# Biocuration:
## Helping Researchers Harness the Data Explosion at TAIR and the Plant Metabolic Network

Kate Dreher

*curator*

TAIR/PMN

Department of Plant Biology

Carnegie Institution for Science

Stanford, California

kadreher@stanford.edu

## Overview

- Biological data **explosion**

- Biocurators want to help!

- Biocuration practices and resources at two plant databases

  - The Arabidopsis Information Resource

  - The Plant Metabolic Network

- Request for *your* help!
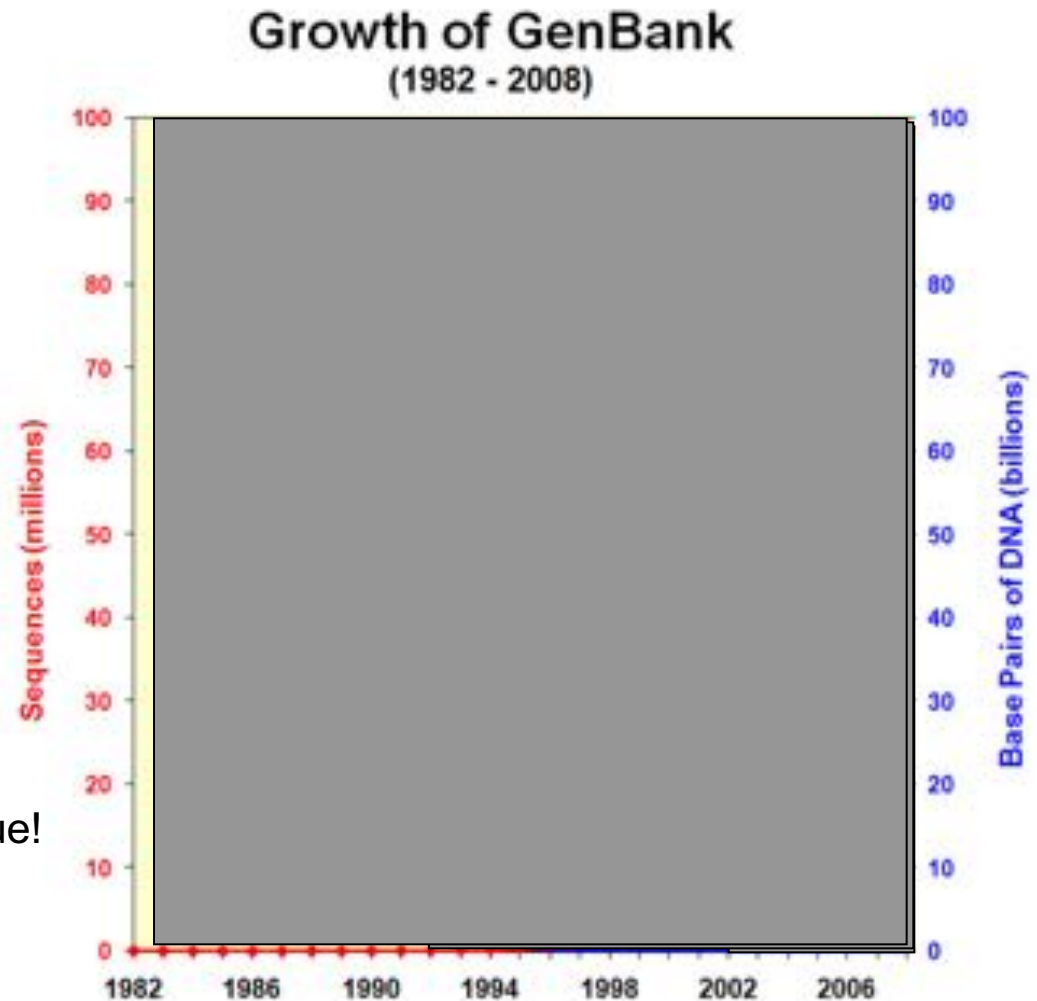
# Growth of biological data

- Over time biological data increases in
  - Quantity
    - Methods improve
    - Costs decrease

# Growth of biological data

- Nucleotide sequences

  - Number of sequences in 1982
    - **606**

  - Number of sequences in 1992:
    - **78,608**

  - Number of sequences in 2002:
    - **22,318,883**

  - Number of sequences in 2008:
    - **98,868,465**

  - And, the acceleration may continue!

Growth of GenBank (1982 - 2008)

# Growth of biological data

- Over time biological data increases in

  - Complexity
    - Protein data
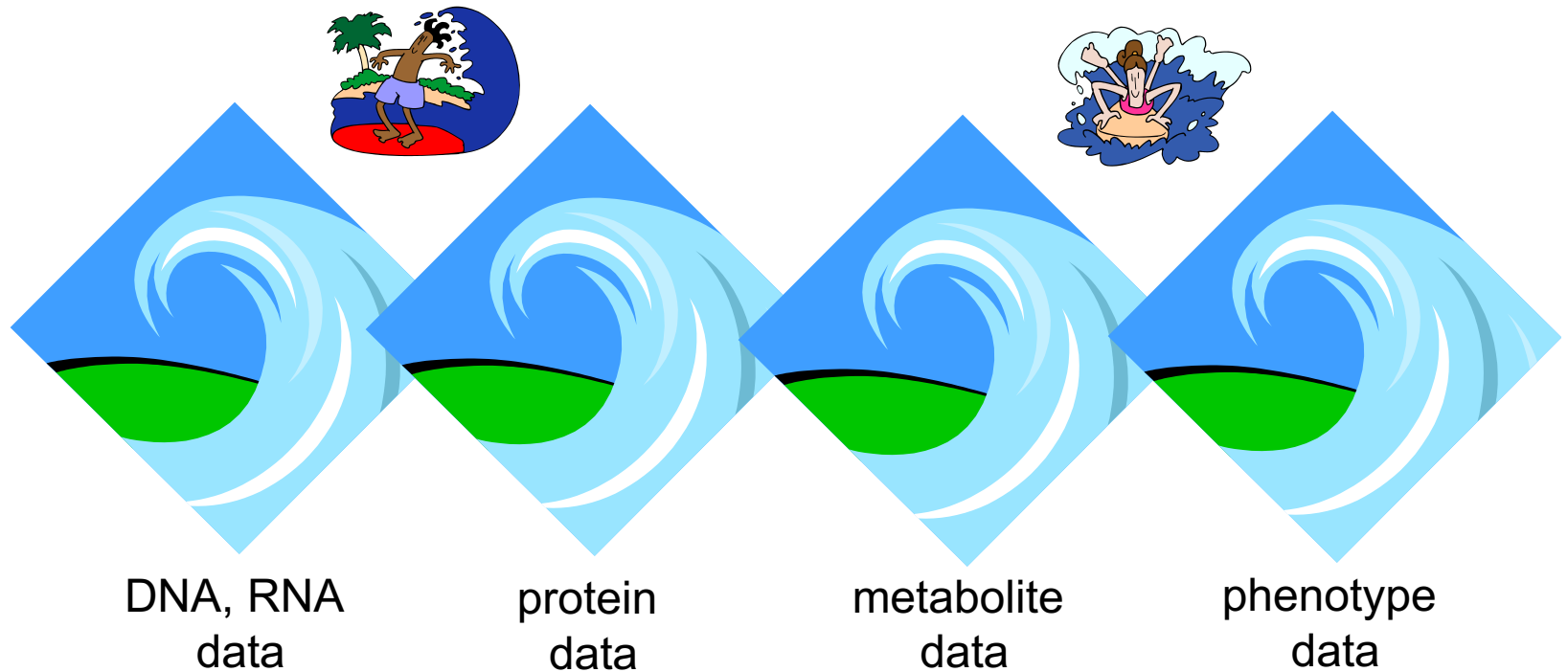      - Primary sequence
      - 3D structure
      - Subcellular localization
      - Rate of degradation
      - Enzymatic activity properties
      - Post-translational modification
        - Phosphorylation
        - Prenylation
        - Methylation
        - Ubiquitination

Is it static or dynamic?

What stimuli cause it to change?

By how much?

# The waves of data keep mounting!



DNA, RNA data     protein data     metabolite data     phenotype data

# Exploring the sea of biological data

- Primary data source
  - Articles published in peer-reviewed journals
  - Over 18 million available through PubMed by 2008!
    - NOT a comprehensive set; many journals are missing

- Answering scientific questions
  - Specific focus:
    - Find *every single piece of information ever discovered* about my favorite gene – XYZ1 – to figure out exactly what it does
  - Broad search:
    - Compare the protein sequence of *every single transcription factor ever discovered* in a prokaryote or a eukaryote to study the evolution of nuclear-localization signals
  - *How do you collect these data from ALL of the relevant research articles?*
    - Data repositories . . . staffed by **biocurators** . . . try to help!
      - Computer scientists and bioinformaticians contribute to these efforts as well!

# Global data repositories / databases

- Centralized data hubs
  - Many data types
  - Many species

  - Asia
    - Several in Japan, e.g. RIKEN, China is adding new ones
  - Europe
    - European Bioinformatics Institute (EBI)
  - USA
    - National Center for Biotechnology Information (NCBI)

# Global data re

# Specialized data repositories / databases

- Model organism databases (MODs)
  - Mouse Genome Informatics (MGI)
  - Flybase (*Drosophila*)
  - Saccharomyces Genome Database (SGD) (yeast)
  - **The Arabidopsis Information Resource (TAIR)**

- Topical databases
  - Worldwide Protein Data Bank (3D structures)
  - miRbase (microRNAs)
  - **Plant Metabolic Network (PMN) (metabolic / biochemical pathways)**

# Roles of biocurators at data repositories

- Organize and process raw data
  - Assign unique stable identifiers for nucleotide sequences submitted by researchers

- Review and improve data to generate **curated** data sets
  - Manually correct errors in raw nucleotide sequences to make RefSeq gene structures

- Develop tools for accessing data
  - Provide a protein interaction viewer

- Train users
  - Present at conferences and universities

- ***Try to help researchers harness the data explosion!***
  - ***TAIR***
  - ***Plant Metabolic Network***

# Introduction to TAIR

- TAIR = **T**he **A**rabidopsis **I**nformation **R**esource

- Why Arabidopsis?

- What does TAIR do?

- What can you do with TAIR?

Arabidopsis

# Introduction to Arabidopsis

- Basic
  - "s...
  - als...
  - ca...
  - an...
  - m...

  - fo...

- Why ...

# Arabidopsis offers some advantages

- "Good" genome
  - very small: 125 Mb - ~27,000 genes
  - diploid
  - 5 haploid chromosomes
  - fewer/smaller regions of repetitive DNA than many plants

- Quite *easily* transformable with *Agrobacterium*
  - NO tissue culture required

- Inertia!
  - A group of scientists lobbied for Arabidopsis
  - The genome was sequenced (2000)
  - **MANY resources have been developed**

# Arabidopsis research can be applied to "real plants"

- Over-expression of the *hardy* gene from Arabidopsis can improve water use efficiency in **<u>rice</u>** (*Karaba 2007)*

- A high throughput screen performed using **<u>castor bean</u>** cDNAs expressed in Arabidopsis found three cDNAs that increase hydroxy fatty acid levels in seeds (*Lu 2006*)

- These experiments and many more benefit from the work of curators trying to help harness the Arabidopsis data explosion . . .

    - ~2400 articles discussing Arabidopsis in PubMed per year!

# What

- Curato... at directors

- TAIR ...

- TAIR ...

- TAIR ...

  - Fu...
  - St...
  - Av...

The Arabidopsis Information Resource

The Arabidopsis Information Resource (TAIR) maintains a database of genetic and molecular biology data for the model higher plant Arabidopsis thaliana. Data available from TAIR includes the complete genome sequence along with gene structure, gene product information, metabolism, gene expression, DNA and seed stocks, genome maps, genetic and physical markers, publications, and information about the Arabidopsis research community. Gene product function data is updated every two weeks from the latest published research literature and community data submissions. Gene structures are updated 1-2 times per year using computational and manual methods as well as community submissions of new and updated genes. TAIR also provides extensive linkouts from our data pages to other Arabidopsis resources.

The Arabidopsis Biological Resource Center at The Ohio State University collects, reproduces, preserves and distributes seed and DNA resources of Arabidopsis thaliana and related species. Stock information and ordering for the ABRC are fully integrated into TAIR.

CARNEGIE SCIENCE — TAIR is located at the Carnegie Institution for Science Department of Plant Biology and funded by the National Science Foundation.

Synteny Viewer at TAIR

A synteny viewer (gbrowse_syn, a GMOD project) is now available at TAIR. This tool allows the user to compare syntenic regions between A. thaliana and A. lyrata. Additional plant genomes will be added in the future.

# Structural curation at TAIR

- Structural curators try to answer the question:

  *What are ALL of the genes in Arabidopsis?*

  - Use many types of data
    - ESTs
    - full-length cDNAs
    - peptides
    - orthology
    - RNASeq data**

  - Determine gene coordinates and features
    - Establish intron, exon, and UTR boundaries
    - Add alternative splice variants
    - Classify genes
      - protein coding
      - miRNA
      - pseudogene

# Structural curation at TAIR

- Even though the genome was sequenced in 2000 . . .
- . . . the work goes on!

  - TAIR9 – released June 2009
  - 282 new loci and 739 new splice variants

  - TAIR10 – on its way
  - 126 novel genes
  - 1182 updated genes
  - 5885 new splice variants added (18% of all loci)

# Structural curation at TAIR

□ Apollo is a program to assist with structural curation

# Functional ...AIR

- Functional cura... questions:
  - *What does e... idopsis do?*
  - *When and w...*
  - We hope tha... orm research in other plants

- Functional cura... *d vocabularies*
  - Allow cross...
  - TAIR curate... d agree upon common terms

**The seed-bearing str...**
**formed from the ova...**

| achene |
| berry |
| capsule |
| caryopsis |
| circumcissile |
| capsule |
| cypsela |
| drupe |
| follicle |
| grain |
| kernel |
| legume |
| loculicidal capsule |
| lomentum |
| nut |
| pod |
| pome |
| poricidal capsule |
| schizocarp |
| septicidal capsule |
| septifragal capsule |
| silique |

→ **FRUIT**

**Plant Ontology:**
*Structure:*
PO:0009001

# Functional curation at TAIR

**Catalysis of the reaction:**
**IAA + UDP-D-glucose = indole-3-acetyl-beta-1-D-glucose + UDP**

IAA-Glu synthetase activity
IAA-glucose synthase activity
IAGlu synthase activity
indol-3-ylacetylglucose synthase activity
UDP-glucose:(indol-3-yl)acetate beta-D-glucosyltransferase activity
UDP-glucose:indol-3-ylacetate glucosyl-transferase activity
UDP-glucose:indol-3-ylacetate glucosyltransferase activity
UDPG-indol-3-ylacetyl glucosyl transferase activity
UDPglucose:indole-3-acetate beta-D-glucosyltransferase activity
uridine diphosphoglucose-indoleacetate glucosyltransferase activity

**indole-3-acetate beta-glucosyltransferase activity**

**Gene Ontology:**
*Molecular function:*
GO:0047215

# Functional curation at TAIR



Gene

# Functional curation at TAIR



48 day old plant.
Image provided by Barry Pogson

*alx8*

1 cm

# Providing access to external tools and data

# Providing Tools at TAIR

# Providing Tools at TAIR

# **Providir**

- Tech team

  - Create

  - Create



**TAIR Gene Search**

[Help]

Genes may be searched by name, keywords, features, and/or location. In TAIR, a Gene Model is defined as any description of a gene product from a variety of sources including computational prediction, mRNA sequencing, or genetic characterization. A locus is defined as the genomic sequence corresponding to a transcribed unit (e.g. AT2G03340) in the genome. In TAIR, many gene models can exist for a given locus, therefore a search for a gene may result in multiple hits for the same gene name.

[ Reset ] [ Submit Query ]

### Search by Name or Phenotype @

| Gene name | ▼ | starts with | ▼ | |

(leaving the input box blank will retrieve all entries)

Include obsoleted genes ☐

### Search by Associated Keyword @

| Keyword Term@ | | starts with ▼ | |
| GO:PO ID (exact match help) | | | |
| Keyword Type | | Any | |
| | | GO Molecular Function | |
| | | GO Biological Process | |
| | | GO Cellular Component | |
| Evidence@ | | Any | |
| | | inferred from direct assay | |
| | | inferred from electronic annotation | |
| | | inferred from expression pattern | |

### Restrict by Features @

| Gene Model Type @ | | Any | |
| | | pre tma | |
| | | transposable element gene | |
| | | protein coding | |
| Advanced | | gene structure predicted | |
| | | has associated literature | |
| | | is sequenced | |
| | | is not sequenced | |
| Time Restriction @ | | ☐ only search last 2 months | |

*tair*

Search

# Other Resources at TAIR



TAIR Community Detail [Help]

| | |
|---|---|
| **Name** | Eleanore Wurtzel |
| **TAIR Accession** | Person:1501423405 |
| **Organisms** | Rice, Maize, Bacteria, Tomato, Wheat, Arabidopsis |
| **Primary Job Title** | Professor |
| **Research Interest** | regulation of carotenoid/provitamin A biosynthesis in cereal crops; evolution of biosynthetic pathways |
| **Keywords** | carotenoids, molecular biology, provitamin A, genomics, genetics, genes, Rice, Maize, Bacteria, Tomato, Wheat, Arabidopsis |
| **Address** | Dept. Of Biological Sciences<br>Lehman College, The City University of New York<br>250 Bedford Park Blvd. West<br>Bronx, NY 10468<br>USA |
| **E-mail** | wurtzel@lehman.cuny.edu |
| **Websites** | http://maize.lehman.cuny.edu |
| **Office Phone** | 718-960-8643 |
| **Lab Phone** | 718-960-4994 |
| **Mobile Phone** | 516-381-5013 |
| **Affiliations** | Organization          Job Title<br>*Eleanore Wurtzel Laboratory<br>* denotes Primary Investigator for this organization |
| **Record last updated** | 10/31/2006 |

| Detail Level: City ⌄ | Visits ↓ | Pages/Visit | Avg. Time on Site | % New Visits | Bounce Rate |
|---|---|---|---|---|---|
| 1. Ithaca | 675 | 6.30 | 00:07:13 | 10.22% | 25.19% |
| 2. New York | 527 | 5.75 | 00:08:14 | 17.08% | 33.02% |
| 3. Cold Spring Harbor | 248 | 6.48 | 00:10:01 | 14.94% | 19.09% |
| 4. Syracuse | 199 | 7.01 | 00:08:48 | 14.00% | 29.00% |
| 5. Upton | 90 | 4.88 | 00:05:43 | 24.44% | 21.11% |
| 6. Stony Brook | 71 | 2.28 | 00:01:07 | 4.23% | 16.90% |
| 7. Bronx | 62 | 7.39 | 00:07:22 | 29.03% | 41.94% |
| 8. Cortland | 38 | 2.53 | 00:02:30 | 0.00% | 84.21% |
| 9. Huntington Station | 36 | 9.69 | 00:22:34 | 0.00% | 19.44% |
| 10. Briarcliff Manor | 18 | 1.22 | 00:00:04 | 0.00% | 94.44% |

| | | | | |
|---|---|---|---|---|
| 5. | Uni | | | 39.59% |
| 6. | Sou | | | 29.02% |
| 7. | Car | | | 33.97% |
| 8. | Fra | | | 31.48% |
| 9. | Ind | | | 40.72% |
| 10. | Au | | | 29.17% |

Visits
1

Visits
1



Bronx
Visits: 62

134,797

**This state sent 2,056 visits via 74 cities**

# How can TAIR contribute to your work?

- If you work on Arabidopsis . . .
  - Find specific information about individual genes and proteins
  - Access large Arabidopsis-specific data sets

- If you work on another species . . .
  - Take your gene / protein of interest and find all the data TAIR contains for its ortholog
  - Look up your favorite:
    - biological process
    - molecular function
    - subcellular compartment
    - organ or tissue
    - developmental stage
    - mutant phenotype

    - Indentify many related genes in TAIR and then find orthologs in your species

- But . . . if you want more on plant metabolism . . . .

# Welcome to the PMN!

- PMN = The **P**lant **M**etabolic **N**etwork
  - Created in 2008
  - Funded by the National Science Foundation
- What is
- What da
- How do
- How ca
- How can you help the PMN to grow?

Sue Rhee
(PI)

Peifen Zhang
(Director)

# What is the PMN?



- Facilitate research that benefits society

# Connecting the PMN to important research efforts

- **More nutritious foods**
  - vitamin A biosynthesis, folate biosynthesis . . .

- **Medicines**
  - morphine biosynthesis, taxol biosynthesis . . .

- **More pest-resistant plants**
  - maackiain biosynthesis, capsidiol biosynthesis  . . .

- **Higher photosynthetic capacity and yield in crops**
  - chlorophyll biosynthesis, Calvin cycle . . .

- **Better biofuel feedstocks**
  - cellulose biosynthesis, lignin biosynthesis . . .

- **Many additional applications relevant to rational metabolic engineering**
  - ethylene biosynthesis, resveratrol biosynthesis . . .

# What data are in the PMN?



**Pathway Tools software provided by collaborators at SRI International**

# PMN databases

- Current PMN databases: **PlantCyc, AraCyc, PoplarCyc**
  - Coming soon: databases for wine grape, maize, cassava, Selaginella, and more . . .

- Other plant databases accessible from the PMN:

| PGDB | Plant | Source | Status |
|---|---|---|---|
| RiceCyc ** | Rice | Gramene | some curation |
| SorghumCyc | Sorghum | Gramene | no curation |
| MedicCyc ** | Medicago | Noble Foundation | some curation |
| LycoCyc ** | Tomato | Sol Genomics Network | some curation |
| PotatoCyc | Potato | Sol Genomics Network | no curation |
| CapCyc | Pepper | Sol Genomics Network | no curation |
| NicotianaCyc | Tobacco | Sol Genomics Network | no curation |
| PetuniaCyc | Petunia | Sol Genomics Network | no curation |
| CoffeaCyc | Coffee | Sol Genomics Network | no curation |

** Significant numbers of genes from these databases have been integrated into PlantCyc

# PMN database content statistics

|  | PlantCyc 4.0 |
| --- | --- |
| Pathways | 685 |
| Enzymes | 11058 |
| Reactions | 2929 |
| Compounds | 2966 |
| Organisms | 343 |

# How does experimentally verified data enter the PMN?

- Biocurators perform manual curation

  - Use journal articles to enter information

  - Receive helpful messages from researchers

  - Request specific data from experts

  - Invite editorial board members to review metabolic domains

# Pathway information

# Pathway information

# Compound information



AraCyc Pathway: choline biosynthesis III

Compound

**Compound: CDP-choline**



Synonyms: citicoline , citicholine , cidifos , cyticholine , cytidine 5'-diphosphocholine , cytidine diphosphate choline

Superclasses: a nucleic acid component -> a base derivative
a nucleic acid component -> a pyrimidine-related compound

Empirical Formula: $C_{14}H_{25}N_4O_{11}P_2$

Molecular Weight: 489.34 daltons

Smiles: C(OP(O)(=O)OP(O)(=O)OCC(N+)(C)(C)C)C1OC(C(O)C1O)n2cc(=O)nc(N)cc2))

Unification Links: CAS:987-78-0

Gibbs Energy of Formation (kcal/mol, estimated): -116.7

In Pathway Reactions as a Reactant:

phospholipid biosynthesis:
a 1,2-diacylglycerol + CDP-choline = a phosphatidylcholine + CMP

choline biosynthesis III:
a 1,2-diacylglycerol + CDP-choline = a phosphatidylcholine + CMP

In Pathway Reactions as a Product:

phospholipid biosynthesis:
phosphoryl-choline + CTP = CDP-choline + diphosphate

choline biosynthesis III:
phosphoryl-choline + CTP = CDP-choline + diphosphate

Synonyms

Classification(s)

Molecular Weight / Formula

Appears as Reactant

Appears as Product

# Enzyme information

# Enzyme information

*Arabidopsis* **Enzyme: phosphatidyltransferase**

# How does computationally predicted data enter the PMN?

**ANNOTATED GENOME**

*Phaseolus vulgaris*

**PlantCyc / MetaCyc**

**DNA sequences**

↓

**Predicted proteins**
Pv1234.56.a

↓

**Predicted functions**
chorismate mutase

**PathoLogic**

**Single species database**

chorismate
mutase
5.4.99.5

chorismate → prephenate

chorismate mutase
Pv1234.56.a

prephenate
aminotransferase

arogenate
dehydratase

91

L-phenylalanine

**PhaseolusCyc**

**+ validation**

# How can researchers use the PMN?

- Learn background information about particular metabolic pathways
    - Utilize simple and advanced search tools

    - Quick search bar



    - Specific search menus

# How can researchers use the PMN?

# How can researchers use the PMN?

- Compare metabolism across species

# How can researchers us

□

Pathway: abscisic acid biosynthesis
Compound: violaxanthin

violaxanthin
|
*trans*-neoxanthin
|
9'-*cis*-neoxanthin

CCD1   NCED9 7.0

NCED2 4.0   NCED6 2.0

NCED3 6.0   AT2G44990

AT1G30100 2.0

xanthoxin

ABA2 5.0

abscisic aldehyde

AT3G43600   AO1

AAO4   AAO3 7.0

(+)-abscisate

Printable version of this pathway diagram

# How will the PMN grow in the future?

- ❑  Help from the research community!!!

- ❑  You are the experts with great knowledge to share!

# Building better databases together

- To submit data, report an error, or volunteer to help validate . . .

  - Send an e-mail: **curator@plantcyc.org**

  - Use data submission "tools"



  - **Meet with me this afternoon**
    - . . . or later this week
    - . . . or later this year

# Building better databases together

- **Details are very, very welcome!!**
  - Reactions:
    - All co-factors, co-substrates, etc.
    - EC suggestions – partial or full

  - Compounds
    - Structure – visual representation / compound file (e.g. mol file)
    - Synonyms
    - Unique IDs (e.g. ChEBI, CAS, KEGG)

  - Enzymes
    - Unique IDs (e.g. At2g46480, UniProt, Genbank)
    - Specific reactions catalyzed

# Community gratitude

□ We

# Biological networking . . .

- Please use our data

- Please use our tools

- Please help us to improve our databases!

- Please contact us if we can be of any help!



curator@arabidopsis.org

www.arabidopsis.org

curator@plantcyc.org

www.plantcyc.org

# TAIR and PMN Acknowledgements

Sue Rhee *(PI - PMN)*
Eva Huala (*PI-TAIR*)
Peifen Zhang *(Director-PMN)*

Current Curators:
- -Tanya Berardini (*lead curator*)
- Philippe Lamesch (*lead curator)*
- Donghui Li (*curator*)
- Dave Swarbreck (*former lead curator*)

- Debbie Alexander (*curator*)
- A. S. Karthikeyan (*curator*)
- Marga Garcia (*curator)*
- Leonore Reiser

Current Tech Team Members:
- Bob Muller (*Manager*)
- Larry Ploetz (*Sys. Administrator*)
- Anjo Chi
- Raymond Chetty
- Cynthia Lee
- Shanker Singh
- Chris Wilks

PMN project post-doc
- Lee Chae

PMN Collaborators:
- Peter Karp (SRI)
- Ron Caspi (SRI)
- Suzanne Paley (SRI)
- SRI Tech Team
- Lukas Mueller (SGN)
- Anuradha Pujar (SGN)
- Gramene and MedicCyc

National Science Foundation
WHERE DISCOVERIES BEGIN

CARNEGIE
INSTITUTION FOR
SCIENCE

# Biological networking . . .

- Please use our data

- Please use our tools

- Please help us to improve our databases!

- Please contact us if we can be of any help!

curator@arabidopsis.org

curator@plantcyc.org

www.arabidopsis.org

www.plantcyc.org

# Out-takes

□ The following slides are relevant but were removed from the presentation due to time constraints

curator@arabidopsis.org

curator@plantcyc.org

www.arabidopsis.org

www.plantcyc.org

# Arabidopsis has good model organism traits

- Fast life cycle (6 weeks)
- Thousands of plants fit in a small space
- Fairly easy to grow
- Thousands of seeds produced by each plant
- Self-fertile (in-breeding)
- Many different subspecies/ecotypes
- Serves as a good model for **crop plants**

- But why Arabidopsis instead of **other plants**?

# Arabidopsis data explosion

- TONS of data are generated about Arabidopsis

  - Over 2400 "Arabidopsis" articles published each year are indexed in PubMed

  - Tens of thousands of mutants have been generated

  - Hundreds of microarray experiments have been performed

  - Proteomics and metabolomics studies are becoming popular

  - "1001" Arabidopsis genomes are being sequenced

  - Large-scale phenotypic studies are scheduled to start soon

- TAIR tries to bring data together to benefit scientists and society

# Pr

T

# What data are in the PMN?

- Plants provide crucial benefits to the ecosystem and humanity

- A better understanding of plant metabolism may contribute to:

  - More nutritious foods
  - New medicines
  - More pest-resistant plants
  - Higher photosynthetic capacity and yield in crops
  - Better biofuel feedstocks
  - Improved industrial inputs (e.g. oils, fibers, etc.)
  - Enhanced ability to do rational metabolic engineering
  - . . . many more applications

- How can the PMN help?

# What metabolites are in the PMN?

- "Primary" metabolites ("essential")

  - sugars
    - glucose, fructose, . . .

  - amino acids
    - tryptophan, glutamine, . . .

  - lipids
    - waxes, phosphatidylcholine , . . .

  - vitamins
    - A, E, K, C, thiamine, niacin, . . .

  - hormones
    - auxin, brassinosteroids, ethylene . . .

# What metabolites are in the PMN?

- "Secondary" metabolites (important, but not "essential")

  - terpenoids
    - orzyalexin, menthol, . . .

  - organosulfur compounds
    - glucosinolates, camalexin . . .

  - isoflavonoids
    - glyceollin, daidzein. . .

  - alkaloids
    - caffeine, capsaicin, . . .

  - polyketides
    - aloesone, . . .

  - many more . . .

# How do computational predictions enter the PMN?

- **New sets of DNA sequences -> predicted proteome**
    - Genomes are sequenced
    - Large RNAseq or EST data sets are created

- **Predicted proteome -> set of predicted enzyme functions**
    - Performed using computer algorithms
    - The PMN is working to develop better algorithms to increase the accuracy of the predictions

- **Set of predicted enzyme functions -> set of predicted metabolic pathways**
    - The PathoLogic program uses a reference database to predict the metabolic pathways for the enzyme sets

- **Set of predicted metabolic pathways -> set of "validated" metabolic pathways**
    - Curators remove incorrect information and add additional data