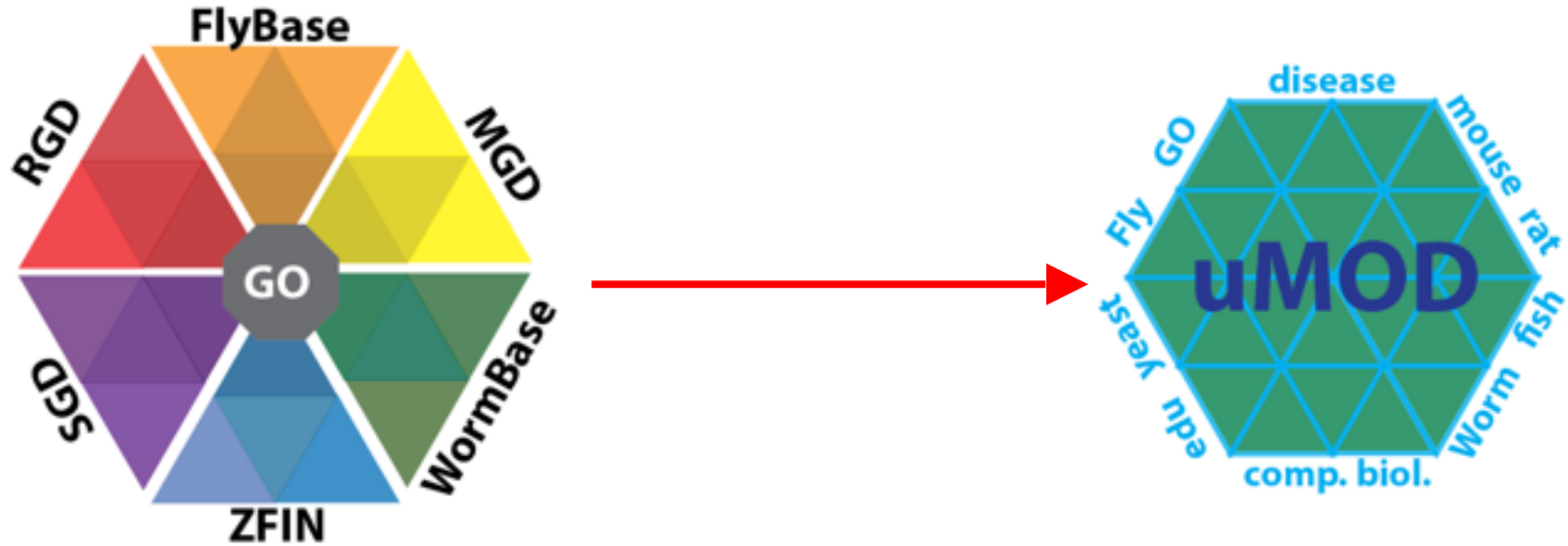


# Bringing the MODs together to create the *Alliance of Genome Resources*



1. Who?
2. Why?
3. How? Long and short term goals
4. Funding implications
5. Should smaller mods to work towards this sort of integration too?

# Impetus for unifying the MODs (from NHGRI)

## User confusion for lack of homogeneity

- User access interfaces
  - need different navigation skills and data access approaches for each resource
  - Semantic inconsistencies and different data structures for the same genomic entities
- Analyses
  - human/model organism association for disease and phenotypes
  - functional annotation
  - Homology representation

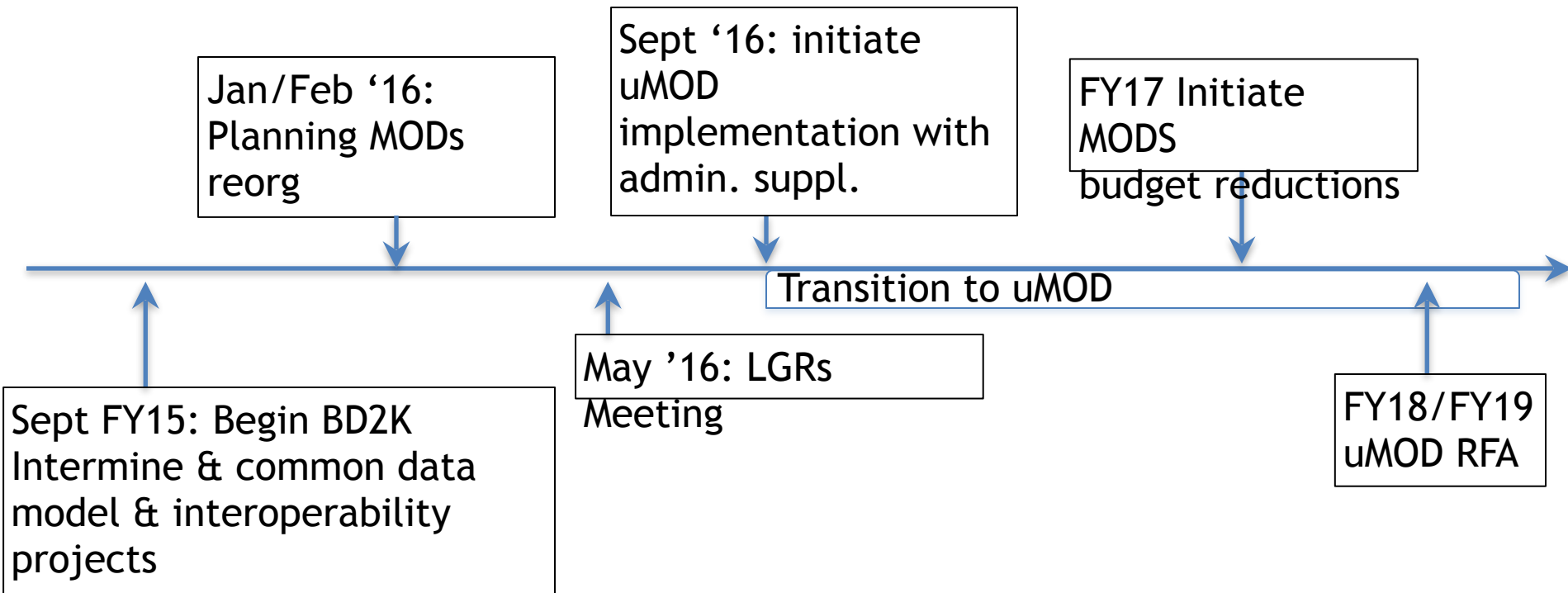
## Redundancy of operations at 6 resources

- Data management systems for related data structures and types
- System administration and IT support
- Technical user support
- Links to the same public resources which need updates and maintenance

# Goals of the Reorganization

- Facilitate access to these resources
- Continue to support the value and services provided by the MODs resources
- Transition the resources to a more effective and sustainable funding model
- Gain flexibility for new informatics program activities at NHGRI

# Timelines



# Unifying the MODs?



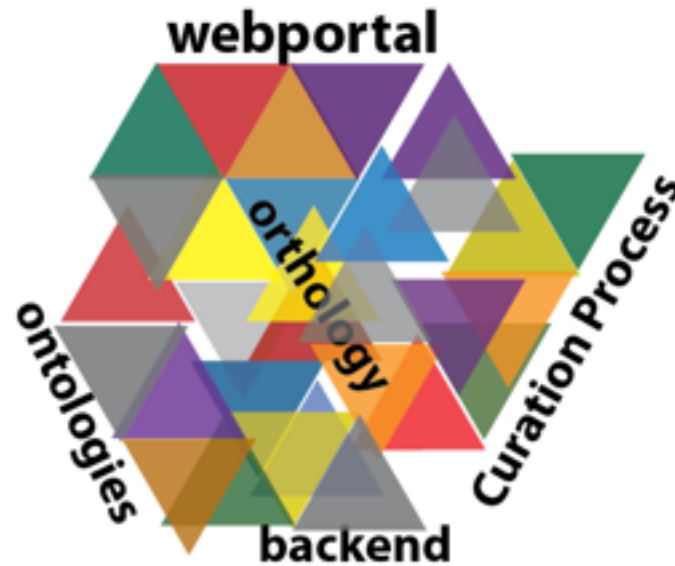
1. Who?
2. Why?
3. How? Long and short term goals
4. Funding implications
5. Should smaller mods work towards this sort of integration too?

# 1. Who are the MODs?



1. Who?
2. Why?
3. How? Long and short term goals
4. Funding implications
5. Should smaller mods to work towards this sort of integration too?

## 2. How



1. Identify and integrate the common tools

# *HOW: Priority Goals for MOD harmonization*

- **Understand existing components of all resources**
- **Understand the strong models/tools within the consortium**
- **Map a clear path to integration and common views to allow federated access to all MOD data**



# Low Hanging Fruit for harmonization

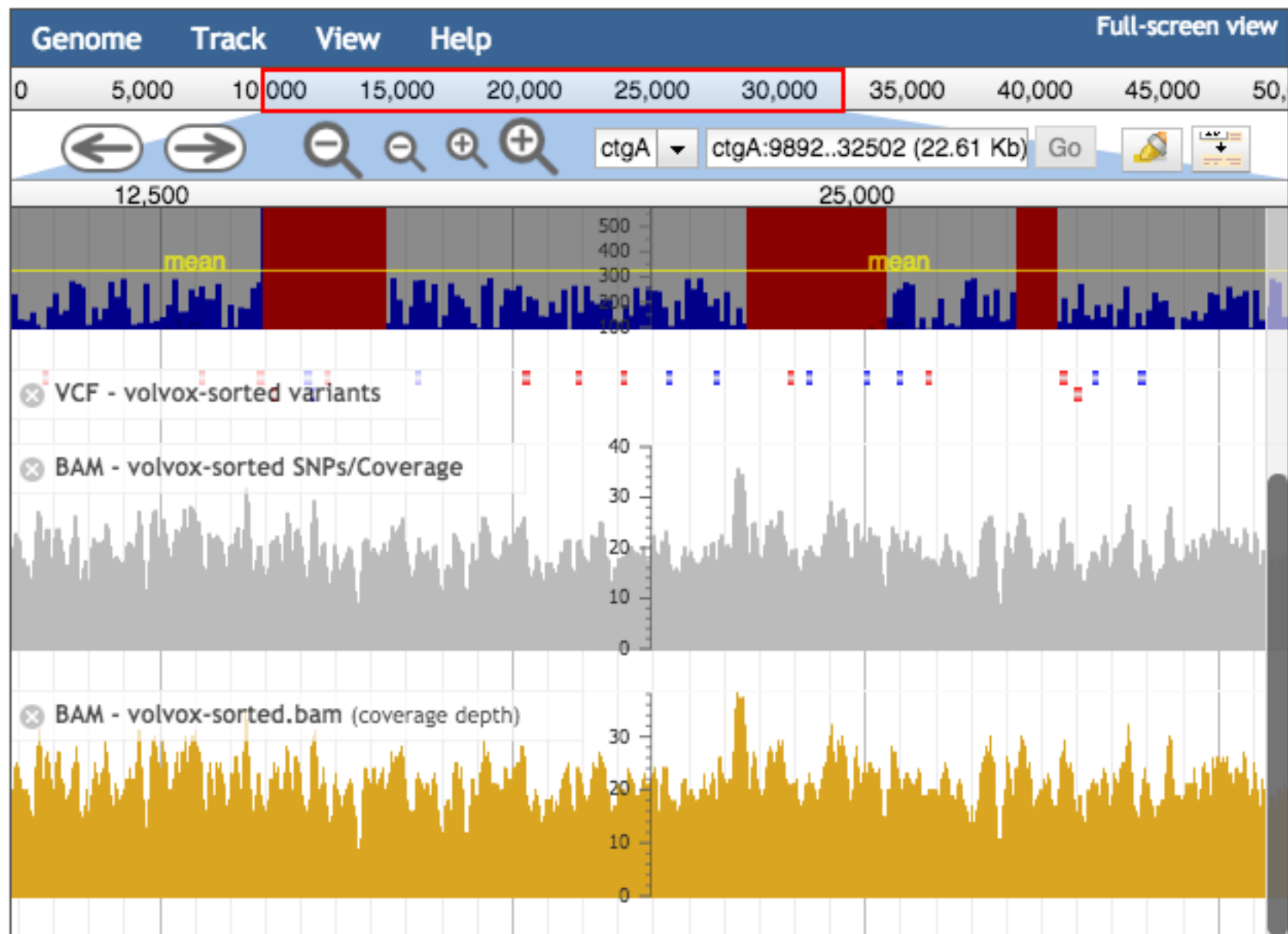
- **Genome Features / Gene Models**
- **Orthology / Homology**
- **Human Disease / Phenotypes**
- **Biomedical Ontologies**



# *Genomes and Genome Features: Two Initial Objectives*

- **Common Genome Browser**
  - JBrowse
  - Already in wide use
- **Genome Features / Gene Detail Pages**
  - Initially protein coding genes
  - Build for all genome features

# Latest Release – [JBrowse 1.12.1](#)



# *Initial Steps for Gene Pages*

- **Support query on any official gene symbol or synonym**
- **Summation page with possible matches and why**
- **Links to current gene detail pages at MOD sites**

# *Gene Model Standardization*

## Core Identifiers

- primaryIdentifier
- secondaryIdentifier
- Symbol
- Name
- Cross references

## Descriptions

- Limited to two fields: description and briefDescription

## Gene Types

- originally represented by several field names, featureType, geneType etc
  - Will be limited to Sequence Ontology terms
-

# *Gene Detail Pages*

- **Protein-coding Genes**
  - Isoforms and modified protein forms
- **RNA genes**
  - Regulatory interactions
- **Two aspects to capture**
  - ‘what does this gene do’ summations
  - Drill down to data detail

# Summations Ex: Phenotype Ribbon

▼ **Phenotype Summary** 191 phenotypes from 6 alleles in 13 genetic backgrounds  
31 phenotypes from multigenic genotypes  
4 images  
193 phenotype references

**All Mutations and Alleles** 8  
Targeted 8  
Incidental Mutations Mutagenetix , APF

**Phenotype Overview** ?

The ribbon chart displays 28 categories, each with a blue bar indicating the number of annotations. The categories are: adipose tissue, behavior/neurological, cardiovascular system, cellular, craniofacial, digestive/alimentary, embryogenesis, endocrine/exocrine gland, growth/size/body region, hearing/vestibular/ear, hematopoietic system, homeostasis/metabolism, immune system, limbs/digits/tail, liver/biliary system, mortality/aging, muscle, nervous system, pigmentation, renal/urinary system, reproductive system, respiratory system, skeleton, taste/olfaction, tumorigenesis, and vision/eye.

Click cells to view annotations.

Homozygous targeted mutants displayed vascular system dysfunctions and thickening of lung aveolar septa from hyperproliferation and fibrosis, ultimately causing the mice physical limitations. Mice also display increased incidence of calcium calculi, kidney stones, and decreased adiposity.

- Click cells to view annotations

# *Orthology / Homology: Two Initial Objectives*

- **Define uMOD orthology set(s) for comparative tools and representations**
- **Incorporate comparative orthology tool to interrogate and visualize orthology data.**



# *Define uMOD orthology set(s) for comparative tools*

- **InterMOD**
  - Each MOD's orthology used with search on human gene
  - But what to do when no human gene in set?
- **Panther / GO \*\*\***
  - Already in deep use in AmiGO; extends to analysis
  - Positive: Quest 4 Orthologs provides community input into definition of reference proteomes to drive algorithm
  - Negative: Updates and completion for some MODs needed
- **Other? InParanoid?**
- **ZFIN - custom set, special purpose**

## *Incorporate comparative orthology tool to interrogate and visualize orthology data.*

- **Model 1: HCOP**
  - Nice visualizations
  - Only from Human, to 17 model organisms
- **Model 2: DIOPT (\*\*\*\*)**
  - Human, fly, mouse, worm, yeast, zebrafish
  - Confusing stats
  - ‘in house’ so further improvements quick
- **Comparative data graphs**
  - <http://www.informatics.jax.org/homology/GOGraph/7247>

# *Phenotype and Disease*

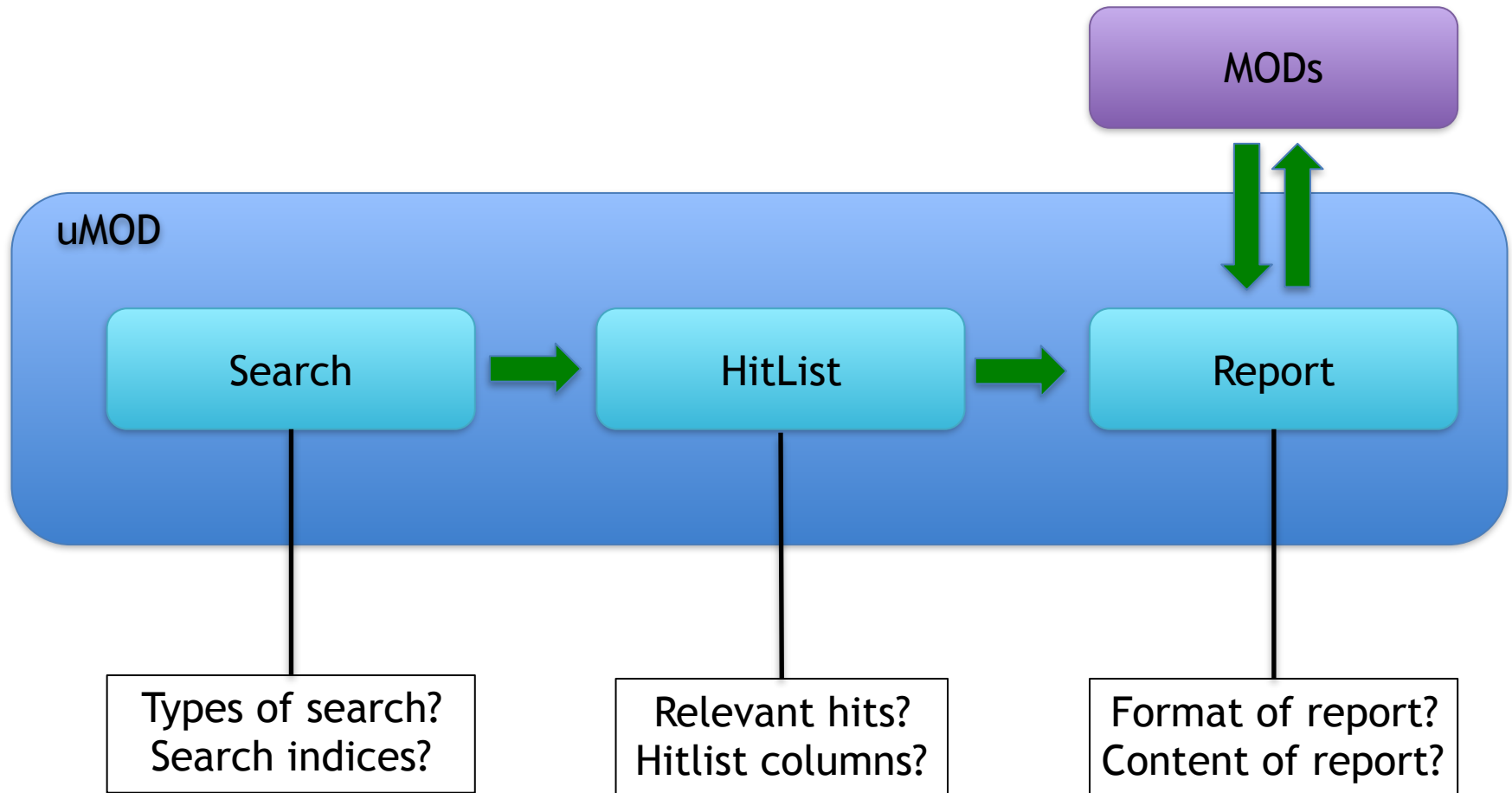
Cindy Smith and Carol Bult (MGD)  
Suzanna Lewis (GO)  
Mary Shimoyama(RGD)  
Monte Westerfield (ZFIN)  
Thom Kaufman(FlyBase)  
H-MOD Disease Working Group

---

# *Two Initial Objectives*

- **Support query by human disease term**
  - OMIM
  - Disease Ontology (DO)
- **Support query by phenotype term**
  - Continued work on alignment of vocabularies

# Objective



# Use cases

## 1. Search with a human gene

- a) Which **diseases** are associated with this human gene?
- b) What are the **orthologs** for this human gene in MOs, and are there **experimental disease models**?
- c) Which **diseases** are associated with the region containing this human gene (GWAS)?

## 2. Search with a Model Organism gene

- a) Which **diseases** are associated with human ortholog(s) of this gene?
- b) Are there existing **experimental disease models** in this MO using this gene, or a human transgene?
- c) What are the **orthologs** for this gene in other MOs, and are there **experimental disease models**?
- d) Which **diseases** are associated with the region containing the human ortholog(s) of this gene (GWAS)?

## 3. Search with a human disease

- a) Which **human gene(s)** are associated with this disease?
- b) Are there **experimental disease models** in MOs for this disease?
- c) Which human **genomic region(s)** (GWAS) are associated with this disease, what human genes are contained therein? (What are there MO orthologs, and are there experimental disease models?)

## 4. Search with a human genomic region?

- a) Which **diseases** are associated with this region, either direct gene-disease associations or via GWAS?
- b) Are there **experimental disease models** in MOs for these diseases?

## 5. Search with a 'phenotype'?

# Disease Model annotation table?

Disease		Human gene		Model Organism gene			
Name (DO)	Subtype (OMIM pheno)	Symbol	Genomic location	Species	Symbol	Ortholog score	Exp. models ?
Alzheimer's disease	ALZHEIMER DISEASE; AD	APP	21q21.3	Mouse	<a href="#">App</a>	10	108
				Rat	<a href="#">App</a>	10	343
				Frog	<a href="#">app</a>	6	-
				Fish	<a href="#">appb</a>	6	-
				Fly	<a href="#">Appl</a>	7	4
				Fly	<a href="#">Hsap\APP</a>	n/a	36
				Worm	<a href="#">apl-1</a>	7	?
		APBB2	4p14-p13	Mouse	<a href="#">Apbb2</a>	8	-
				Frog	<a href="#">apbb2</a>	5	-
				Fish	<a href="#">apbb2b</a>	3	-
	Worm			<a href="#">feh-1</a>	6	-	
	ALZHEIMER DISEASE 2	APOE	9q13.32	Mouse	<a href="#">ApoE</a>	8	-
				<i>etc...</i>	<i>etc...</i>	<i>etc...</i>	<i>etc...</i>

# *Biomedical Ontologies*

GOC group

---





Free-text filtering X

Bmp4

Your search is pinned to these filters

+ document\_category: annotation

No current user filters.

Source

Assigned by

Ensembl	(688)	+	-
UniProt	(533)	+	-
GO_Central	(269)	+	-
RGD	(266)	+	-
MGI	(241)	+	-
ZFIN	(43)	+	-
BHF-UCL	(41)	+	-
AgBase	(9)	+	-
DFLAT	(8)	+	-
InterPro	(8)	+	-
UniProtKB	(6)	+	-
GOC	(3)	+	-
WormBase	(3)	+	-
CACAO	(2)	+	-
IntAct	(2)	+	-
HGNC	(1)	+	-
Reactome	(1)	+	-

Ontology (aspect)

Evidence type

## Found entities

Total: 2124; showing 131-140

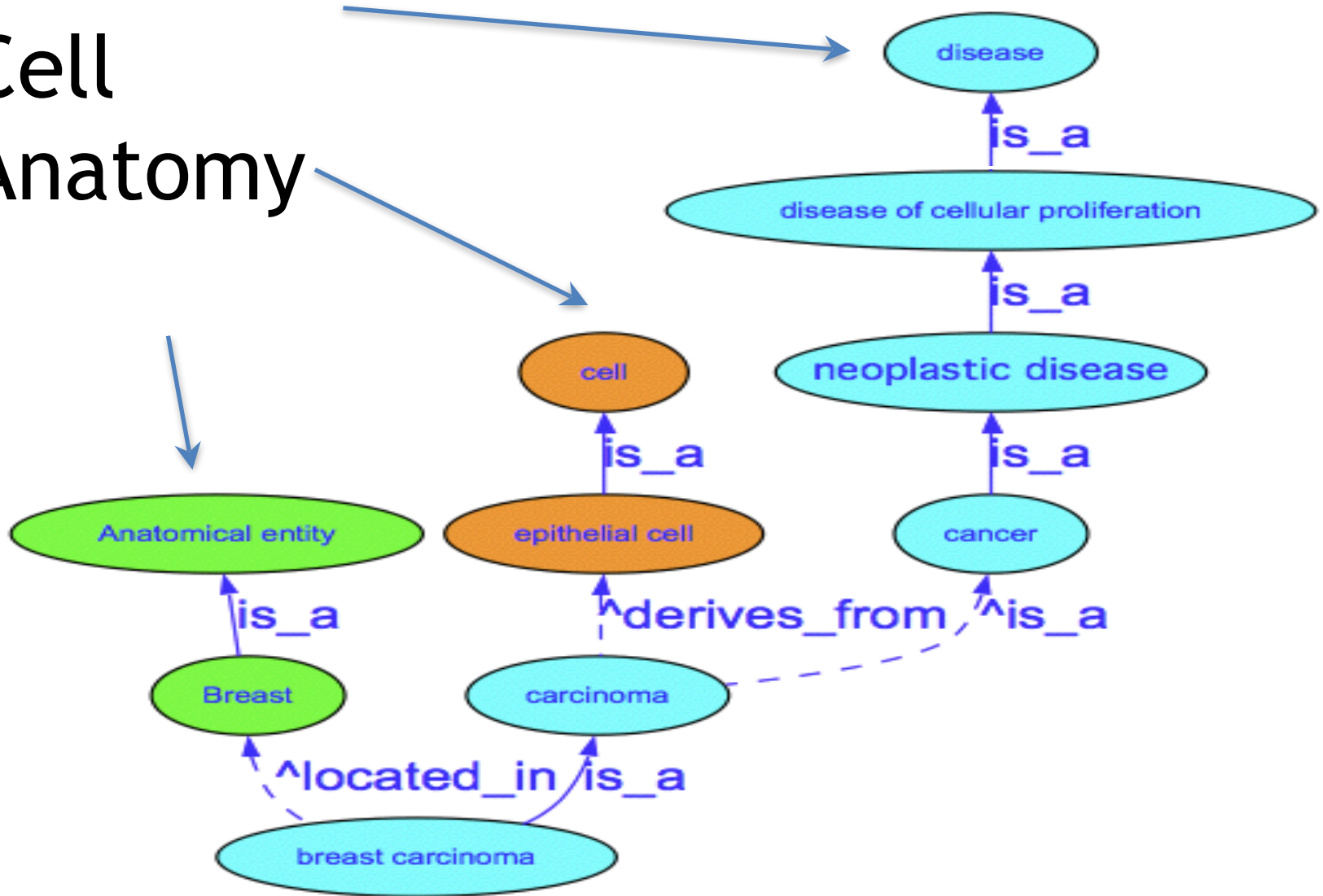
Results count 10

<input type="checkbox"/> Gene/product	Gene/product name	Qualifier	Direct annotation	Annotation extension	Assigned by	Taxon	Evidence	Evidence with	PANTHER family	Is
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		extracellular region		ZFIN	Danio rerio	IEA	UniProtKB-KW:KW-0964	tgf-beta family pthr11848	
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		growth factor activity		ZFIN	Danio rerio	IEA	UniProtKB-KW:KW-0339	tgf-beta family pthr11848	
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		pronephros development		ZFIN	Danio rerio	IMP	ZFIN:ZDB-MRPHLNO-050429-3	tgf-beta family pthr11848	
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		pronephros development		ZFIN	Danio rerio	IMP	ZFIN:ZDB-MRPHLNO-080919-1	tgf-beta family pthr11848	
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		digestive tract development		ZFIN	Danio rerio	IMP	ZFIN:ZDB-MRPHLNO-050228-1	tgf-beta family pthr11848	
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		heart development		ZFIN	Danio rerio	IEP		tgf-beta family pthr11848	
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		proepicardium development		ZFIN	Danio rerio	IMP	ZFIN:ZDB-GENO-080213-1	tgf-beta family pthr11848	
<input type="checkbox"/> <a href="#">bmp4</a>	bone morphogenetic protein 4		embryonic hemopoiesis		ZFIN	Danio rerio	IMP		tgf-beta family pthr11848	

# *Biomedical Ontologies*

- **GO:** (function, process, cellular location)
- **SO:** (sequence features)
- **PRO:** (specific proteins by species/strain)
- **MP, HPO, others:** (phenotypes)
- **Anatomies / Homologies** (morphology)
- **DO:** (diseases, not phenotypes; definitions not diagnoses)
- **CL:** (cells and their lineages)

Disease  
Cell  
Anatomy



# *Summary*

- **Start with compelling use cases that join models and build collaborative environment**
- **Small work groups and teams define requirements and bring to steering committee for refinement and agreement**
- **Groups exist now, focus and interactions strong**

# Acknowledgements

- Karen Yook and the other curators at WormBase for helping me sort some of this out.
- Lincoln Stein, Paul Sternberg, Paul Kersey and Matt Berriman , co-PIs of WormBase.

# A Few Questions

- 1) Specifically, for each of the 5 MODs and GO, what activities and data would be supported in common, and which would be unique to each site?
- 2) What is the likely platform(s) for the common components? For the individual sites?
- 3) Which common elements/activities might be phased in by what order?
- 4) How will issues of scalability and variety of data types be addressed?
- 5) Would there be incremental deliverables along the way. What might these specific milestones be?
- 6) Are there any components or steps for which there is already agreement among the groups?
- 7) What are the components or steps which are not yet resolved?
- 8) Are there any components or elements where there are clear conflicts, or differences of opinion among the sites now?
- 9) In the proposed new model, where exactly will the savings be realized? Staff reductions? Which staff?
- 10) Will any current activities by any of the participants stop? Which ones?
- 11) Will any new activities not currently done by any participant be started?

# uMOD Functions

Data Ingest

Integration

Services

Interfaces

Literature  
prioritization  
and  
markup

Textpresso

Online  
Curation

Apollo

Noctua

Phenote

Large scale  
data  
acquisition

Orthology  
calling and  
inferencing

Quality  
Control

Syntax  
validators

Identifier  
resolution

Data  
delivery

Query  
Engines

Enrichment  
analysis

Web  
portal

community  
outreach and  
support

Third  
party  
tools

SHARED SEMANTIC FRAMEWORK

# SHARED FRAMEWORK

UBERON + Cell

Assays

Evidence

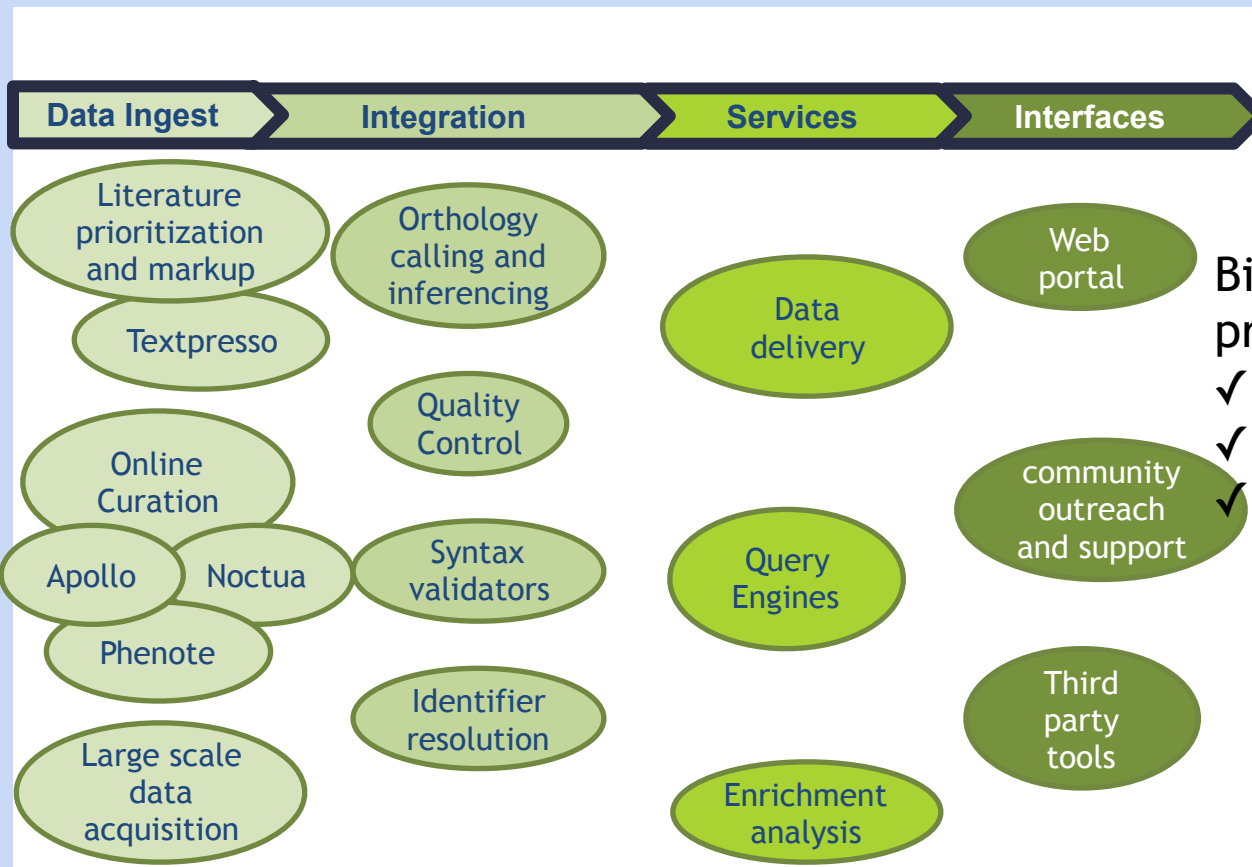
Phenotyp

- Type
- ✓ Human
  - ✓ Vertebrate anatomy
  - ✓ Mouse anatomy
  - ✓ Fly
  - ✓ Etc.

NCBI taxonomy

Environment

- ✓ HPO
- ✓ WP
- ✓ MP
- ✓ Yeast
- ✓ Etc.



Cellular component

- Biological process
- ✓ Physiological
  - ✓ Behavior
  - ✓ Population

Molecular function

Sequence features

Molecular process pathway



# Recommendation: Standard APIs

- **Shared Web APIs for MODs and related resources**
  - Standard **URL structure** to programmatically fetch database objects
  - Standard **JSON format** for returning data
- **Wins**
  - Loose coupling, lightweight, no/few backend changes required
  - Shared UI components
    - Steps toward more consistent front ends
  - (Power) User efficiencies
  - Steps towards tighter integration

# Standardized exchange formats

- **Pre-requisite for Web API** (previous slide)
  - Can be implemented independently
  - Replace or extend current ad-hoc database dumps available from most MODs FTP sites
  - Enables aggregators and bioinformatics power users
    - Aggregators: Intermine, Monarch, ClinGen
- **Recommendations:**
  - JSON or JSON-LD
    - Not XML, TSV
  - Leverage semantics of ontologies rather than reinvent
- **Track with existing efforts**
  - GA4GH: G2P, metadata, ...
  - FHIR, Bio2RDF
  - PhenoPackets
  - Domain-specific efforts
    - Pathways, interactions, pathways, ...

# Best behavior for Identifiers

1. Use established identifiers
2. Design unique identifiers for use by other groups
3. Help local identifiers travel well: document Prefix and Namespace
4. Opt for simple durable web resolution
5. Avoid embedding meaning
6. Make URIs clear and findable
7. Implement a version management policy
8. Do not re-assign or delete identifiers
9. Document the identifiers you issue and use
10. Reference responsibly

# Shared UI components

- Current shared components
  - Industry-standard 3<sup>rd</sup> party frameworks (e.g. jquery)
  - Occasional use of shared bioinformatics components
    - E.g. JBrowse genome browser
  - Sometimes: intermine, biomart

# Standardized Curation (Tools)

- Begin with curator developed curation standards
- Move to shared curation tools, for example:
  - Apollo for genome features
  - Noctua GO curation tool
    - Being extended to model phenotypes
  - TextpressoCentral (integrated with Noctua now!)
- Wins
  - Curator efficiencies
  - Developer efficiencies (less duplication of effort)
  - Community contribution

# A new way of developing software: Commons based Peer Production

Distributed  
Version  
Control

# GitHub

Social  
Coding

1



## Set up Git

A quick guide to help you get started with Git.

2



## Create repositories

Repositories are where you'll work and collaborate on projects.

3



## Fork repositories

Forking creates a new, unique project from an existing one.

4



## Work together

Send pull requests, follow friends. Star and watch projects.

Graph-based audit trail of every contribution

Issue tracking

Radical transparency

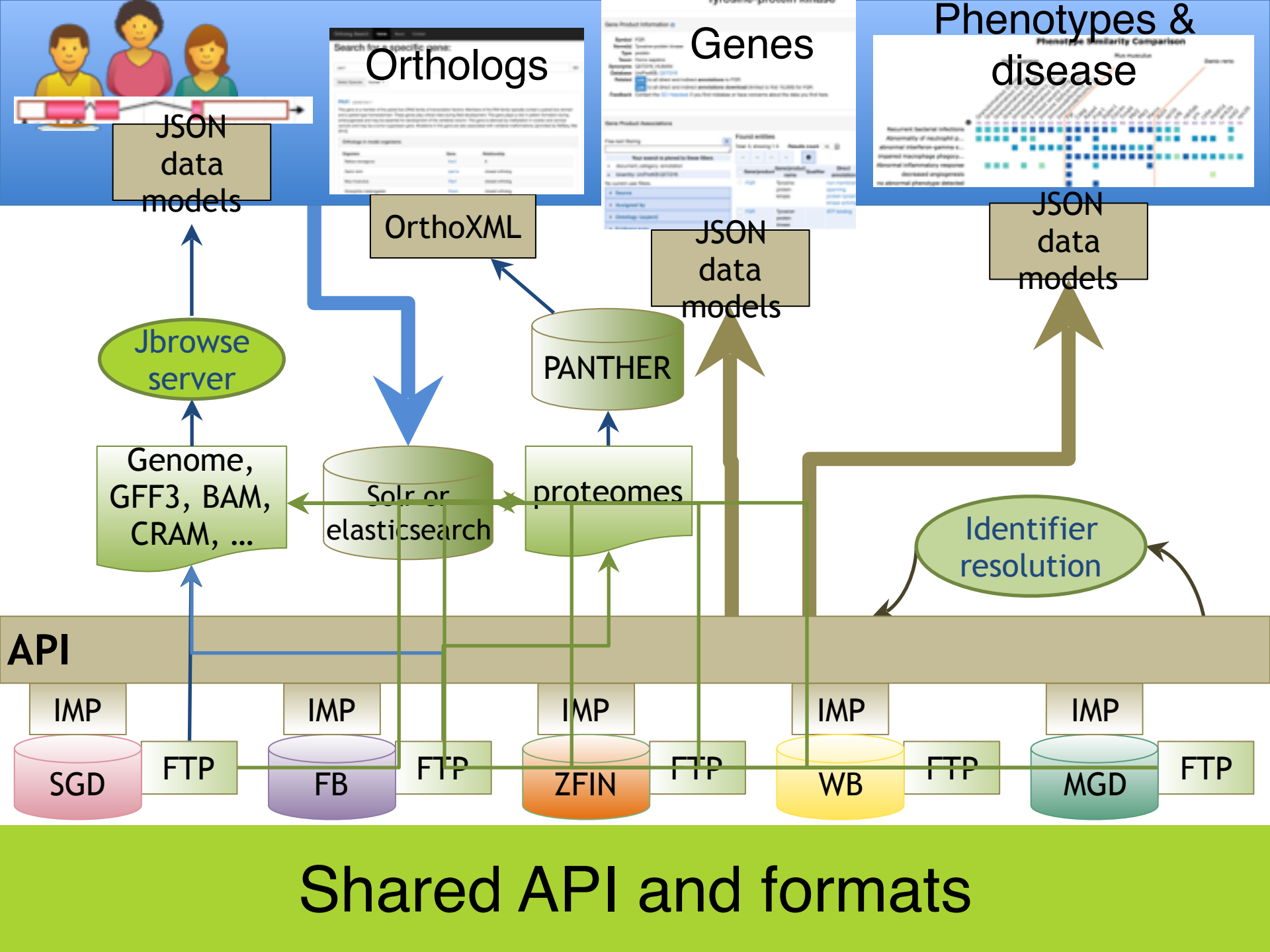
# Commons based Peer Development for community standards

- Does the GitHub model make sense for community standards development?
  - Yes!
  - Watch this space: Global Alliance for Genomic Health (GA4GH)
    - forks, branches, +1s, pull requests, attribution, immediacy, dynamism, experiments and crazy ideas

# Shared Hygiene

- Use of version control systems
  - Open hosted repos (GitHub, GitLab)
  - Trackers for features and bugs
- Unit tests and Continuous Integration
- Common language-level APIs





# Potential end goal: Shared full stack

- A unified database/schema and architecture
  - Long term developer efficiencies
  - Reuse of core backend tools, e.g. curation tools
  - Structural analysis of the relationships between different entities organized in huge networks of graph-like structures enables data science
- Recommendations:
  - Incremental steps: Start with common exchange format (bootstrap), more shared tooling and APIs
  - Exploratory task force to evaluate Cost/benefit tradeoffs



# MODs have similar user

*human geneticists* who want access to all model organism data which are the main source of experimental annotation of human genes

*basic science researchers* who use specific model organisms to understand fundamental biology

*computational biologists and data scientists* who need access to standardized, well-structured data, both big and small

*educators and students* who want to teach and learn

# uMOD Consortium Mission

**Develop and maintain a unified information resource that facilitates the use of diverse model organisms in understanding the genetic and genomic basis of human biology and health. This understanding is fundamental for advancing genome biology and for translating genome data into clinical utility**

# The MATRIX

**Matrix of shared infrastructure modules**

comprising

**outwardly facing components** such as common web portal, APIs, hosted online curation and data analysis tools, web services, and community outreach;

**internal data management applications** in concert with biological expertise to enable comprehensive integration across biological domains including, orthology analysis, a shared ontology framework, and data consistency verification.

**Expert teams that select or develop modules** from sources within and outside the consortium.

**Organism-specific working groups** that use shared



**A unified resource that assembles, integrates, mines and disseminates information about intensively studied organisms**

<input type="text" value="search any term"/> options for searches	<b>human</b>	<b>genomes</b>
<b>gene or gene list</b>	<b>mouse</b>	<b>gene products</b>
	<b>rat</b>	<b>homology</b>
<input type="text" value="full text search of papers"/>	<b>zebrafish</b>	<b>variants</b>
	<i>Drosophila</i>	<b>function</b>
<b>data mining</b>	<i>C. elegans</i>	<b>pathways</b>
	<i>Saccharomyces</i>	<b>networks</b>
<input type="text" value="help!"/>	<b>Gene Ontology</b>	<b>physiology</b>
	<b>Clinicians</b>	<b>anatomy</b>
	<b>submit data</b>	<b>computational bio</b>
		<b>education</b>



A unified resource that assembles, integrates, mines and disseminates information about intensively studied organisms

**search any term**

options for searches

- gene
- variant
- phenotype/disease
- anything

---

**gene or germline**

options for

- ID mapping
- GO enrichment
- sequence
- phenotype./disease
- tissue enrichment

---

**full text search of papers**

---

**data mining**

---

**help!**

**human**

**mouse**

**rat**

**zebrafish**

*Drosophila*

*C. elegans*

*Saccharomyces*

**Gene Ontology**

**Clinicians**

---

**submit data**

**genomes**

**gene products**

**homology**

**variants**

**function**

**pathways**

**networks**

**physiology**

**anatomy**

**Computational Bio**

**Education**



**human**

**mouse**

**rat**

**zebrafish**

*Drosophila*

*C. elegans*

*Saccharomyces*

**Gene Ontology**

**Clinicians**

**Links to community-specific portals that replace and enhance existing MODs or new views**

Perform simple searches on data

gene
variant
phenotype/disease
anything

Perform queries and analyses of gene lists

ID mapping
GO enrichment
sequence
phenotype /disease
tissue enrichment

Textpresso searches of PMC

interMOD interMine

**search any term**  
options for searches

**gene or gene list**  
options for list of genes

**full text search of papers**

**data mining**

**help!**

# InterMine platform from InterMOD

Will enable federated queries between MOD data and the network of linked data

- InterMOD standardizes access to MOD data with InterMine software.
- InterMine creates data warehouses that integrate data and enable bioinformatics analyses.
- 28 InterMines are available for different community databases:

- **YeastMine** (SGD), **WormMine** (WB), **FlyBaseMine** (FB), **ZebrafishMine** (ZFIN), **MouseMine** (MGI), **ThaleMine** (AraPort), **HumanMine** (Micklem), **RatMine** (RGD), **XenMine** (Cherry), **BovineMine** (USDA), **SoyMine** (USDA), **modMine** (modENCODE), **GrapeMine** (EU), **LegumeMine** (NSF), **PlanMine** (planarian, EU),

...

- Core data object model (commonly used) + specific MOD customizations for flexibility



**genomes**

**gene products**

**homology**

**variants**

**function**

**pathways**

**networks**

**physiology**

**anatomy**

**Computational Bio**

**Education**

**Links to topic-specific portals that provide global entry into the combined data.**

**New pathway viewer + reactome + etc.**

**Anatomy ontology browsers and links to atlas sites**

# Link to community annotation systems

**submit data**

**about**

**people**

**funding**

**friends**

**downloads**

**news**

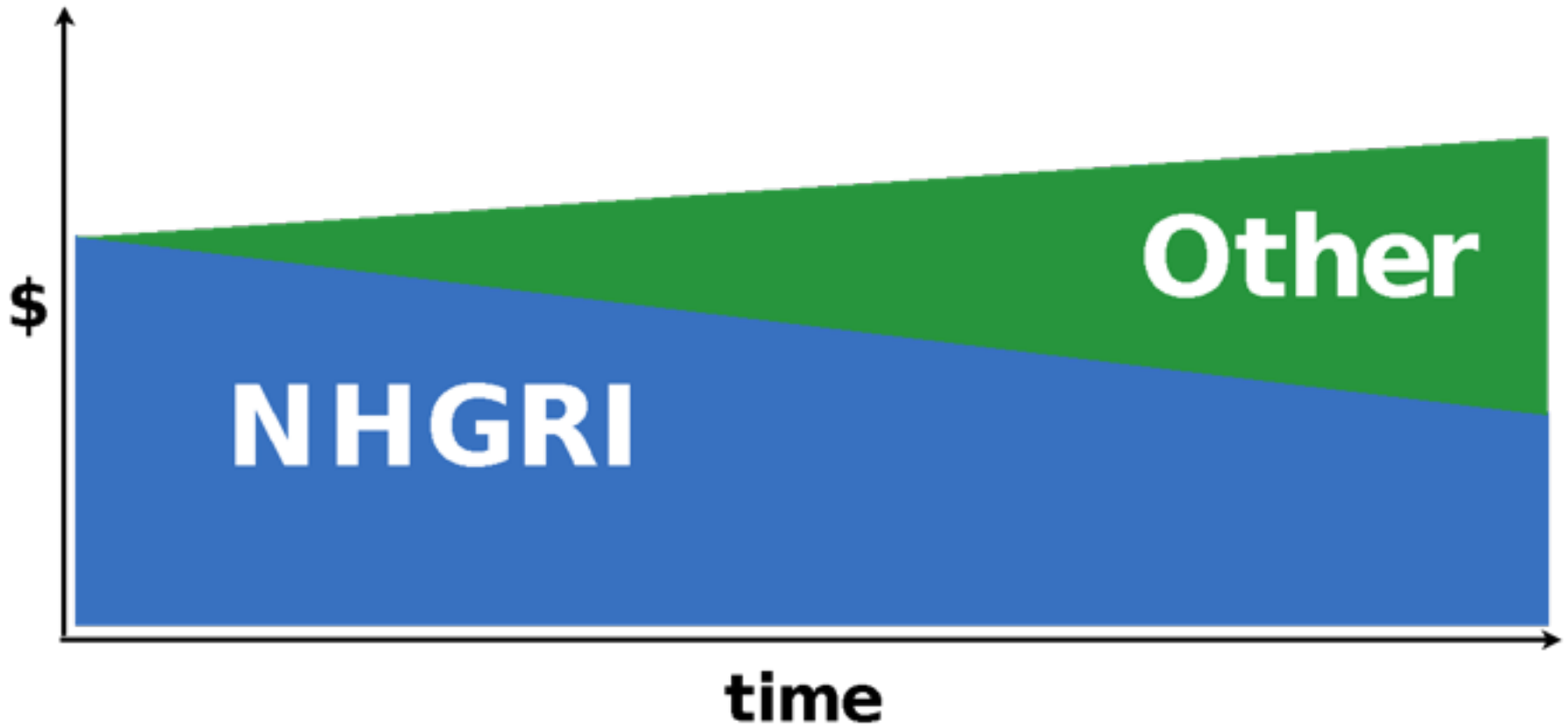
**community**

**support**

**Link to Standard website  
pages, ftp sites, help desk,  
etc.**

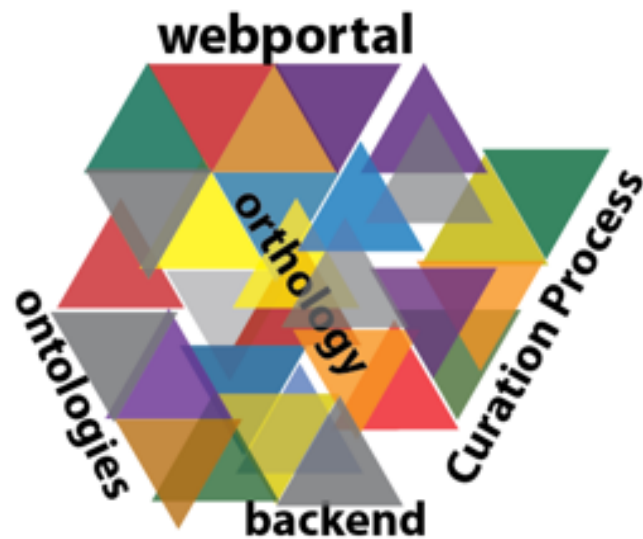
# The Finances

## Funding Plan





**2016**



**2018**



**2020**

**Sept  
2017**

Web Portal sufficient for use and as target of further development

**Sept  
2018**

Backend sufficient for use and as target of further development