# KBase
## PREDICTIVE BIOLOGY

# DOE Systems Biology Knowledgebase

Adam Arkin, KBase PI
2023 Spring BERAC

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

Office of Biological and Environmental Research

*COMMUNITY-DRIVEN*

*PREDICTIVE BIOLOGY*

# What is KBase?

KBase enables users to analyze, share, and collaborate using data and tools designed to help build increasingly realistic models for biological function.

# A Virtual Cycle for  Genotype->Phenotype Prediction from Genes to Ecosystems



Image Credit: LLNL Microbes Persist SFA

**KEY**
- Data or Model
- Analysis Step
- Sampling
- Analyze

**Bin and assemble isolates**

Isolate genomes

**Design and organize laboratory phenotyping and engineering**

Metagenomics → Assembly + Binning → Metagenome-assembled genomes (MAGs)

**Bin and assemble metagenomes**

**Predictive modeling of community/ecological dynamicss**

Reactive transport simulations

Quantitative community models

Quantitative data fitting

Qualitative community models

**Design and organize field sampling and isolation  and metadata**

**Organize multi-omic data sampling, metadata, and public data integration**

Multi-omics data

**Data- and model-driven probabilistic annotation of genes and genomes**

Annotating core metabolism

Energy pathway reconstruction

Annotation model building

Annotated genomes and models

**Predictive modeling of organismal function/dynamics**

Mechanistic microbiome analysis ← Predictive models ← Phenotype prediction & fitting

KBase
PREDICTIVE BIOLOGY
Systems Biology Knowledgebase

# KBase ensures all data & analyses are FAIR *and* credited

Team member shares an analysis

'Narrative' of analysis has everything
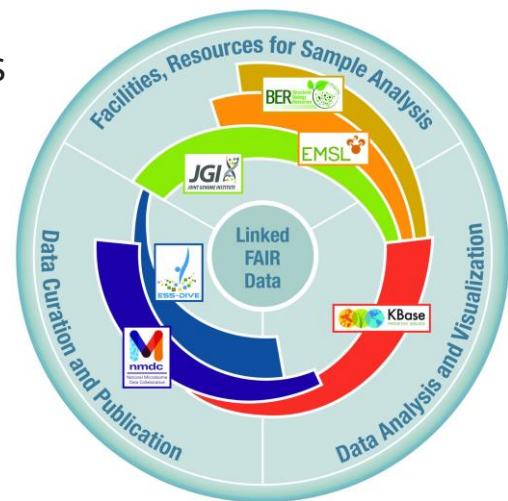
Every object fully 'provenanced'

# Why do people use KBase?

- Combine different types of data from many resources
- Create unique, complex, reproducible analyses
- Contribute new knowledge back to the research community

# But it takes a (BER) village to enable FAIR data at scale

- Data (and all research products) must be FAIR across programs
- FAIR data must have provenance across the data life cycle
- BER program can create the infrastructure, but culture change requires *trust*
  - Where the data came from
  - How it was processed
  - Ability to explore quality
  - Ability to see impact and effect of combined power
- People *want* to be FAIR, but it takes support to do well, and they *must get credit for it*
- **Designed KBase to connect everything *and* give credit**



U.S. DOE. 2021. Biological Systems Science Division Strategic Plan, DOE/SC-0205. U.S. Department of Energy Office of Science
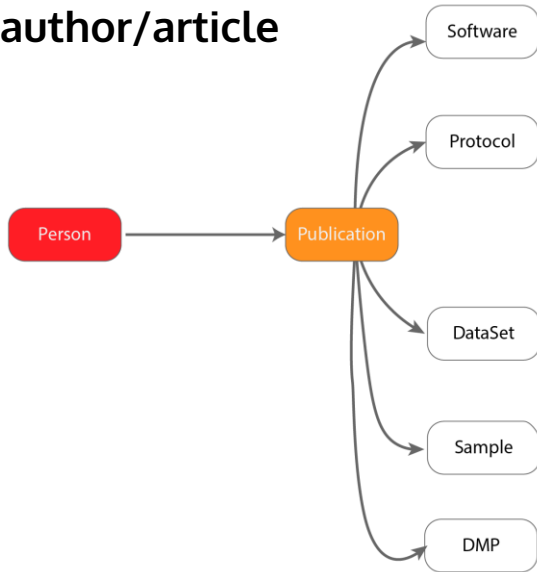
U.S. DEPARTMENT OF
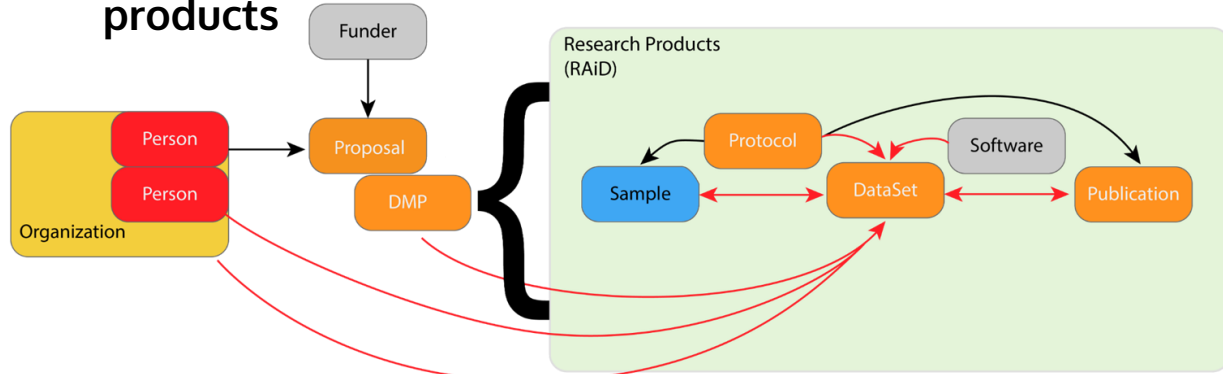**ENERGY**
Office of Science

**KBase**
PREDICTIVE BIOLOGY
**DOE Systems Biology Knowledgebase**

# Goal: FAIR, trackable research outputs



**Now: Focus on author/article**

**Future: Credit to author/funder for all research products**

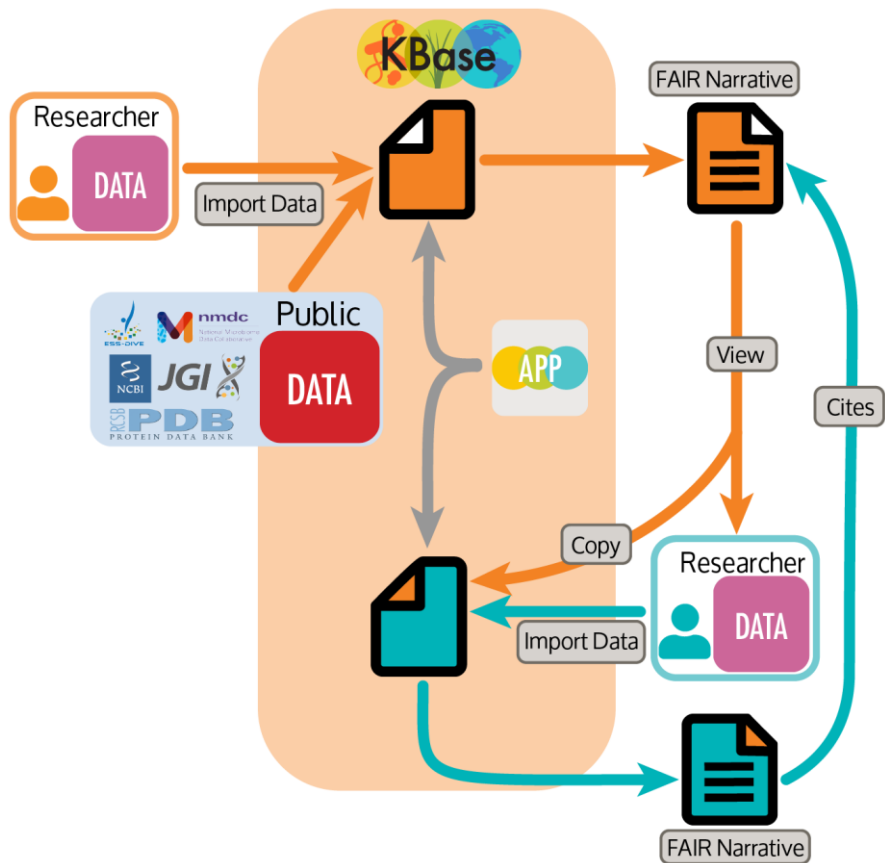*Red arrows: KBase connections between data sets and other PIDs*

*Glossary of Persistent Identifiers (PIDs)*
- DOI - digital object identifier
- RAiD - research activity identifier
- IGSN - international geo/general sample number
- ORCiD - open researcher and contributor identifier
- ROR - research organization registry

Legend

**Sample Sets Organize Complex Data**

# Tracking impact requires provenance and PIDs



- Provenance has always been integral to KBase's design
- We are not alone! KBase is positioned to effectively leverage external resources to ensure BER research products are FAIR, visible, and trackable.

- FAIR data is best reused when truly understood - knowing what it is, who generated it, and how it was generated.

# Moving FAIR data requires alignment between platforms

JGI-KBase Co-development "Data Transfer Service"

**Data transfer**
-file
-type
-md5



**KBase Credit Engine Metadata Fields**
- citation (person, title, date, version, etc)
- license
- funding
- related IDs
    - DOIs (proposal, protocols, data, etc)
    - Sample ID (IGSN, GOLD, etc)
    - ORCIDs
    - etc.

SCHEMATRON 9000

**KBase Data Object**

Document of Provenance

This is a genuine KBase object, imported into the narrative at 03:57:19 PST on Thursday October 20th in the year of our lord 2022. Day after day, day after day, We stuck, nor breath nor motion; As idle as a painted ship Upon a painted ocean. Water, water, every where, And all the boards did shrink; Water, water, every where, Nor any drop to drink. The very deep did rot: O Christ! That ever this should be! Yea, slimy things did crawl with legs Upon the slimy sea. About, about, in reel and rout The death-fires danced at night; The water, like a witch's oils, Burnt green, and blue and white. And some in dreams assured were Of the Spirit that plagued us so; Nine fathom deep he had followed us From the land of mist and snow. And every tongue, through utter drought, Was withered at the root; We could not speak, no more than if We had been choked with soot. Ah! well a-day! what evil looks Had I from old and young! Instead of the cross, the Albatross About my neck was hung.
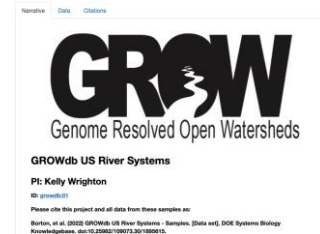
Credit Metadata

There passed a weary time. Each throat Was parched, and glazed each eye. A weary time! a weary time! How glazed each weary eye, When looking westward, I beheld A something in the sky. At first it seemed a little speck, And then it seemed a mist; It moved and moved, and took at last A certain shape,

*FAIR Narrative*

Index by schema.org

Mikayla Borton   Kelly Wrighton

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.

Hypothesis: Spatial and geochemical features influence river microbiomes



Como Creek

Mississippi River

Results: Stream order, geochemistry, and temperature correlate with community structure

# Data integration - provenance, credit



ESS-DIVE: sample metadata and biochemistry

JGI: metaG and metaT
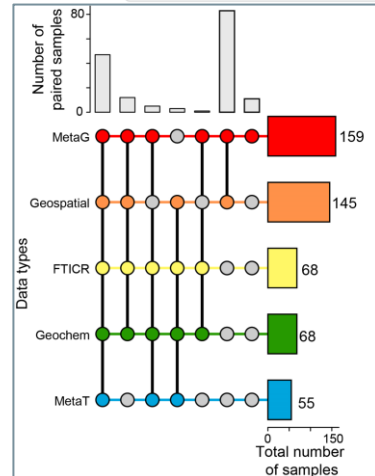
EMSL: metaB

NMDC: study information, standardized analyses

KBase: combines everything to generate community models

# GROW Science

Metabolic models of metagenome-assembled genomes (MAGs) explore energy biosynthesis across major river phyla

# Publishing to share and get credit

*Samples, data, analyses - all FAIR and free!*

Citation: Borton M, *et al.* (2022) GROWdb US River Systems - Samples. [Data set]. DOE Systems Biology Knowledgebase. doi:10.25982/109073.30/1895615.



https://kbase.us/n/109073/41

**GROWdb US River Systems - Samples**

GROWdb data account, Boris Sadkhin, Christopher Henry, Elisha WC, Filipe Alexandre Wang Liu, Janaka Edirisinghe, José Pedro Lopes Faria, Mikayla Borton, Mikayla A Borton, Shane Canon, Zach Crockett

Generated November 3, 2022

Narrative | Data | Citations

**GROWdb US River Systems**

**PI: Kelly Wrighton**

**ID:** growdb:01

Please cite this project and all data from these samples as:

Borton, et al. (2022) GROWdb US River Systems - Samples. [Data set]. DOE Systems Biology Knowledgebase. doi:10.25982/109073.30/1895615.

---

**Google Scholar** — GROWdb

Articles — About 35 results (0.05 sec)

Any time
Since 2023
Since 2022
Since 2019
Custom range...

**GROWdb** US River Systems-Samples
M Borton, K Wrighton, B Sadkhin, C Henry... - 2022 - osti.gov
... an ecosystem-specific publicly available genome database (**GROWdb**) that will be a resource for ... Dataset Acknowledgement **GROWdb** contains data from various research campaigns, ...
☆ Save  99 Cite  »

https://orcid.org/
0000-0001-8037-4253

Name
**Mikayla Borton**

ORCID

Activities — Collapse all

∨ Works (6)   ☰ Sort

**GROWdb US River Systems - Samples**
KBase
2022 | Data set
DOI: 10.25982/109073.30/1895615
CONTRIBUTORS: Mikayla Borton; Kelly Wrighton; Boris Sadkhin; Christopher Henry; Elisha Wood-Charlson; Filipe Liu; Janaka Edirisinghe; Jose Faria; Shane Cannon;   Show more detail
**Source**: Mikayla Borton

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

**KBase** — PREDICTIVE BIOLOGY — DOE Systems Biology Knowledgebase

# Current and Future Efforts

Completing the KBase Credit Engine Inside and Out

Maturing our connection to the publishing infrastructure including 'pushing' MRA to ASM journals.

Working with large DOE team on a more universal data transfer system with long term goals of:

- Unified authorization validation
- Universal query with key terms across the resources
- Common identifiers, credit, etc.
- Common Data Transfer System: to uniformly move data among systems and validate on both ends. (near term)

# The KBase Team