# WormBase

## Todd Harris, PhD

todd@wormbase.org          @tharris
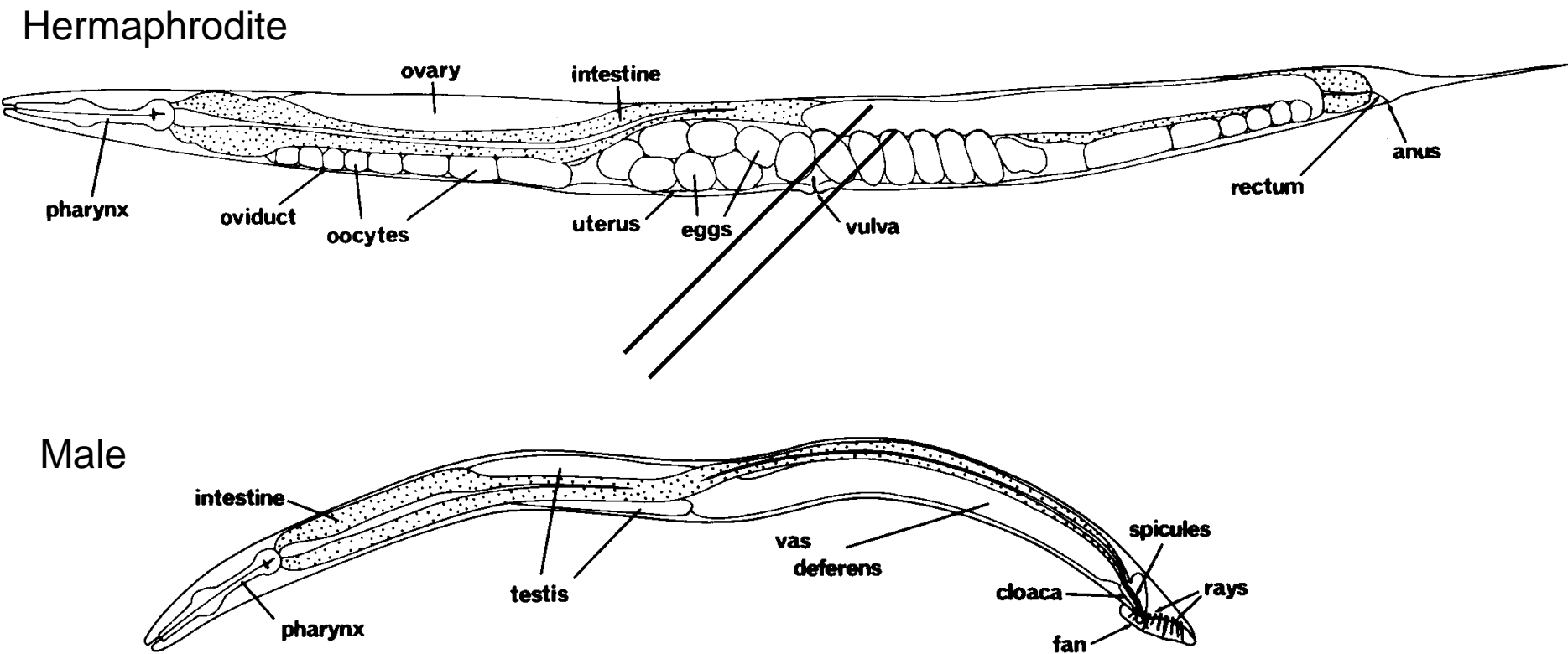
CBPSS Mini Symposium

# Mission

Provide the biomedical research community with **accurate**, **current**, and **accessible** information on the genetics, genomics, and biology of the model system *Caenorhabditis elegans* and related nematodes.

# *C. elegans* in 30 seconds

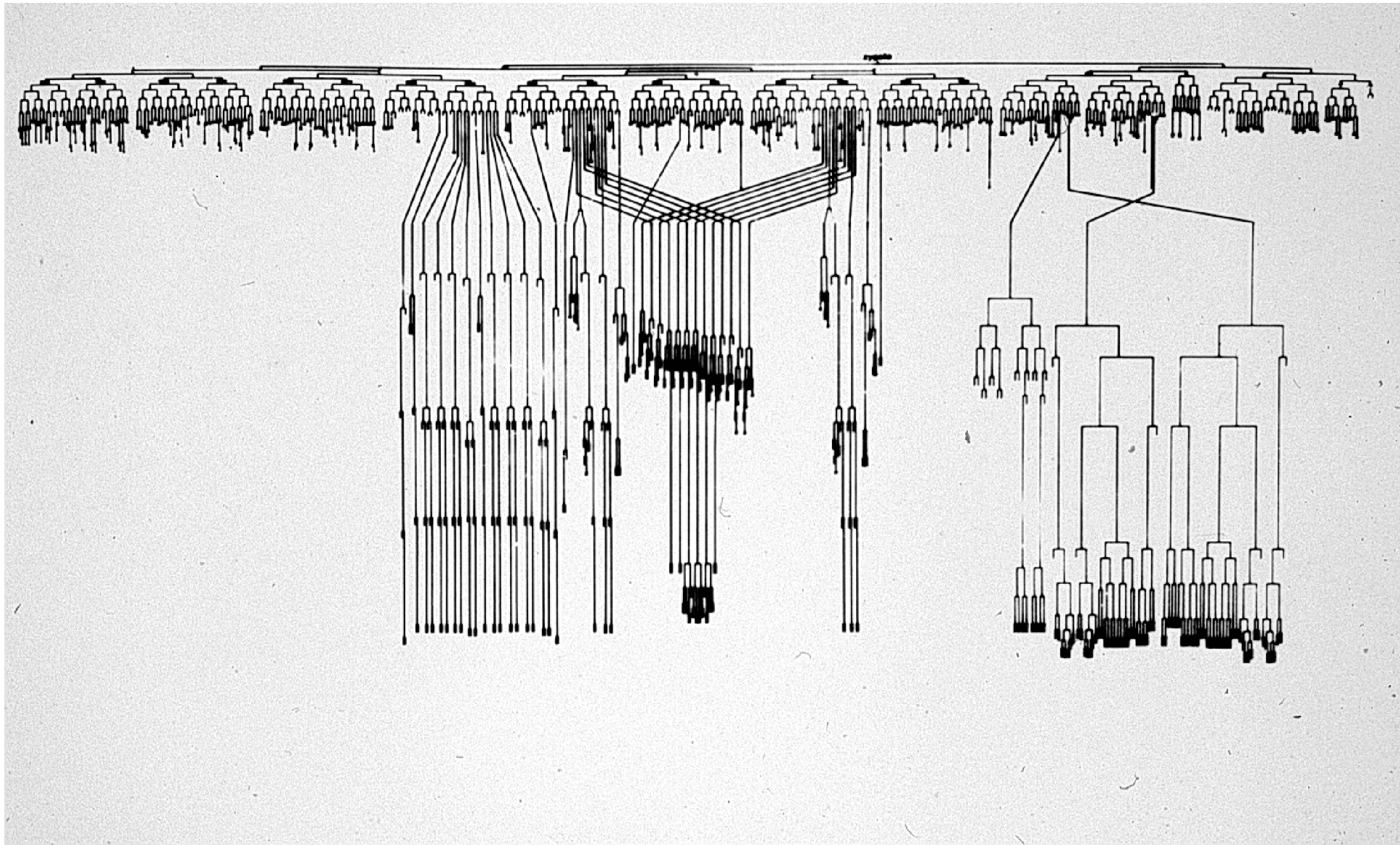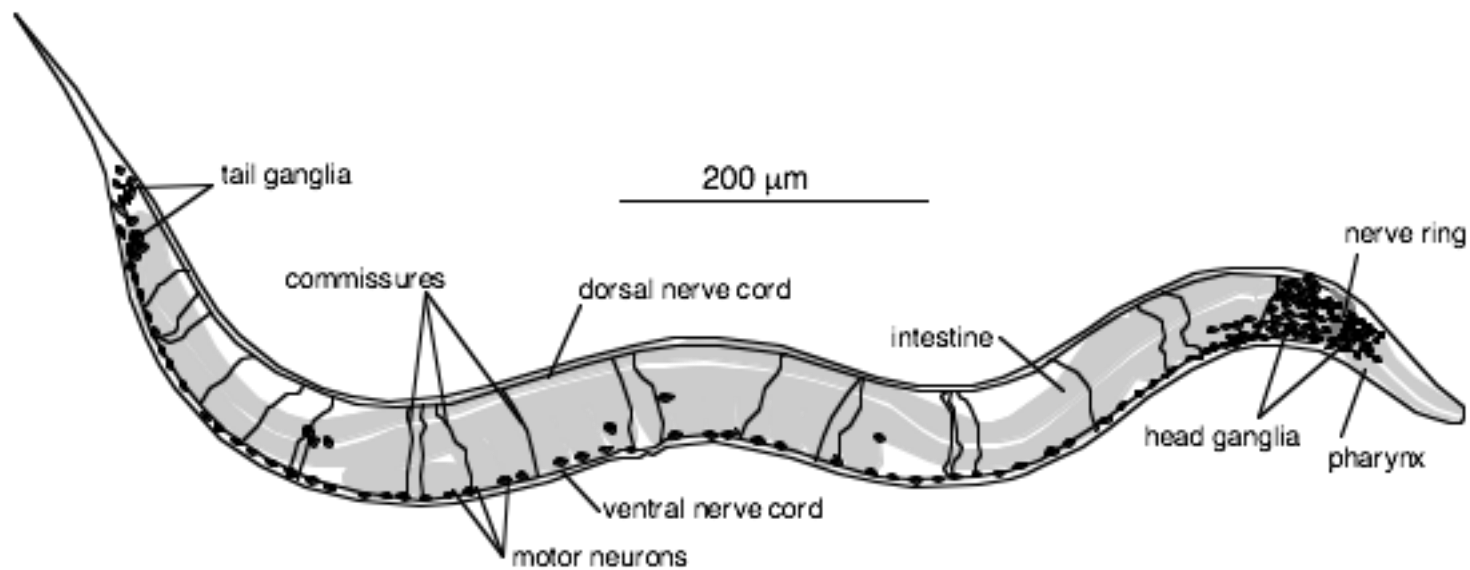Relatively simple organism, advanced genetic system.



Hermaphrodite

ovary — intestine — pharynx — oviduct — oocytes — uterus — eggs — vulva — rectum — anus

Male

intestine — pharynx — testis — vas deferens — cloaca — fan — spicules — rays

1mM

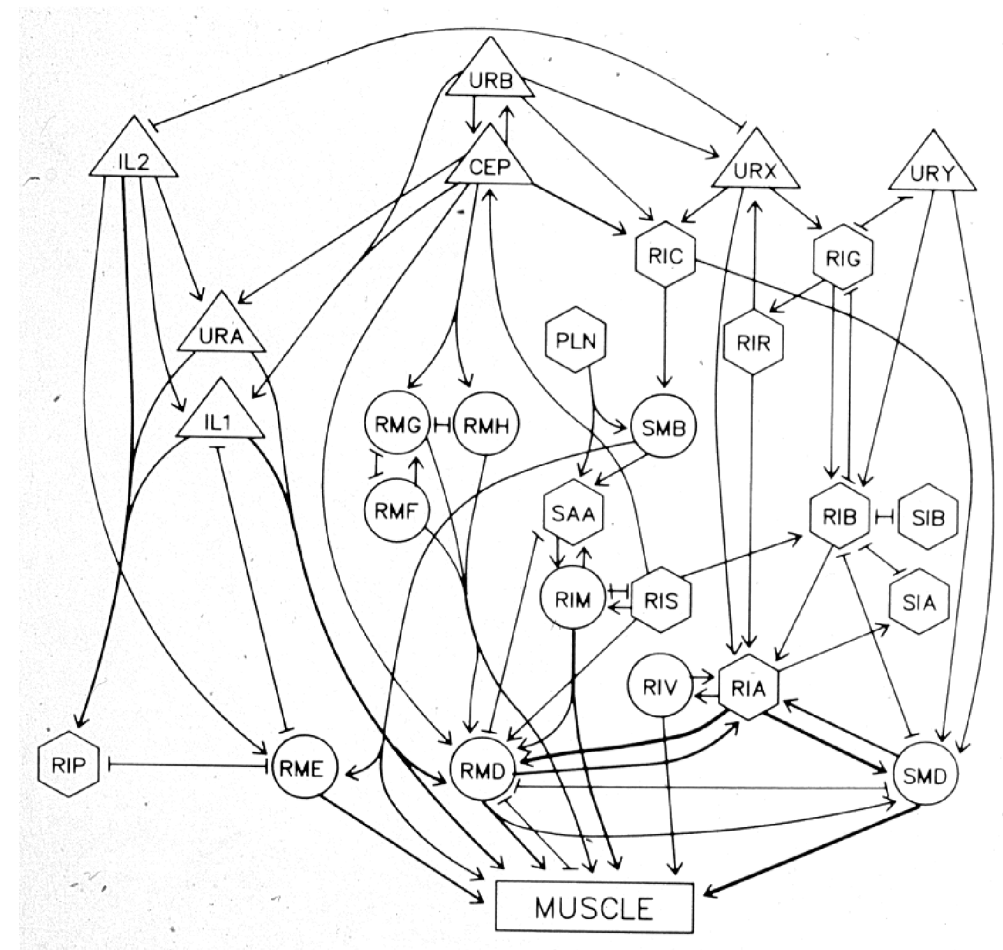# *C. elegans* in 30 seconds

Invariant lineage

# *C. elegans* in 30 seconds
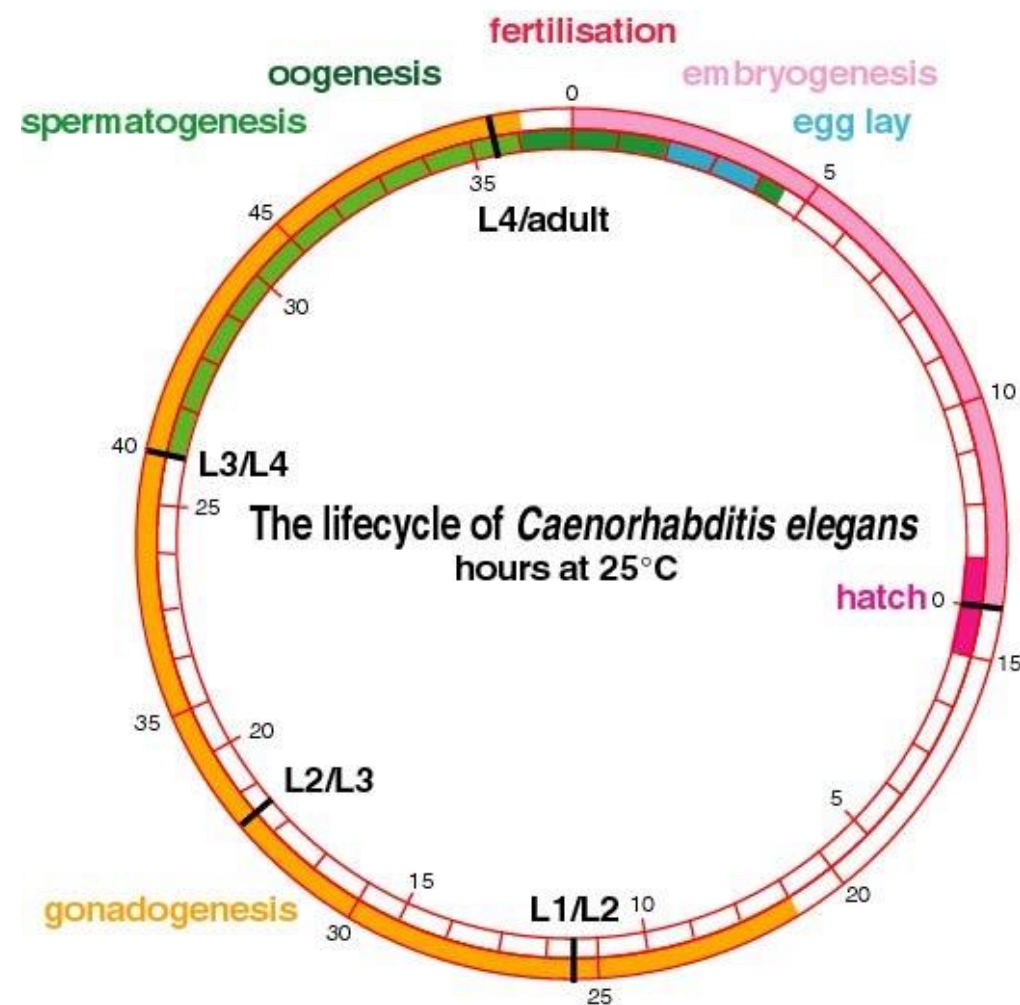
Simple nervous system



302 neurons

Described connectivity

# *C. elegans* in 30 seconds

Rapid generation time
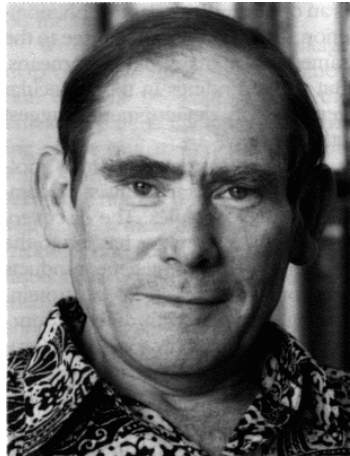
A frozen *C. elegans* library

# *C. elegans* in 30 seconds

100 MBp Genome

~20K genes

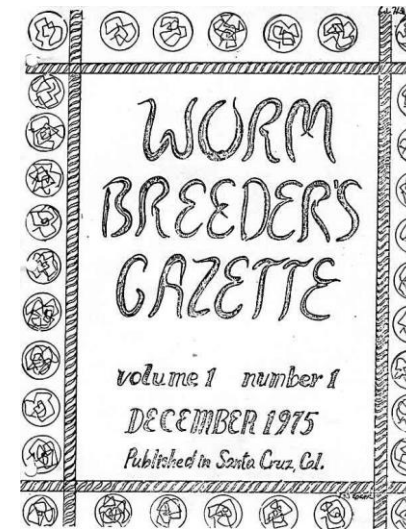1998 (!)

# A tradition of Open Science
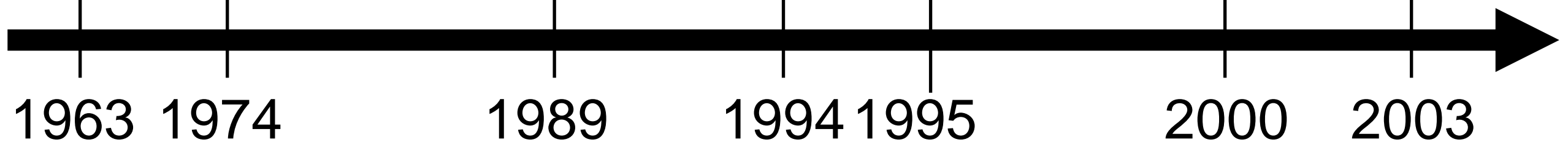


Brenner's
Letters

1st genetic screen
published

AceDB
development
begins

BioNet
www
gopher

Gazette

WormBase

**WormBook**
THE ONLINE REVIEW OF *C. elegans* BIOLOGY

1963  1974          1989          1994 1995              2000      2003

# The WormBase Consortium

# User Community

## Registered *C. elegans* laboratories

| Country | Labs |
|---|---|
| United States | 594 |
| Canada | 62 |
| United Kingdom | 60 |
| Japan | 58 |
| Germany | 48 |
| France | 31 |
| China | 28 |
| Spain | 20 |
| Switzerland | 20 |
| The Netherlands | 16 |

**1106 laboratories**
**53 countries**
**3000 researchers**

# User Community

Biomedical researchers studying aging, neurobiology, cancer, etc.

185 **countries**

37K **unique users/month**

5.5M **page views / month**

**wormbase.org**

# Contents & Features

28 Species
Genomes
Genes
Orthology / Homology / Paralogy
Comparative Genomics
Strains / Antibodies / Oligos
Expression
Lineage & Connectivity
Authors & Publications
Labs

Reports
Genome Browsers
Alignment Tools
Query Tools
APIs
Data Mining Platforms
Social Features
FTP
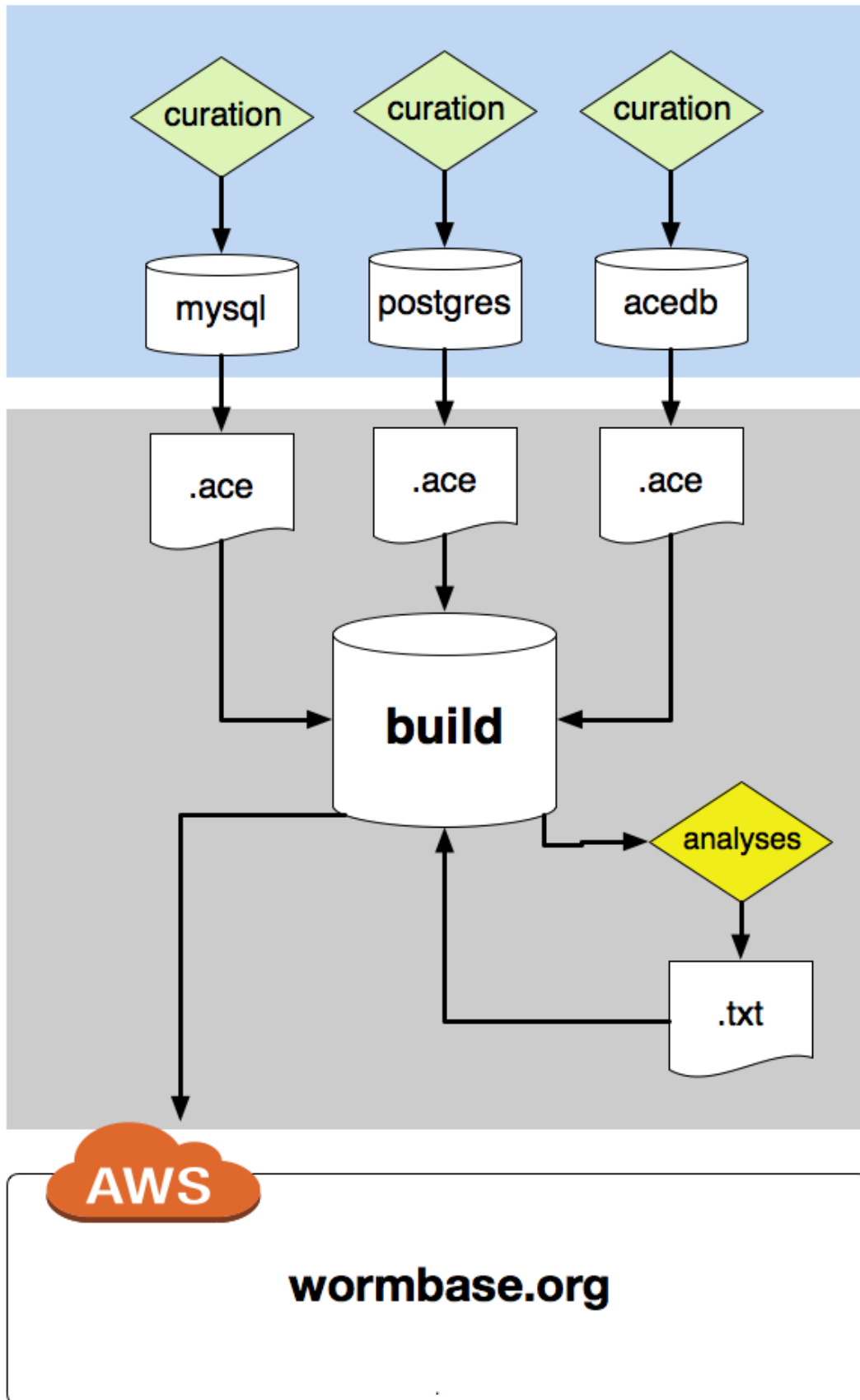Forums, Wikis, Blogs

# Workflow



**1. Curation**

**2. Integration & analysis**

**3. Presentation**

# Curation Goals

1. Extract data from the scientific literature.

2. Develop standards to structure data.

3. Facilitate new insights by making prose observations computable.

# Curated Sources

Scientific literature (~30K papers)

User submissions

Genomic sequences (gene models)

3rd party datasets

# Early Realizations

**Curation is hard and time-consuming!**

Requires automation.

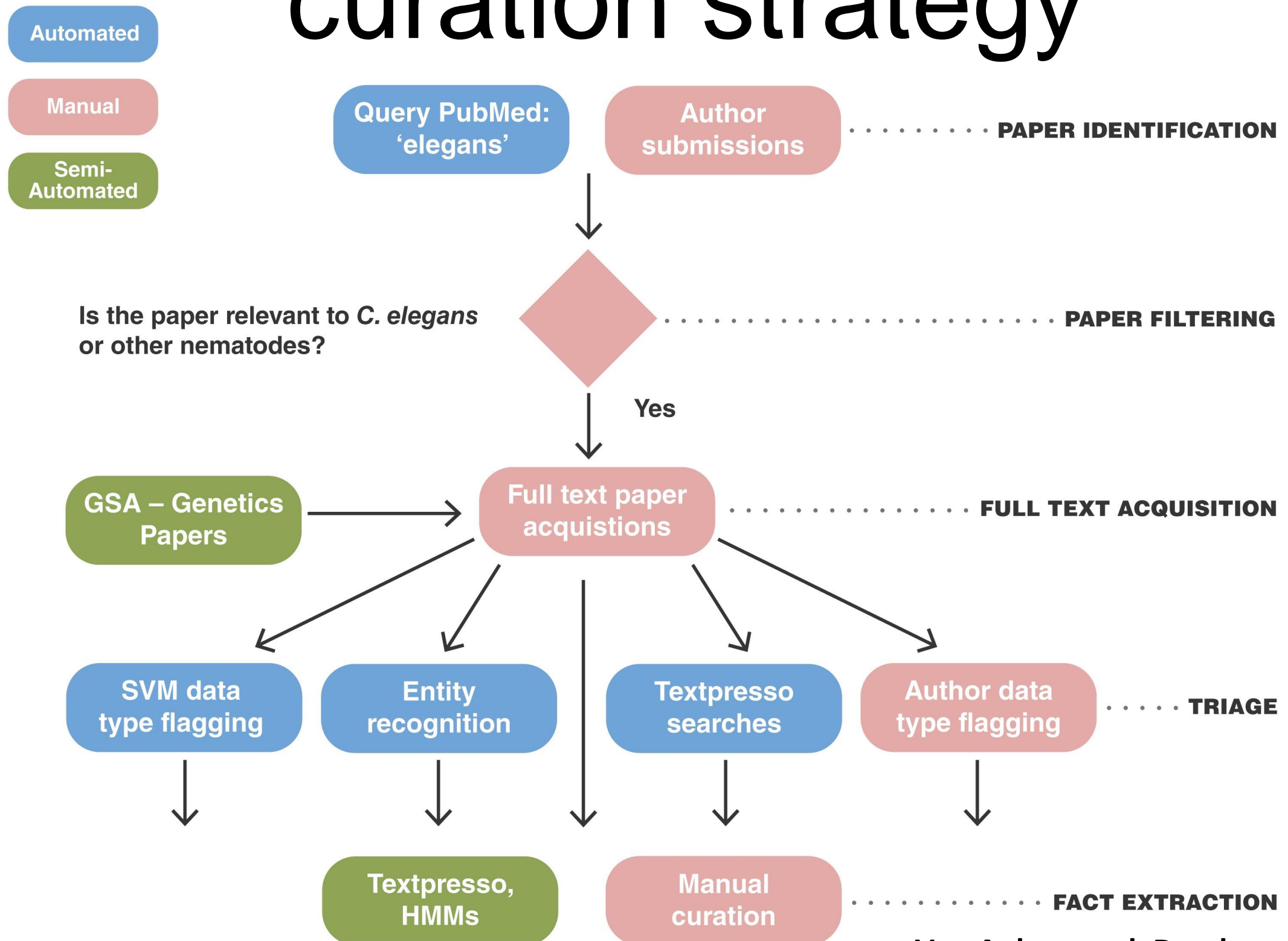Need tools to facilitate.

**Many data types.**

Prioritization is key.
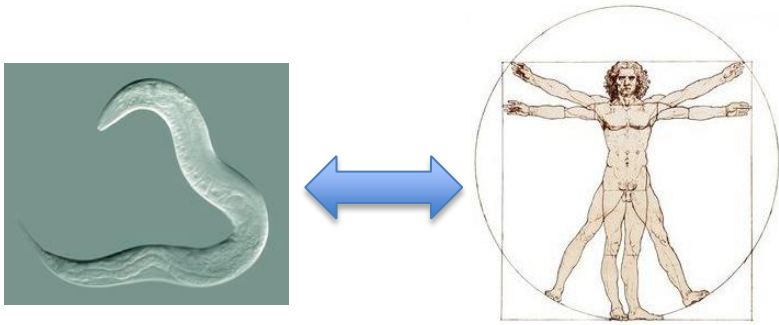
Work procedurally through data types.

Balance of breadth **and** depth critical for making useful community resource.
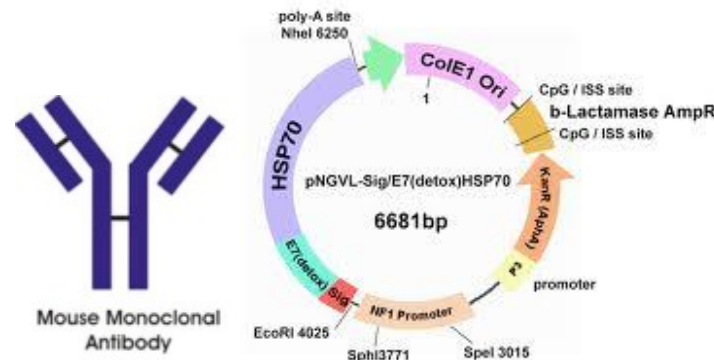
# Hybrid automated/manual curation strategy



Automated

Manual

Semi-Automated

Query PubMed: 'elegans'

Author submissions

· · · · · · · · · · PAPER IDENTIFICATION

Is the paper relevant to *C. elegans* or other nematodes?

· · · · · · · · · · · · · · · · PAPER FILTERING

Yes

GSA – Genetics Papers

Full text paper acquistions

· · · · · · · · · FULL TEXT ACQUISITION

SVM data type flagging

Entity recognition

Textpresso searches

Author data type flagging

· · · · TRIAGE

Textpresso, HMMs

Manual curation

· · · · · · · · FACT EXTRACTION

**Van Auken et al, Database, 2012**

# Curated data types


Human Disease Relevance


Reagents


the Gene Ontology


Anatomy Function


Expression
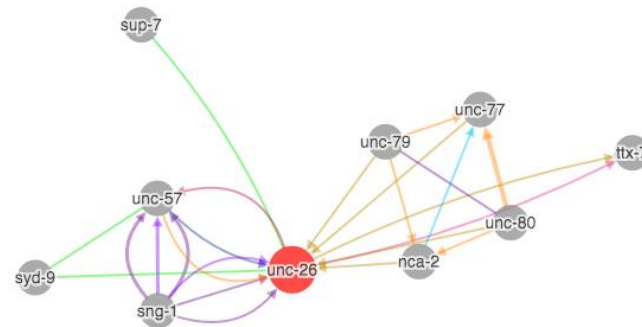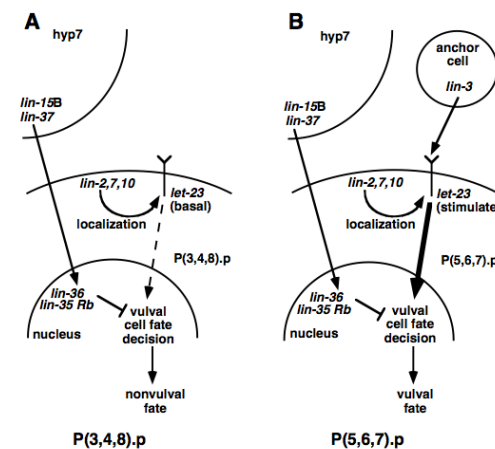

Phenotypes


Gene Interactions


Pathways


Sequence Features

# Reference datasets

**Large scale data at WormBase**

- Proteomics (mass spec)

- Transcriptomics (splicing, UTRs)

- Expression (microarray, *in vivo* imaging)

- Interactions (physical, genetic)

- Perturbation: RNAi, systematic mutation

- Lineage and connectivity

# Reference datasets

**Broad reference data sets can fill knowledge gaps**



- Verification can be difficult

- Relevance?

- Utilization varies greatly. Confidence?

# Do we assess the quality of…

**experimental design? external data?**

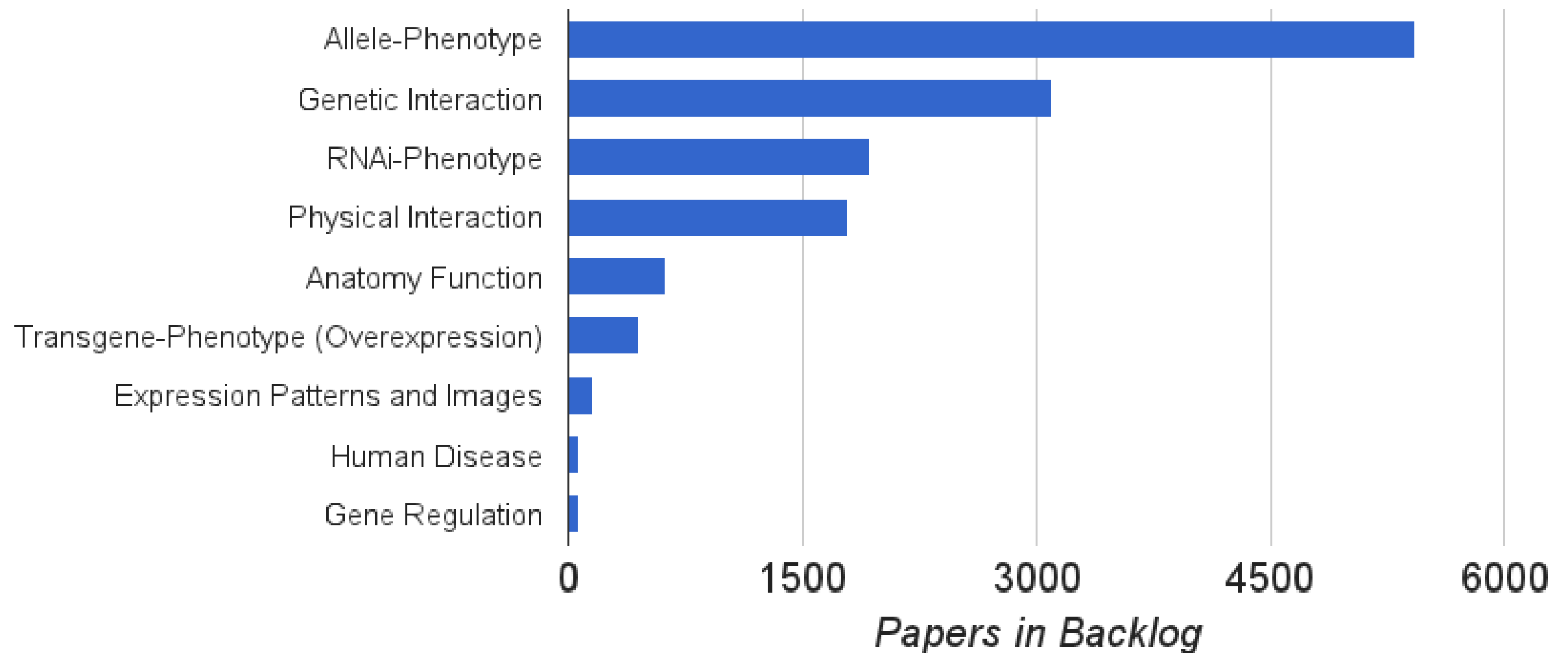**Publication** is the gold standard.

Revisit: erroneous data

Request corrections or clarifications when warranted

# Remaining backlog

# Curation: Lessons Learned

- **harder** and **consumes more time** than expected

- more **enriching** to the final product than expected

- curation ensures data integrity and builds **trust** in the resource
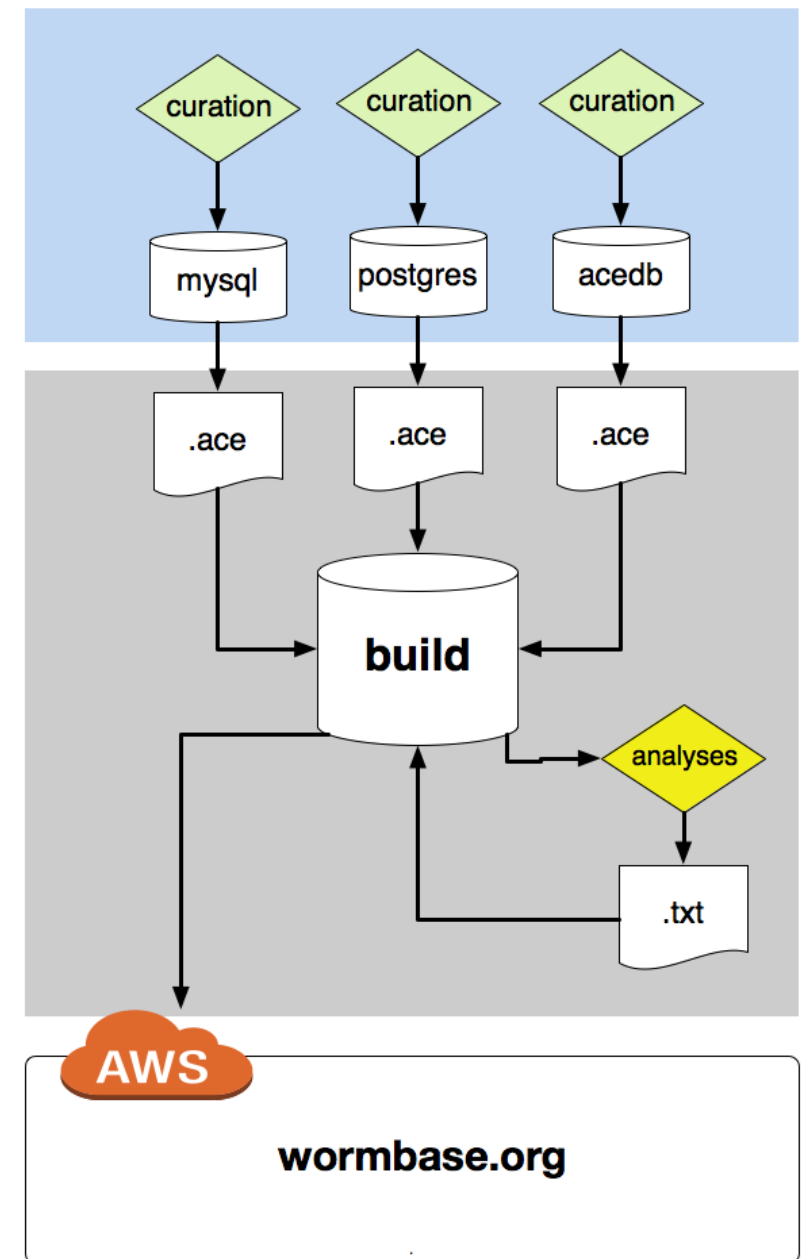
# Curation: Suggestions

- Start early to develop best practices.

- Automate as much as possible.

- Employ domain experts for high value **manual curation** and to confirm **precision** of automated curation.

- **Expect** publication rate and new data types to **exceed** manual curation capacity (10% Y-o-Y).

- **Refining** curation will be an ongoing enterprise.

# What fundamentals have driven our workflow design?

# What fundamentals have driven our design?

## 1. Ease of data modeling and loading



*Emphasis on collecting and sharing data.*

# What fundamentals have driven our design?

**2. Handling unknown unknowns**

Yet-to-be-discovered …

- datatypes

- data relationships

*Data model must be able to evolve.*

# What fundamentals have driven our design?

**3. Ability to track supporting evidence, metadata, and provenance**

*Reproducibility and accountability.*

# What fundamentals have driven our design?

## 4. Coping with high-connectivity data



*eg: What happens to downstream annotations if gene merge? Orthology, proteomics, expression, etc…*

# What fundamentals have driven our design?

**5. Finding a suitable refresh rate**

Datasets evolve. New data becomes available. Analyses need to be updated.

*How often will you update analyses?*

*How tolerant will your community be of **stale data**?*

# What fundamentals have driven our design?

**5. Finding a suitable refresh rate**

1 week -> 2 weeks -> 3 weeks -> 1 month -> 2  months

2001          2002          2005          2008          2011

*Balance of stability, rate of new data, cost/time of analysis, churn.*

# Design: Lessons Learned



1. **A flexible model/workflow is essential.**

2. **Evidence and metdata collection needs to be central to process.**

3. **High connectivity data presents unique challenges.**

4. **Needed to adjust release frequency.**

# Design: Suggestions

1. Build flexibility into both the data model and workflow.

2. Be aware of consequences of changing high connectivity data.

3. Refresh frequency is a balance of user needs, resources, and rate of change.

# Integration & Interoperability

# Suggestions for integrating with organismal databases (easy)

- Liaise with organismal databases early and often!

- Use **stable identifiers**! Most organism databases have them. Please?

# Suggestions for integrating with organismal databases (harder)

**Reciprocal data exchange and cross links**

*Crosslinks alone are boring and do not engage users.*

*Without some supporting context, crosslinks do not increase interoperability.*

# Suggestions for integrating with organismal databases (hardest)

**Avoid direct data import**

*Except for core scaffolding features (genomes, genes, eg), use **APIs** to fetch and embed functional data.*
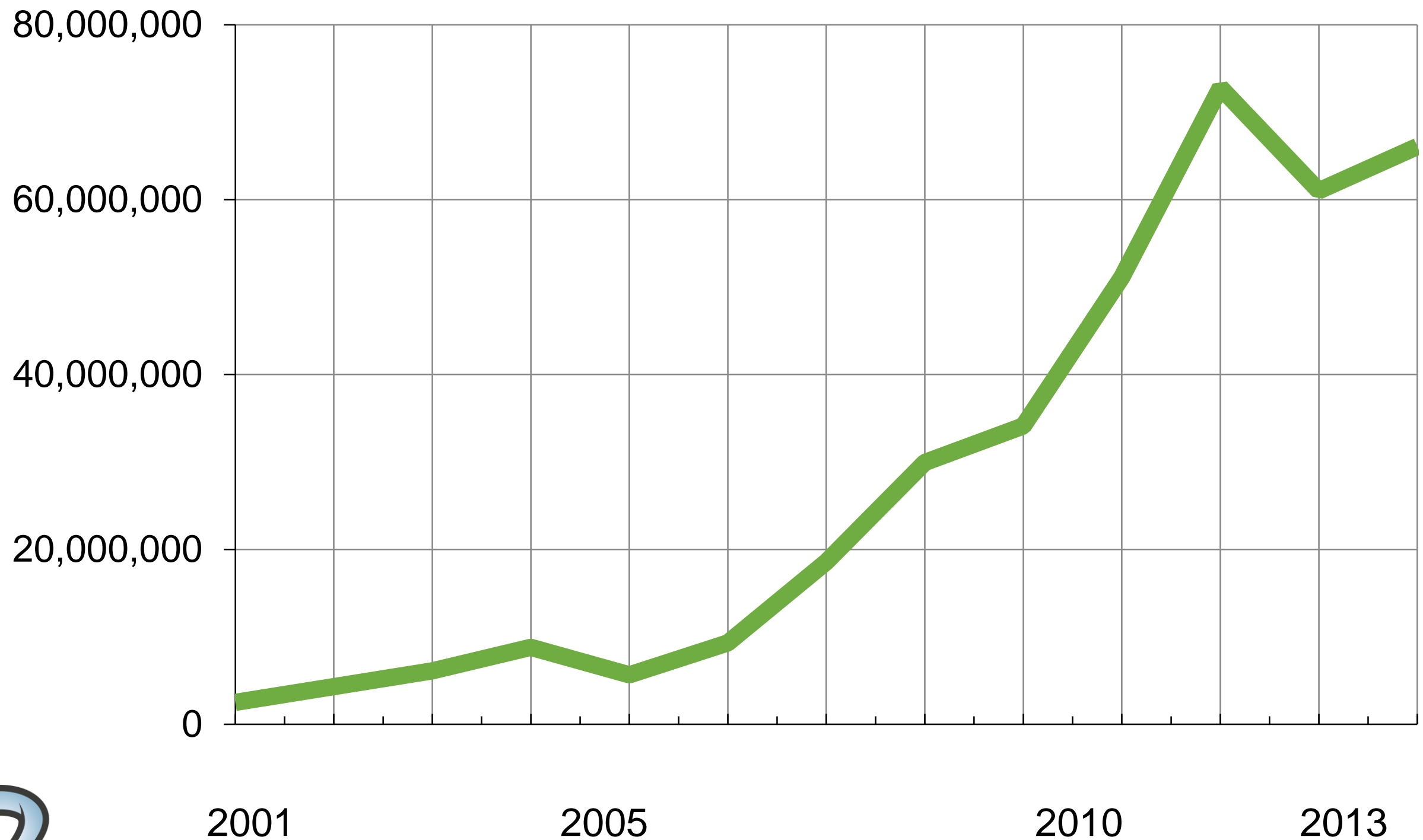
# Interoperability Suggestions

1. Provide data in (multiple) common formats

2. API (RESTful) with JSON and XML delivery

3. Data files programmatically accessible — simple is better (FTP), no registration barrier or fancy web-based download scheme.

4. Consistent, shared identifiers

If you build it, will they come?

Pageviews vs time

# Nurture Your Community

## Collect feedback

Chat, Twitter, Google Alerts, mailing lists, conferences, webinars, surveys.

## Measure

Web logs, CloudWatch, Google Analytics

## Set standards

Data quality, curation, submission, help desk response times.

# Metrics of success

**Not easy to measure.**

Small user communities, niche domains.

Providing annotation or feedback is a low priority for busy scientists.

Positive feedback rare, but you'll **know** when users don't like something!

# Suggested Metrics

- Page Views

- Citation Rate

- Downloads

- Queries & Resolutions

- Rate / precision of curation

- Database size / objects / submissions

# Performance Metrics

# Acknowledgments