# The Genome Citation Service:

**Capturing JGI data citations for comprehensive impact assessment**

**Neil Byers**

**Scientific Impact Analyst**

- Part of the **Lawrence Berkeley National Laboratory**
- Located in **Berkeley, California**
- **~400 Staff**, including graduate students and postdocs
- **2,000+** Active PIs/collaborators per year
- **10,000+** Online users per year



**U.S. DEPARTMENT OF ENERGY** | Office of Science

**BERKELEY LAB**
Bringing Science Solutions to the World

**BioSciences**

**UNIVERSITY OF CALIFORNIA**

# JGI User Programs

# JGI Public Resources

## How do JGI user groups compare?



FY21 Publications
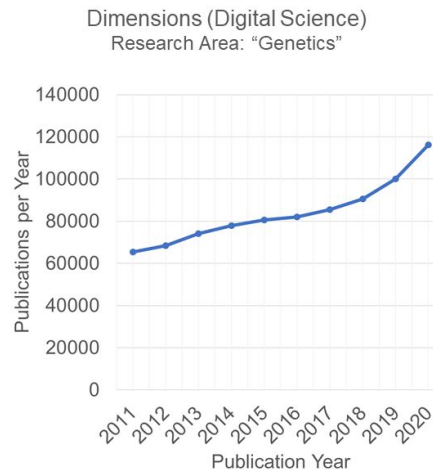
JGI
JOINT GENOME INSTITUTE

## Scalability: Publications, Data, & Metadata

- **Large amounts of data**
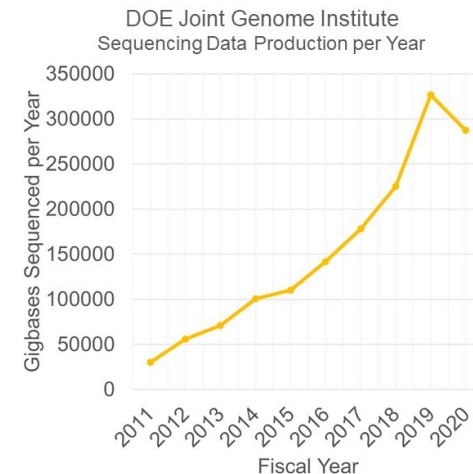- **Rich and diverse metadata**
- **Growing body of literature**

*– Too many citations –*
*– Using too many means –*
*– Of too many products –*

*To identify manually*



Dimensions (Digital Science)
Research Area: "Genetics"

Net increase in yearly output: 2011-2020: 77.53%
Publication total 2011-2020: 840,780



DOE Joint Genome Institute
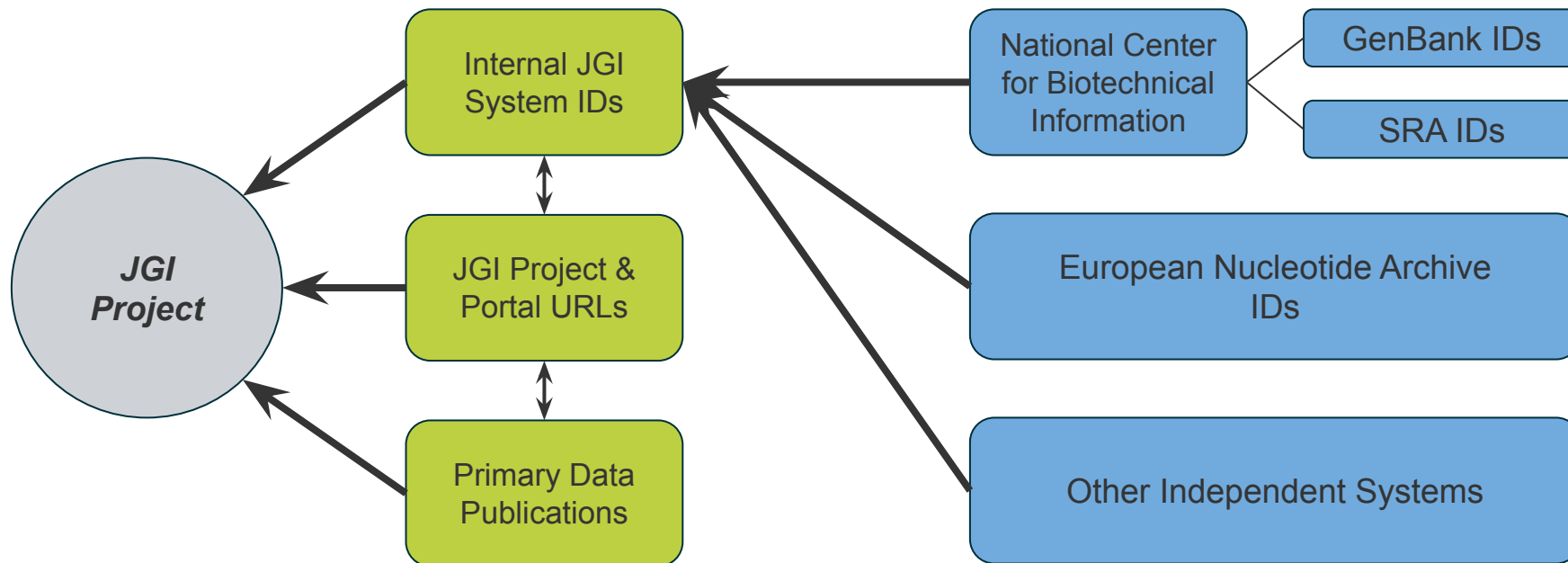Sequencing Data Production per Year

*FY2020 sequencing output affected by the COVID-19 pandemic.
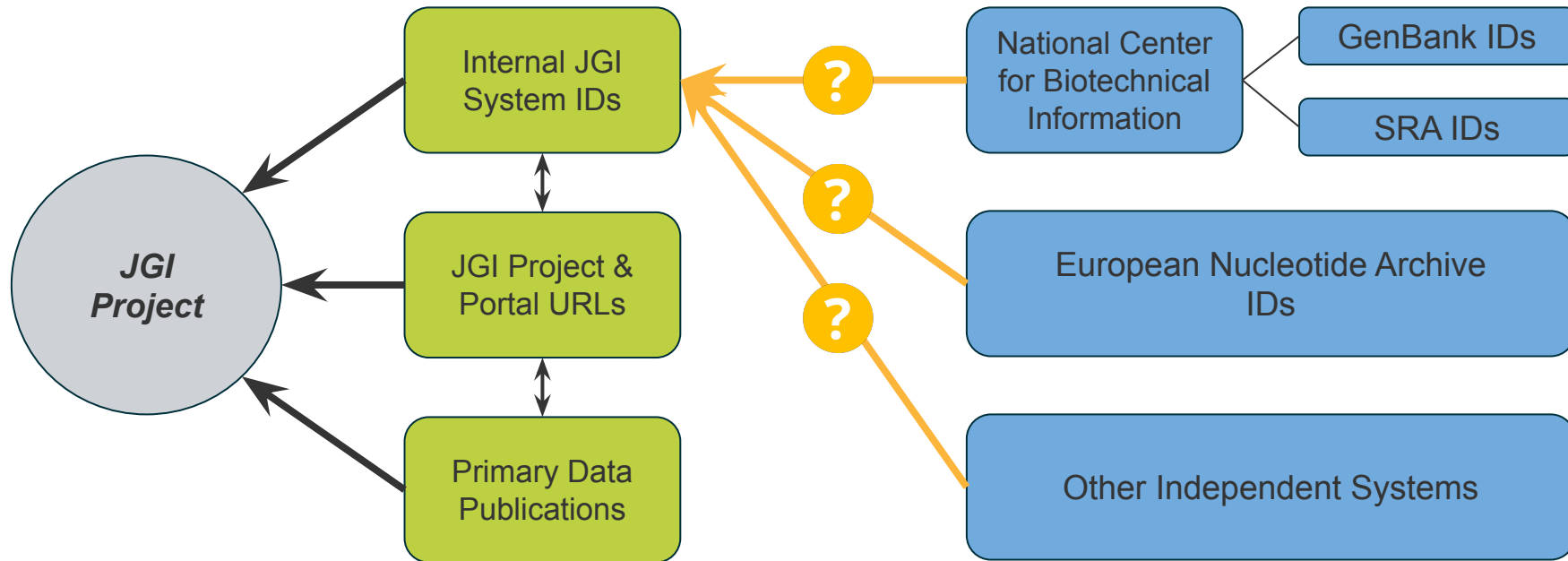
Raw data source:

Dimensions

## Incomplete Metadata

# The Problem

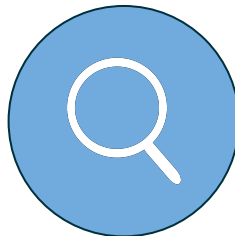## Is a cited identifier linked to JGI? What project(s) is it linked to?

## Automation: The Genome Citation Service (GCS)
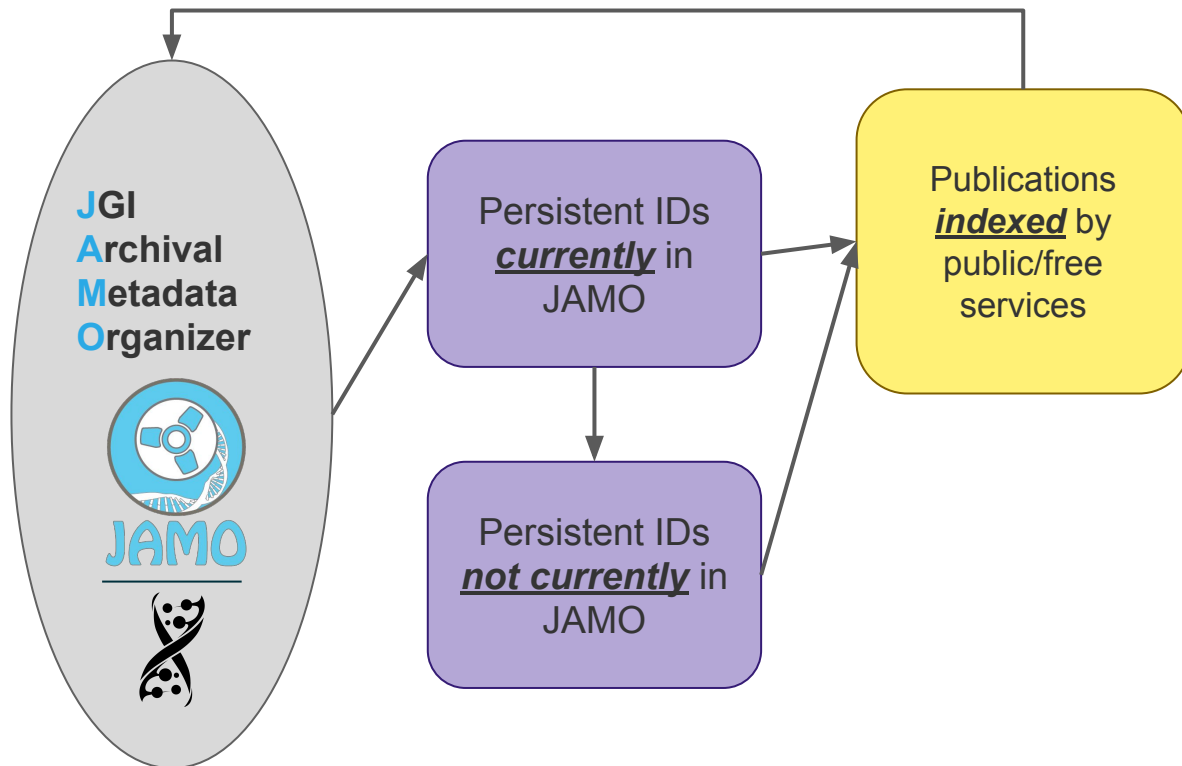
**Problem:**
*Incomplete Metadata*

**Solution:**
*Automatic, incremental discovery and traversal of additional linked data resources*

**Problem:**
*Scalability*

**Solution:**
*Automatic queries of public literature and cataloging of pathways back to JGI data resources*
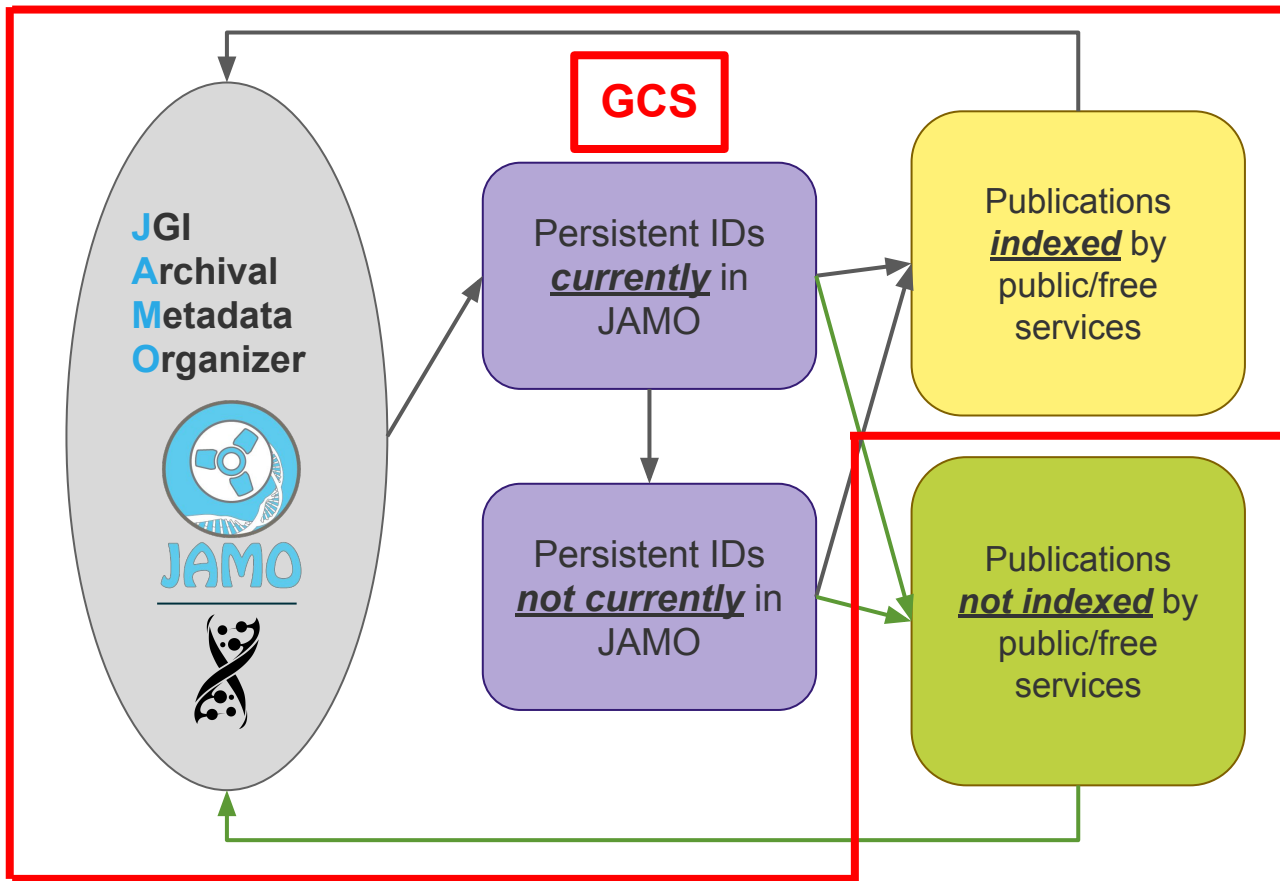
# The Genome Citation Service

## How well did the GCS perform?

- **Methodology**
  - 300 JAMO records (stratified, random)
  - Search public and subscription sources for linked publications
  - Manually evaluate hits
- **Results**
  - >98% Precision
  - 1234 total publications
  - 90% previously unidentified

Publications identified through the NamesforLife / JGI "Genome Citation Service" using two distinct publication databases and a stratified random sample of 300 JAMO records

1234 total
unique publications

NCBI
(PubMed / PubMed Central)
(Approx. 7 million
searchable
full-text articles)

207 publications

369 publications

658 publications

Dimensions
(Approx. 75 million
searchable
full-text articles)

National Center for Biotechnology Information
NCBI

>98% validity rate
for both sources

**How does JGI benefit from the Genome Citation Service?**

1. **Thousands of JGI data citations now potentially identifiable**
2. **Little or no manual effort required**

**How can we use this information for impact assessment?**

- **Community Analysis**
  – Coauthor and JGI user networks
- **Topic Analysis**
  – Data use trends
  – DOE goal alignment
- **New researcher metrics**
  – Equitable credit attribution

# Community Analysis

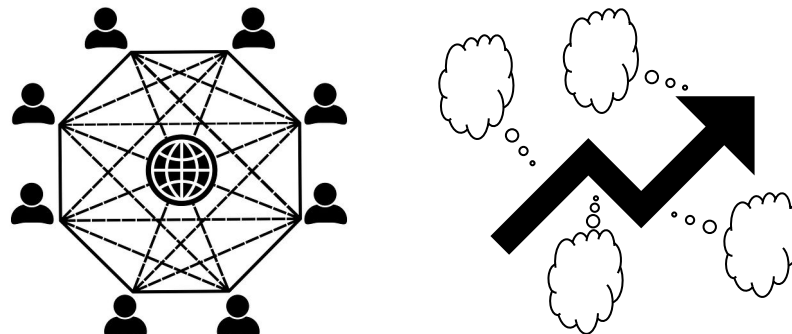## Which PI/Collaborator groups are producing data of scientific interest?



**JGI Proposal Contributors**

**Group A**
- X proposals with
- Y datasets receiving
- Z total data citations
- Consists largely of DOE Bioenergy Research Center (BRC) personnel

**Group B**
- X proposals with
- Y datasets receiving
- Z total data citations
- Contains many contributors to JGI's 1000 Fungal Genomes Project

# Community Analysis
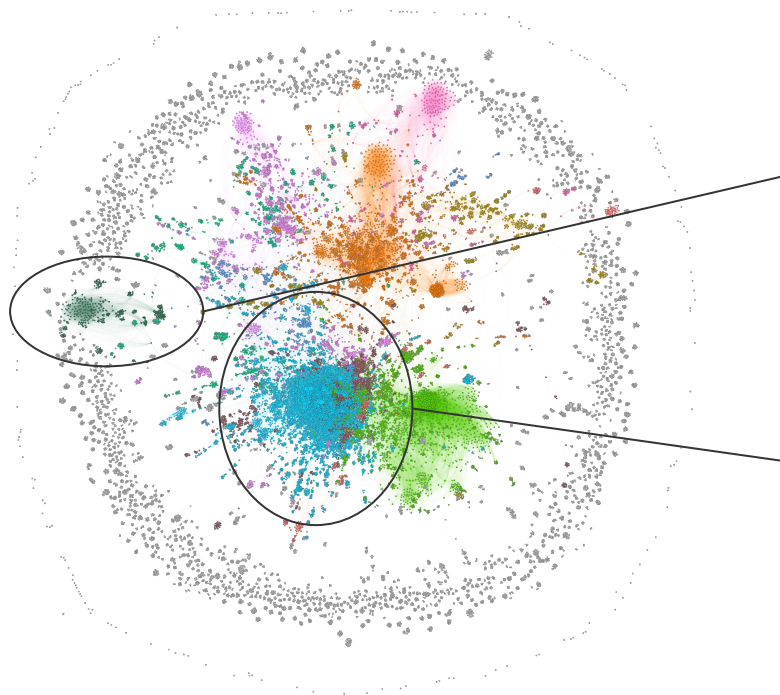
**Can we identify distinct communities among users of public JGI resources?**



**Blue Group*** vs. **Teal Group**
- Total Authors:
  - **5015**
  - **378**
- Top Institutions:
  - **DOE Labs, American & European Universities**
  - **European universities**
- Top concepts:
  - **Evolution, Pathogens, Biomass degradation**
  - **Bioinformatics, Ontology**
- Publications
  - **1774 publications (~18 average citations)**
  - **123 publications (~30 average citations)**

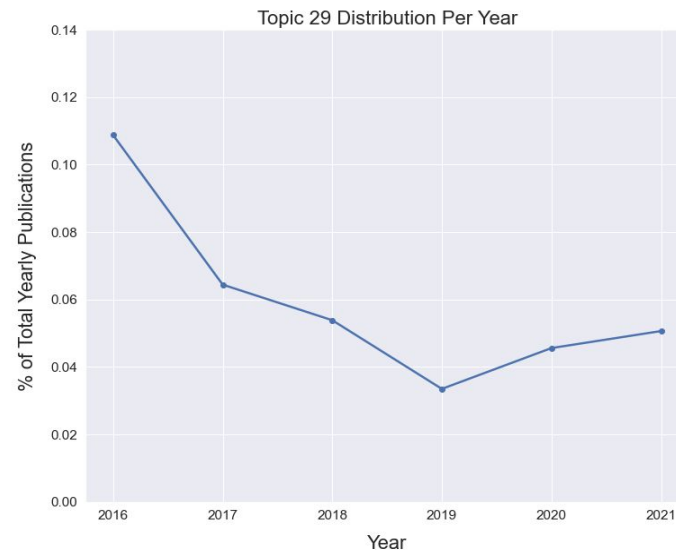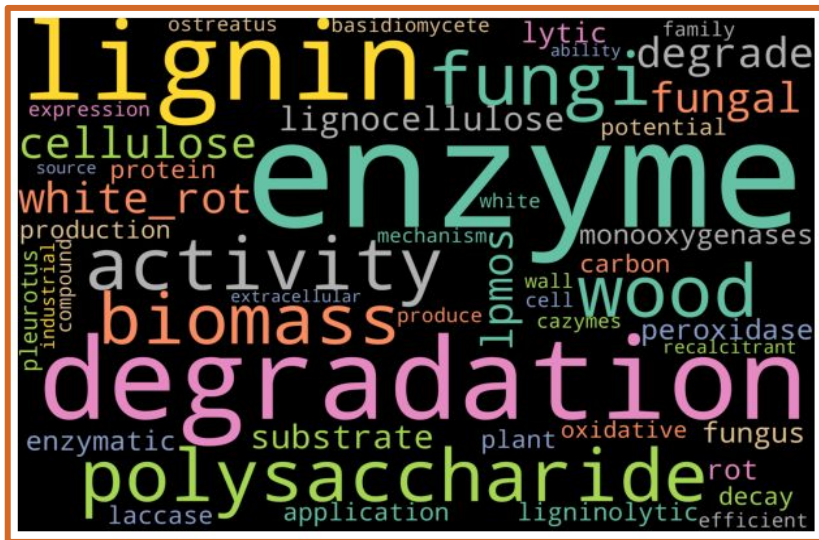**JGI Users**
**(mock subset)**

Raw data source:

Dimensions

*Top authors in Blue Group are JGI PIs & Collaborators

## What 'topics' are trending amongst our user base?



Topic 13 Distribution Per Year

This 'topic' seems to be on the rise among citing publications. It consists of 235 publications from our test set, and is increasing in total publication share over time.

# Topic & Concept Analysis

## What 'topics' are trending amongst our user base?



Topic 29 Distribution Per Year

This 'topic' seems to be losing prominence among citing publications. It consists of 341 publications from our test set, and is decreasing in total publication share over time.
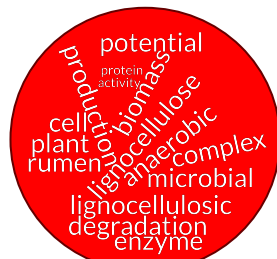
# Topic & Concept Analysis

## Are there counter-intuitive use cases for JGI data?

**JGI Proposal:**
**"Transcriptomic Characterization of Anaerobic Gut Fungi", 2014**
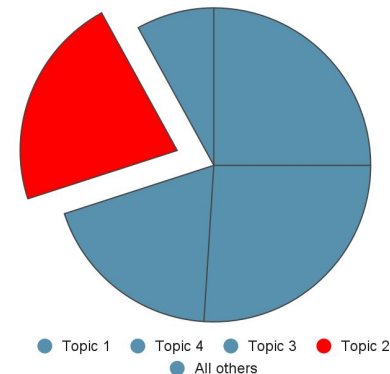549 Downstream Citations

**Proposal Text:**

*"The overall goal of our efforts is to develop the tools to engineer anaerobic gut fungi as novel platform organisms for **biofuel production** from **lignocellulosic biomass**."*

Topic 2
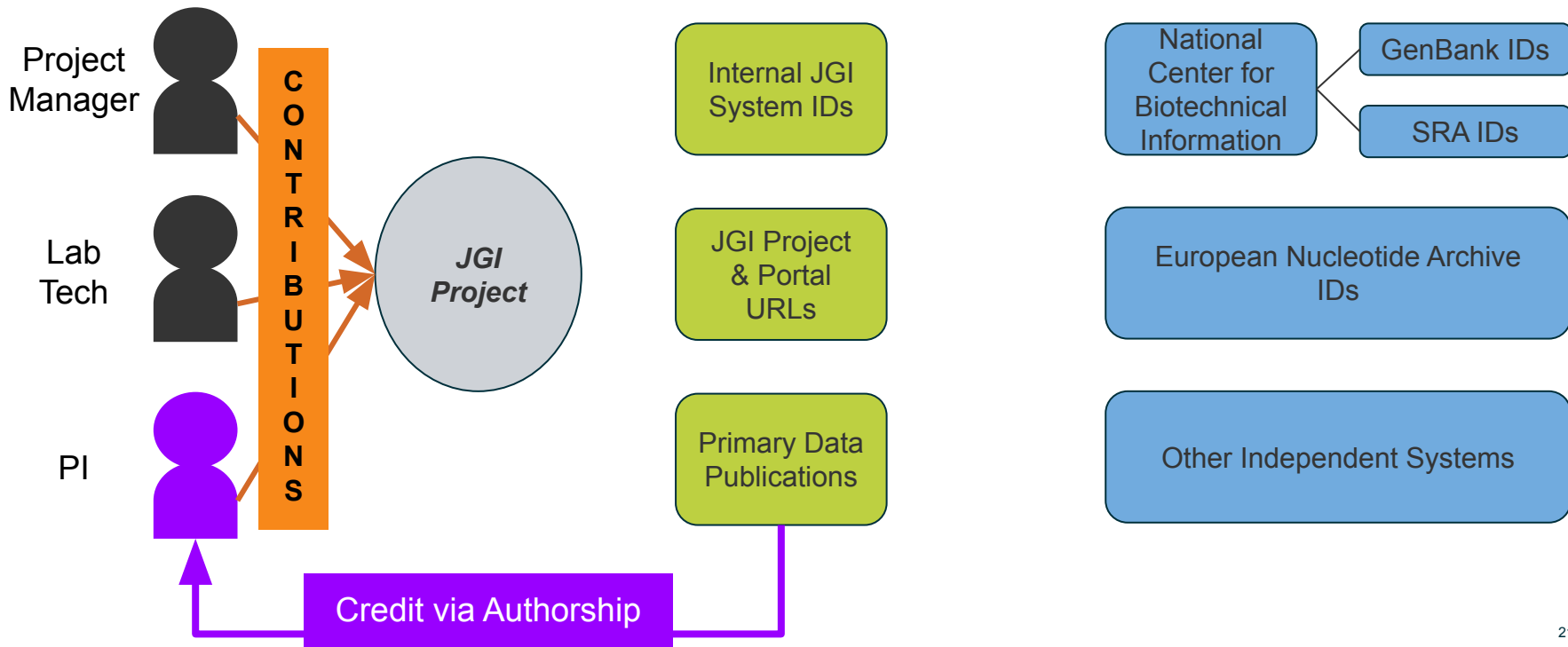


- 137 (25%) Publications
- 2.3k citations
- Top Journals: *Biotechnology for Biofuels; Bioresource Technology*

Citing Literature Share per Topic

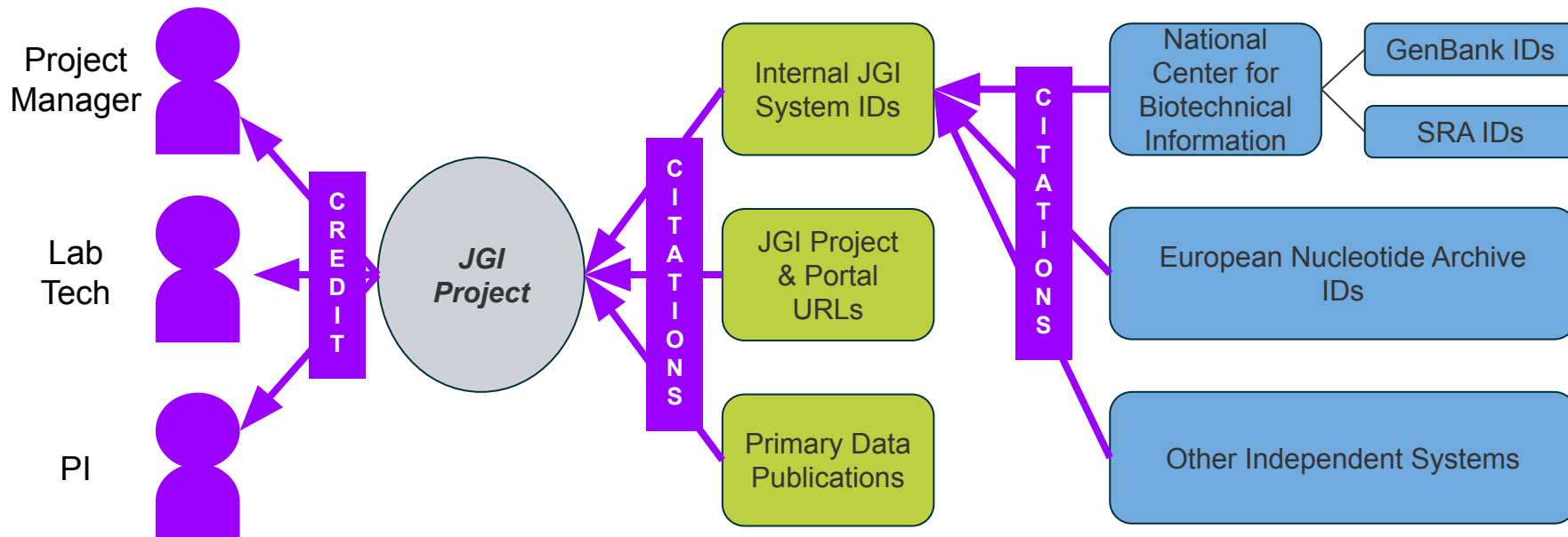

Topic 1 ● Topic 4 ● Topic 3 ● Topic 2
● All others

# Equitable Contributor Metrics

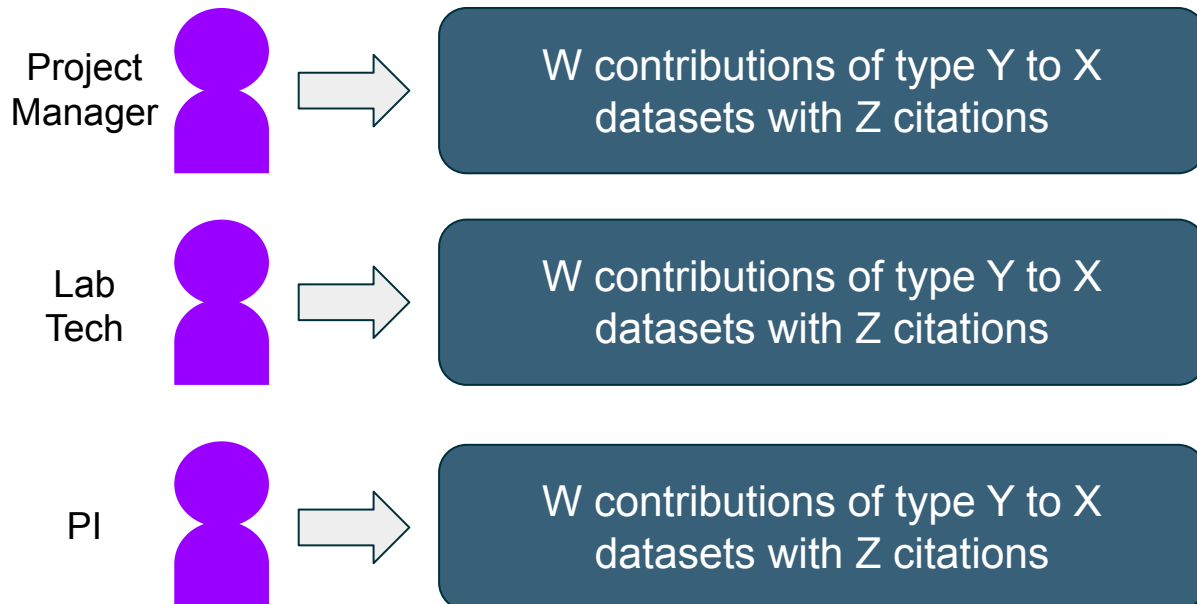**Authorship Model - Often arbitrary and uneven credit attribution**

# Equitable Contributor Metrics

**Data Citation Model - Equitable attribution of credit**

# Equitable Contributor Metrics

**Data Citation Model - Equitable attribution of credit**

Project Manager → W contributions of type Y to X datasets with Z citations

Lab Tech → W contributions of type Y to X datasets with Z citations

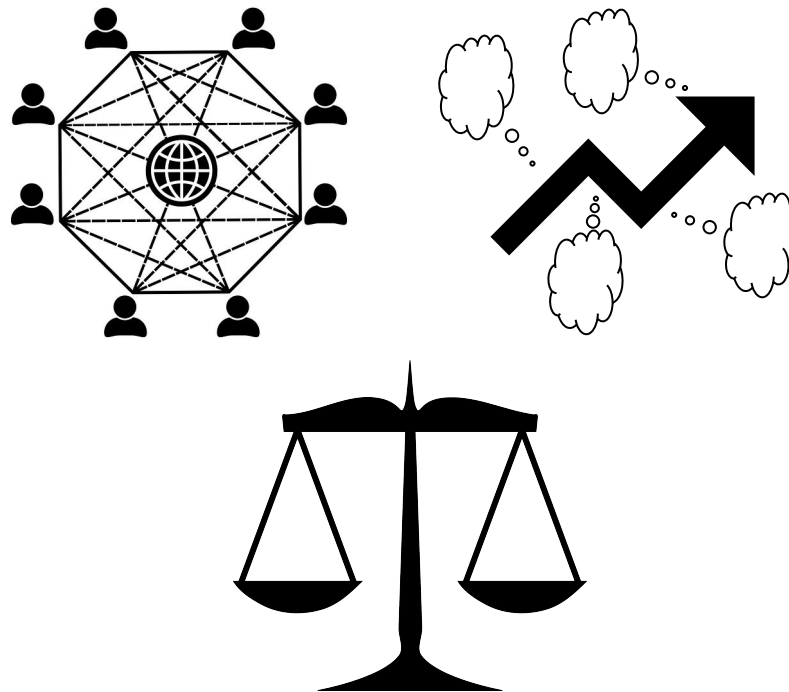PI → W contributions of type Y to X datasets with Z citations

Example:

Hood & Sutherland,

"The data-index: An author-level metric that values impactful data and incentivizes data sharing"

*Ecology and Evolution* Oct. 2021

- **GCS Goals**
  - Link metadata and capture data citations at scale
- **Results**
  - Many JGI data citations identified with high validity rates
- **Impact Implications**
  - More comprehensive picture of JGI community impact
  - Equitable contributor metrics

# Thank you!

## Genome Citation Service Team:

JGI | JOINT GENOME INSTITUTE

**Chris Beecroft**
**Hugh Salamon**
**Kjiersten Fagnan**
**TBK Reddy**
**Neil Byers**

*NamesforLife*
*Bringing meaning to life...*  |  **Charles Parker**

MICHIGAN STATE UNIVERSITY  *NamesforLife* *Bringing meaning to life...*  |  **George Garrity**



Kjiersten Fagnan

Chris Beecroft

Neil Byers

Hugh Salamon

TBK Reddy

Charles Parker

George Garrity

**Contact:** Neil Byers, npbyers@lbl.gov