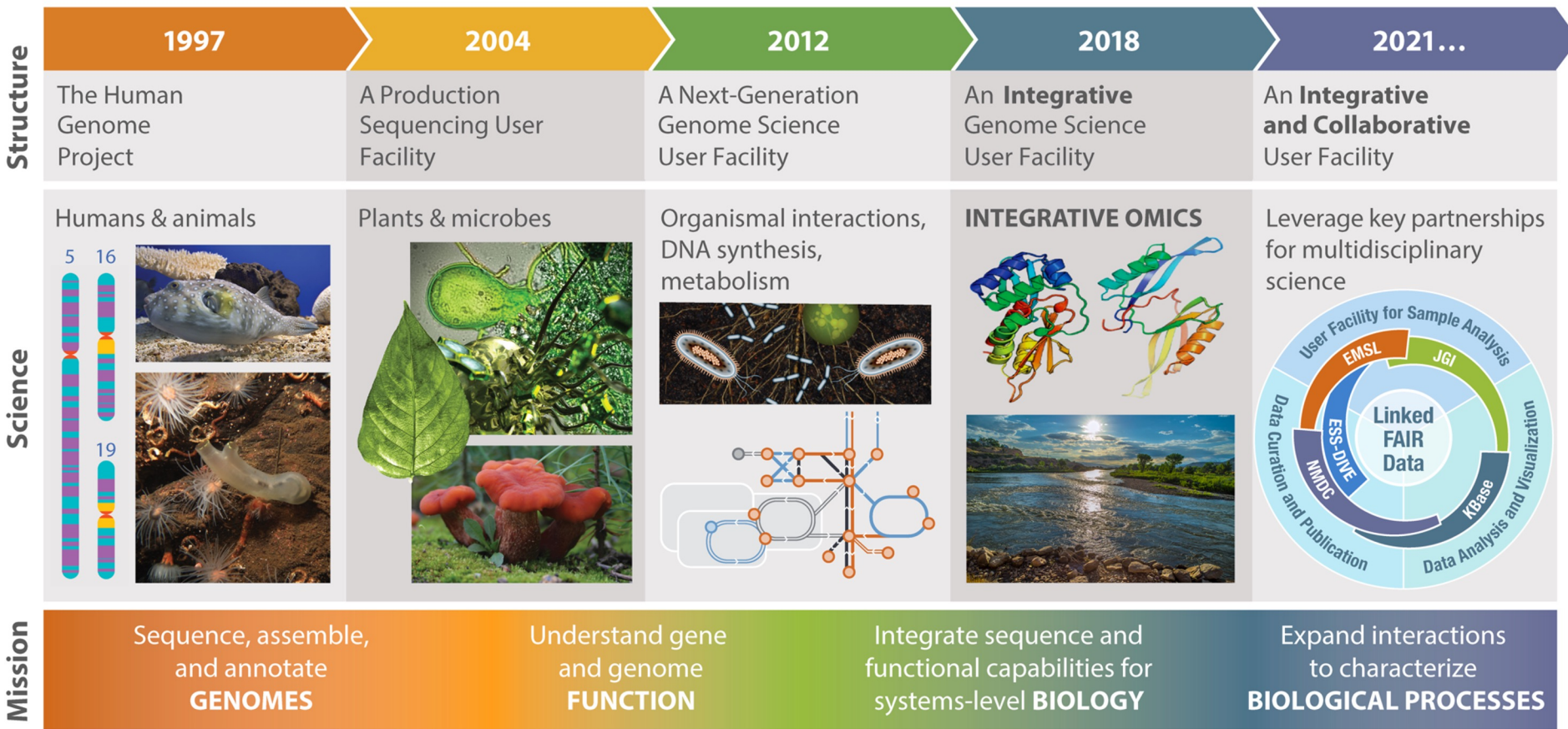


JGI-NERSC Partnership: lessons in data-intensive computing at scale

NERSC Celebrates 50 years
July 22, 2024



Continued Evolution into an Integrative and Collaborative User Facility



User Programs and Science Programs

SCIENTIFIC PEER REVIEW

40%
Community
Science
Program (CSP)

10%
Facilities
Integrating
Collaborations for
User Science
(FICUS)



30%
Bioenergy
Research Centers
(BRC)

10%
Director's Science

10%
Biological and
Environmental Research
Support Science (BERSS)

Plant



Fungal, Algal



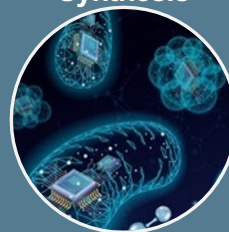
Metagenome



Microbial



DNA
Synthesis



Metabolomics



2^o Metabolites



2023: JGI by the Numbers



259 Staff · 25 New Hires
20 Grad Students · 18 Postdocs



Total Funding

DOE BER \$89.4M

2,373 Primary Users

22,262 Secondary Users

238 Publications



- 716.9 Terabases sequence generated
- 11 Megabases DNA synthesized
- 11.56K metabolomics analyses runs
- 161 proposals submitted
- 64 proposals approved

- Total files requested: 7.9M
- JGI Archive size grew to:
15.2 million file records
- 15.95 Petabytes (PB) of data



8.2K total podcast downloads

Engaged Reach:

X 24.5K

in 3.3K

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

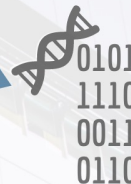
Data Generation and Reuse



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



Samples become data



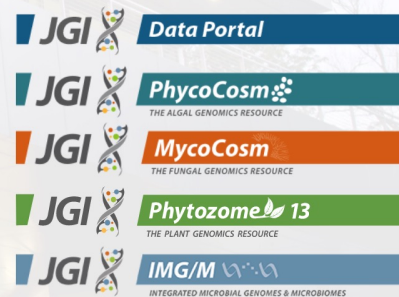
FY22 Downloads

4.1M Files

1.2PB Data



Primary and Secondary Users leverage data through JGI Flagship Science Gateways



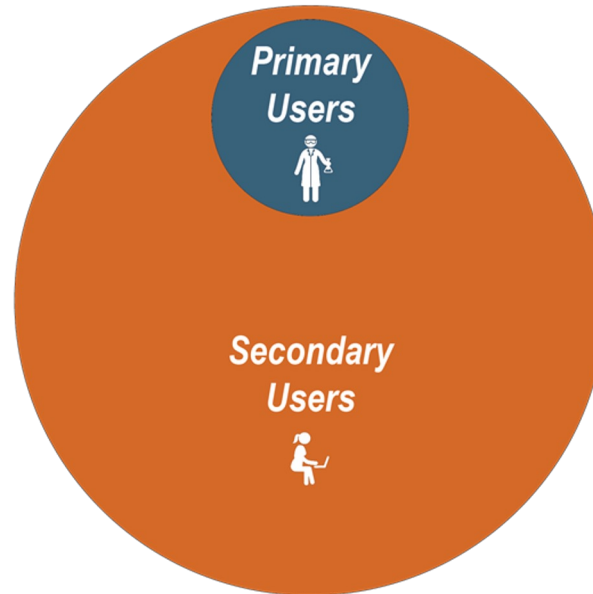
Primary and Secondary Users

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.



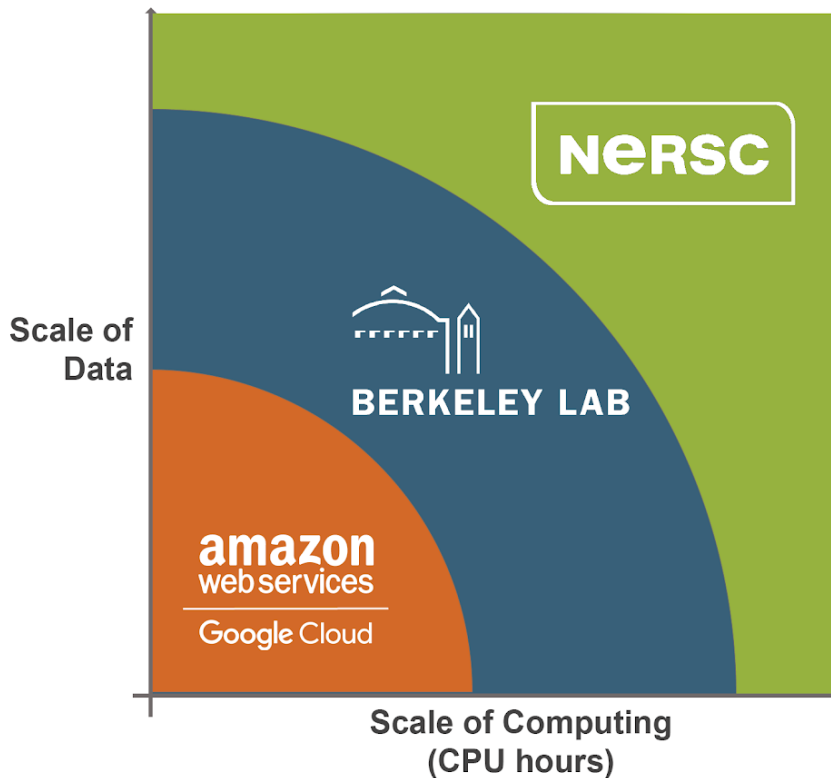
Example Outcomes

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

- 2010 – JGI data and computing was growing too fast for the Walnut Creek facility
- New hardware was redirected to NERSC facility in Oakland
- 2011 - JGI hardware unified to become the Genepool system
- 2012 – Consultants hired to help support the JGI use of shared storage and computing resources
- 2013 – Mendel deployed (consolidate Genepool and PDSF)
- 2019 – Mendel retired (moved to LBL IT), JGI had a cabinet of Cori
- 2023 – Cori retired, JGI stands up Dori at LBL IT



JGI's Computing Infrastructure Spectrum



Binning JGI Compute Infrastructure Requirements



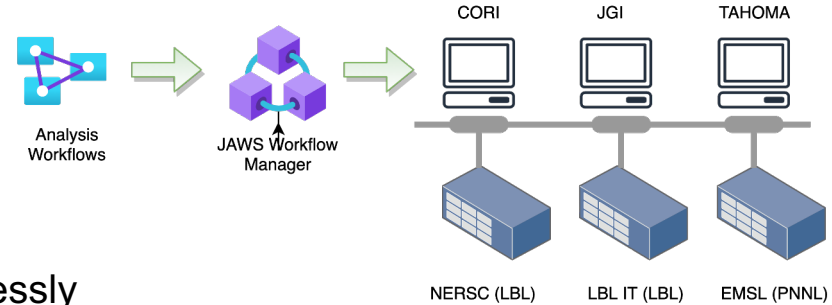
Prototyping,
exploratory
analysis, small-scale
production work



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

Unifying Workflow Execution Across JGI Resources

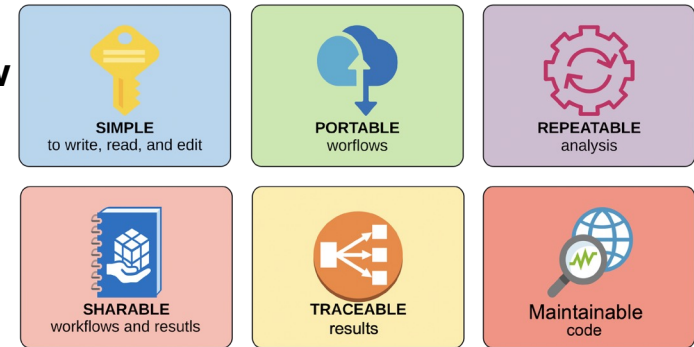
- Developed a workflow manager called **JGI Analysis Workflow Service (JAWS)** to run complex computational workflows with support for distributed computation across multiple HPC enabled sites.



- Provides a **user-friendly common interface** to seamlessly route jobs and data across multiple sites.

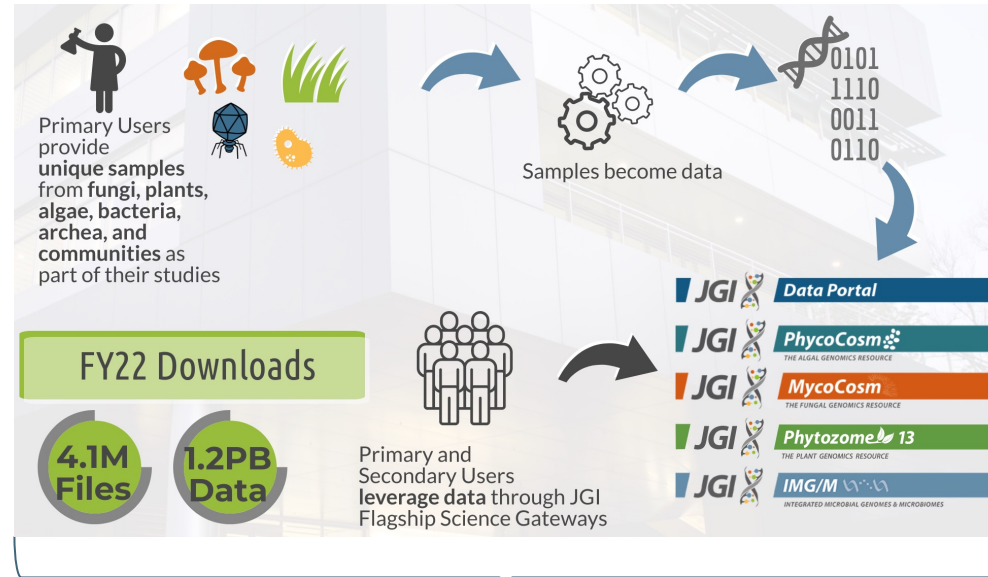
- Uses **Cromwell** to execute workflows in a common **Workflow Description Language (WDL)**, standardizing the workflow language.

- Improves the **reusability** and **robustness** of bioinformatics workflows in evolving and/or diverse high-performance computing (HPC) and cloud environments.



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**



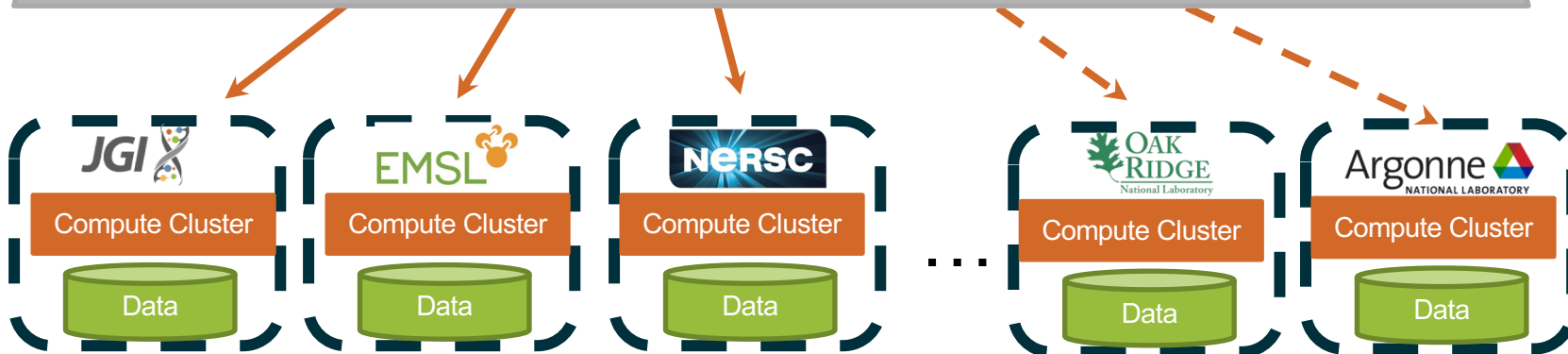
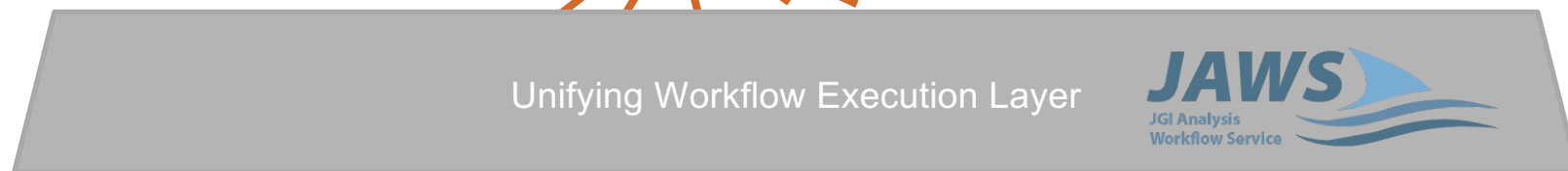
JAWS: Write Once, Run anywhere



JGI staff analyze JGI data on distributed resources



Containerized workflow, executable anywhere



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 notification banner indicates system maintenance on Sunday, April 11, 2021. The main heading reads 'Top-quality genomic data, open to all researchers' with a subtext 'Explore and download invertebrate genomes and metagenomes'. A search bar is present with the placeholder text 'Search by genome, metagenome, project, or ID'. Below the search bar are two columns: 'Search' and 'Download'. The 'Search' column describes the search capabilities and provides a link to 'More about search'. The 'Download' column describes bulk download options and provides 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.

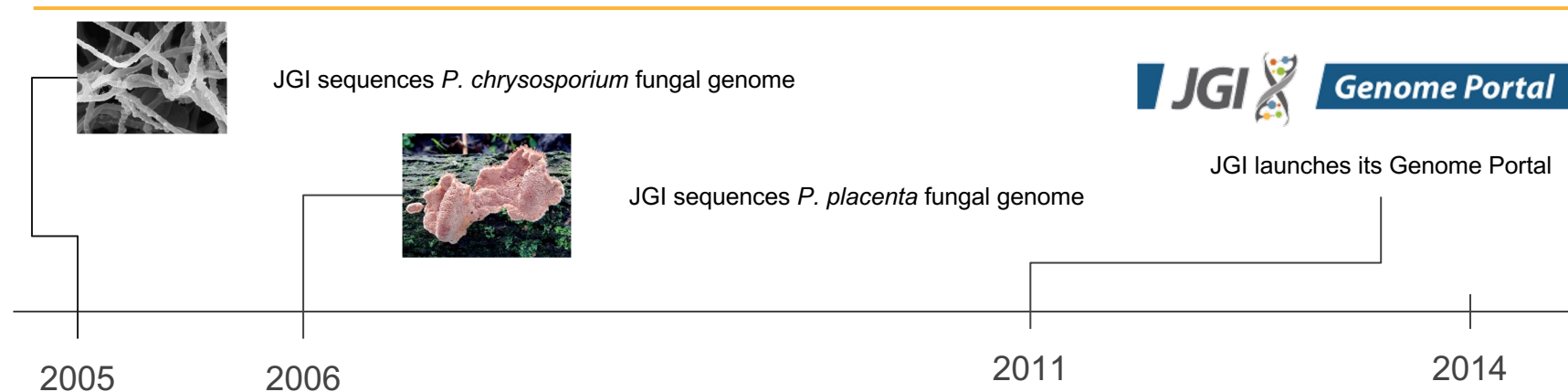
[Contact Us](#) | [Site Us](#) | [Accessibility/Section 508](#)
[Disclaimer](#) | [Credits](#)

© 1997-2020 The Regents of the University of California.
Genome Portal version: 8.16.44 content:1875004897_jgi_portal-web-1 Release Date: 23-Mar-2020 15:05:00.223 PST Current Data: 04-Mar-2020 08:57:24.987 PST



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

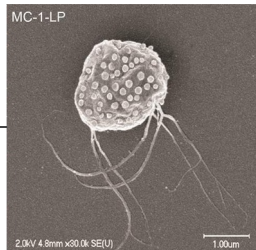
How does 'impact' begin?



Using data from the above genomes retrieved from JGI's Genome Portal, university and BP researchers [patent processes](#) to improve cellular sugar transportation for **production of biofuels** like ethanol.



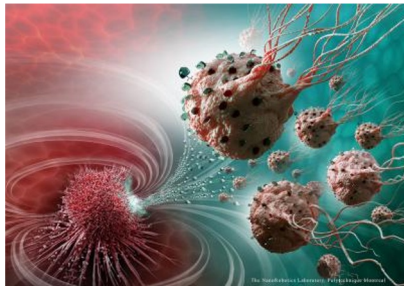
How does 'impact' begin?



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

2007

2016



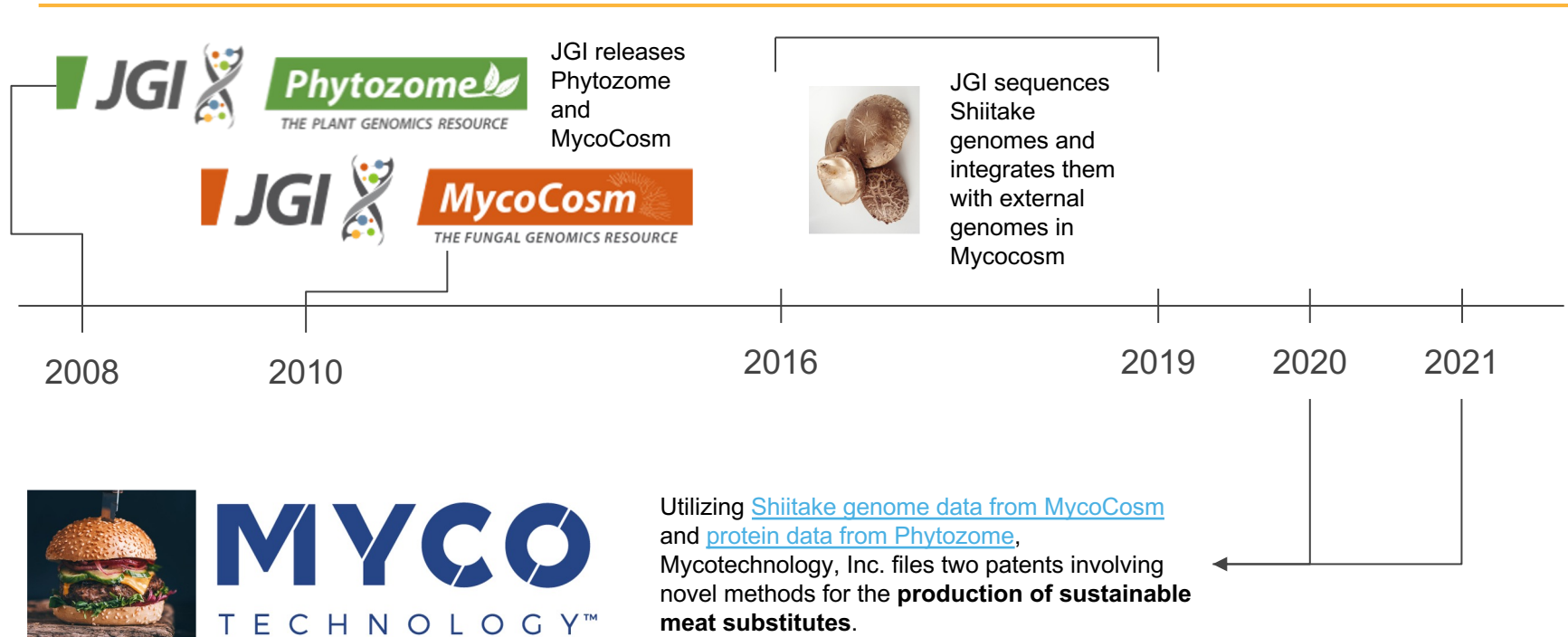
Aided by available genomic information, [researchers determine](#) that this bacteria is a very effective **medication delivery tool** for tumors in hard-to-reach areas of the brain.



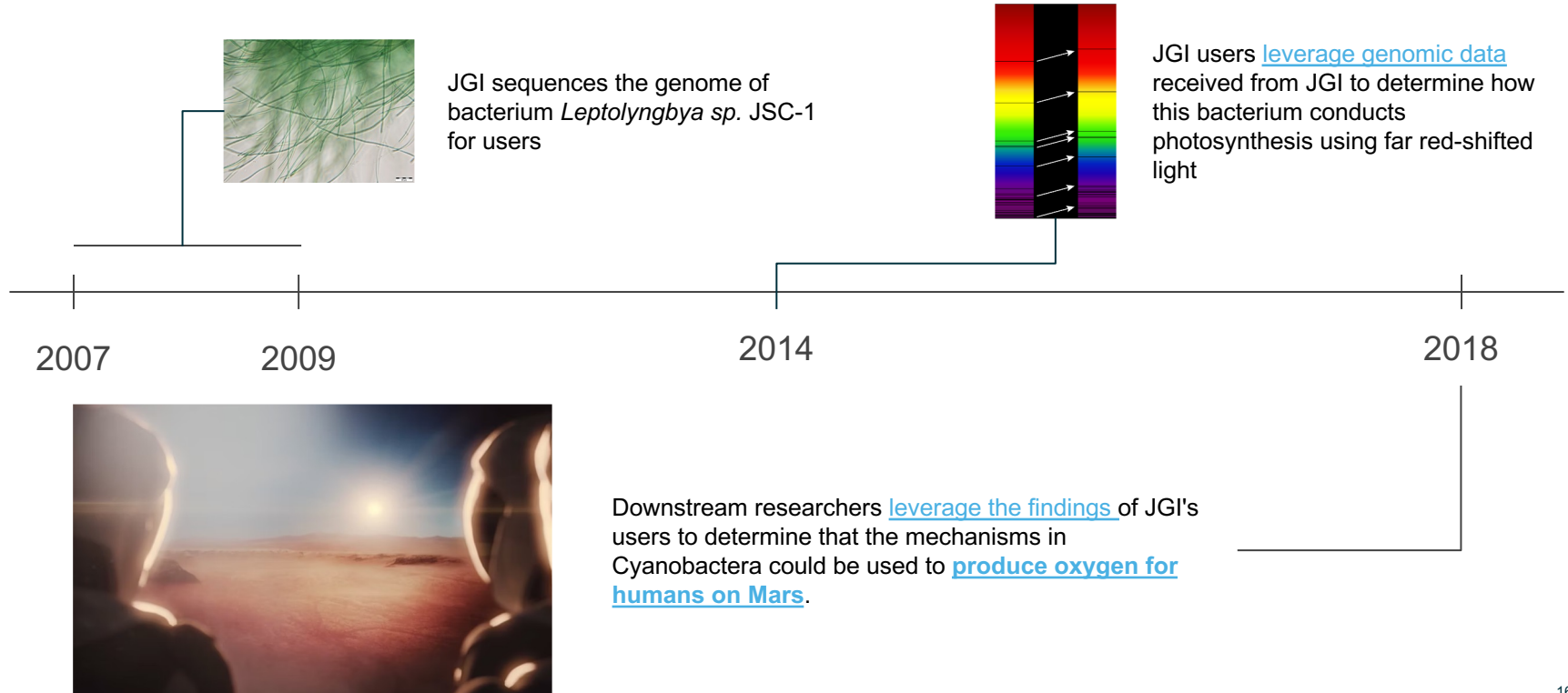
**POLYTECHNIQUE
MONTRÉAL**

UNIVERSITÉ
D'INGÉNIERIE

How does 'impact' begin?



How does 'impact' begin?



Highlight: The Megadata of Lake Mendota

- 3-part podcast arc released November-December 2023
- 20-year field samples from UW-Madison
- Highlighted JGI, NERSC, ExaBiome resources and capabilities harnessed to sequence and assemble 25Tb of metagenome data
- Episode interludes highlight that these are projects approved through JGI proposal calls



Host: Menaka Wilhelm



U.S. Department of Energy and 7 others



Trina McMahon



Host: Menaka Wilhelm

Season 4 Eps 6-8 highlight the combination of supercomputing and sequencing to understand environmental microbiomes!

Supercomputing, MetaHipMer and the Mega Dataset of Lake Mendota



David Buoy sits at the deepest part of Lake Mendota, where microbial samples have been collected for over two decades.

Sequencing technology has changed dramatically over the last [25-plus years](#) since the JGI's inception, making it possible for researchers to get a close look at more ecosystems and organisms than ever before. In 2006, the JGI produced 33 billion base pairs of sequence; by 2023, that number was almost 717 trillion. Last year, the JGI surpassed three Petabases of data sequenced – that's three-quadrillion base pairs of DNA sequence!