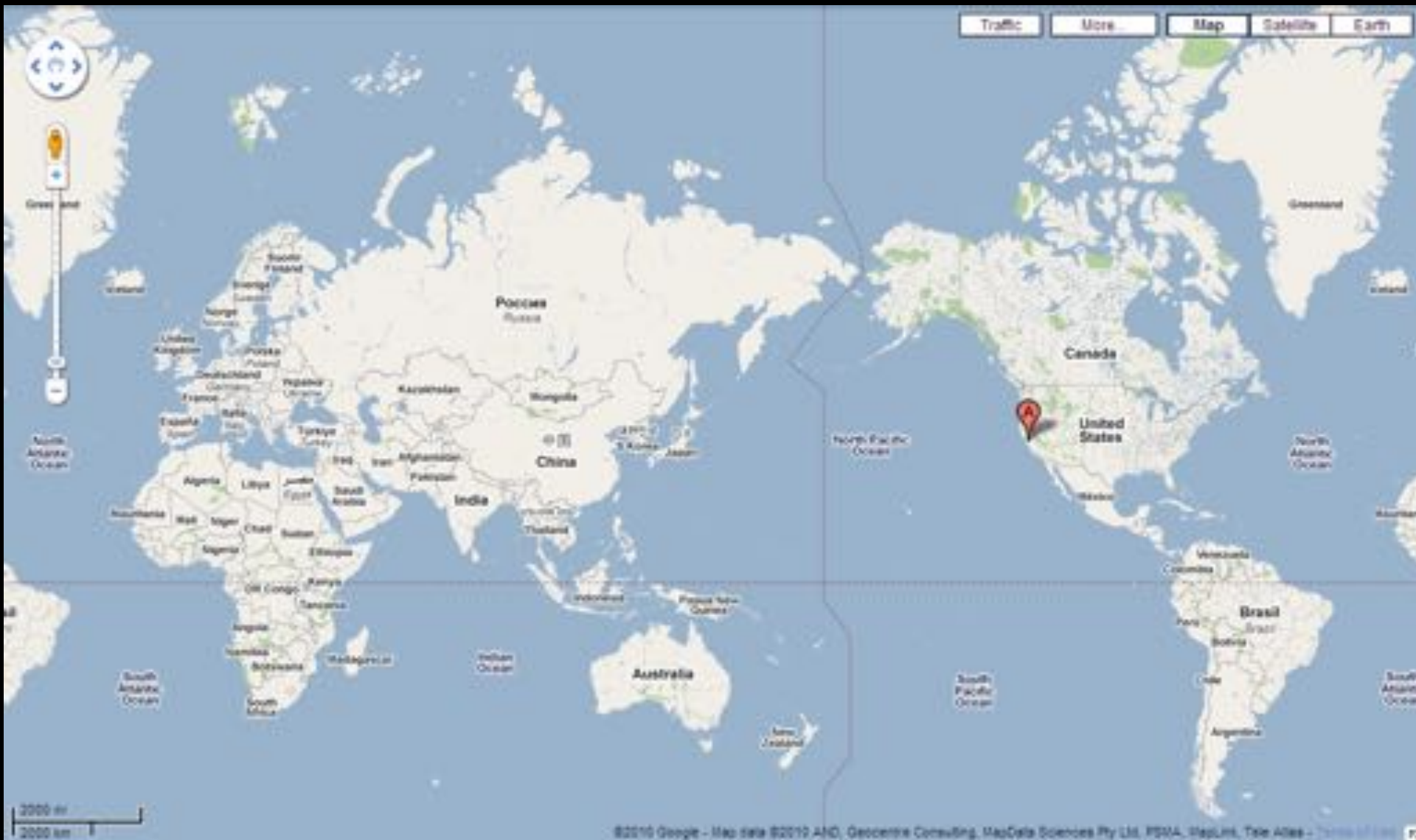# Combining Computational Prediction and Manual Curation to Create Plant Metabolic Pathway Databases

Peifen Zhang

Carnegie Institution For Science

Department of Plant Biology

# Where We Are

# Who We Are

PMN:

- Sue Rhee (*PI*)

- Kate Dreher (*curator*)

- Lee Chae (Postdoc)

- Anjo Chi (*programmer*)

- Cynthia Lee (TAIR *tech team*)

- Larry Ploetz (TAIR *tech team*)

- Shanker Singh (TAIR *tech team*)

- Bob Muller (TAIR *tech team*)

Key Collaborators:

- Peter Karp (MetaCyc, SRI)

- Ron Caspi (MetaCYc, SRI)

- Lukas Mueller (SGN)

- Anuradha Pujar (SGN)

National Science Foundation
WHERE DISCOVERIES BEGIN

http://plantcyc.org

# Outline

- Introduction and database snapshot
- Pathway database creation pipeline
- Manual curation
- Future work

# Introduction

- Background and rationale
  - Plants (food, feed, forest, medicine, biofuel…)
  - An ocean of sequences
    - More than 60 species in genome sequencing projects, hundreds in EST projects
  - Putting individual genes onto a network of metabolic reactions and pathways
    - Annotating, visualizing and analyzing at system level
  - AraCyc (Arabidopsis thaliana, TAIR/PMN)
    - predicted by using the Pathway Tools software, followed by manual curation

# Browsing Pathways

Pathways
- Biosynthesis (613 instances)
  - Amines and Polyamines Biosynthesis (14 instances)
  - Amino acids Biosynthesis (52 instances)
  - Aminoac...
  - Aromatic...
  - Carbohyd...
  - Cell struc...
  - Cofactor...
  - Fatty Ac...
  - Hormone...
  - Metaboli...
  - Nucleosi...
  - Other Bi...
  - Seconda...
  - Sideroph...
- Degradation...
- Detoxificati...
- Generation...
- Superpathw...
- Transport P...

## *PlantCyc* Pathways Class: Alkaloids Biosynthesis

Summary:
This class contains biosynthetic pathways of alkaloids. Most alkaloids contain cyclic nitrogen. T function as defense compounds. Many alkaloids, including morphine and cocaine, have a high receptors of neurotransmitters and have pharmacological activities.

Parent Classes:
Nitrogen-Containing Secondary Compounds Biosynthesis

Child Classes:
Betalaine alkaloids (8) ,
Indole alkaloids (5) ,
Isoquinoline and Benzylisoquinoline alkaloids (9) ,
Peptide alkaloids (0) ,
Polyketide alkaloids (0) ,
Purine alkaloids (4) ,
Pyrrolidine, Piperidine and Pyridine alkaloids (3) ,
Pyrrolizidine alkaloids (0) ,
Quinoline alkaloids (2) ,
Quinolizidine alkaloids (1) ,
Terpenoid Alkaloids Biosynthesis (3) ,
Tropane alkaloids (4)

Instances:
(S)-reticuline biosynthesis I ,
berberine biosynthesis ,
bisbenzylisoquinoline alkaloid biosynthesis ,
dehydroscoulerine biosynthesis ,
laudanine biosynthesis ,
magnoflorine biosynthesis ,
morphine biosynthesis ,
palmatine biosynthesis ,
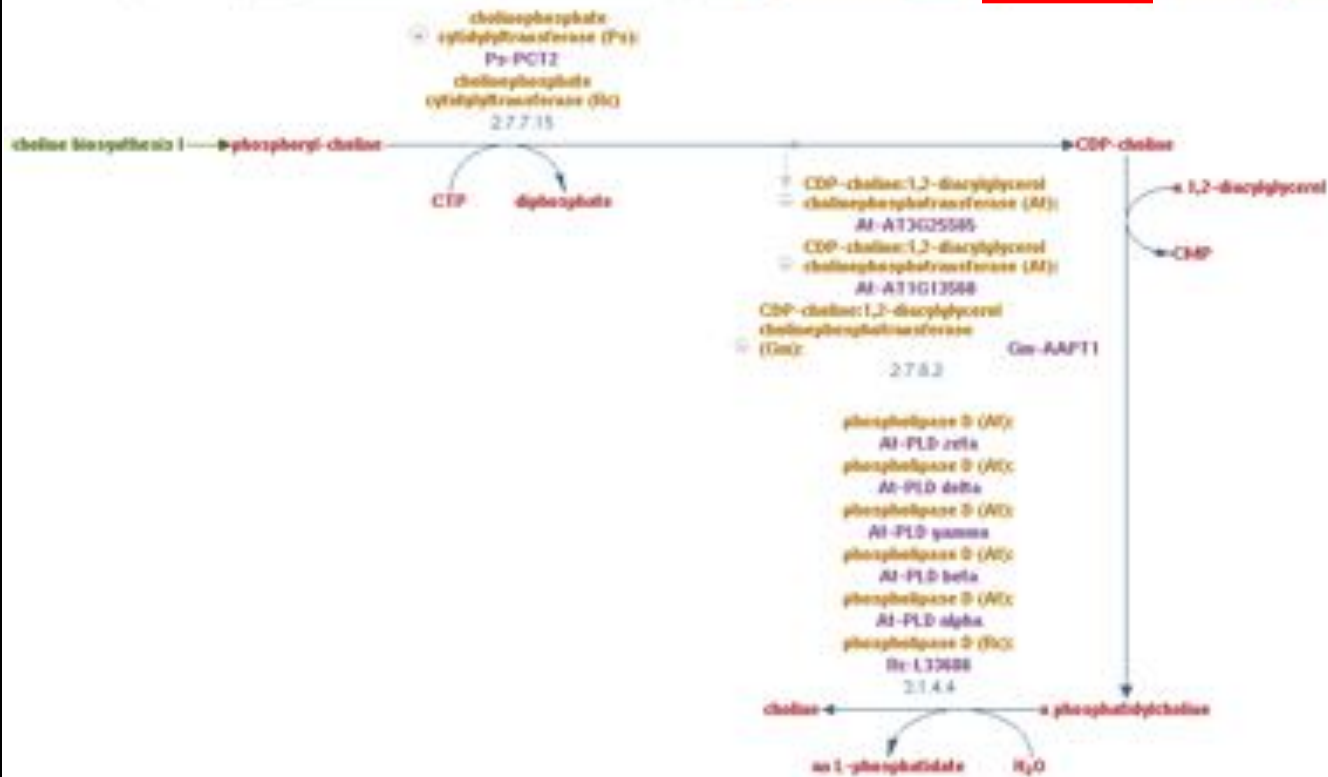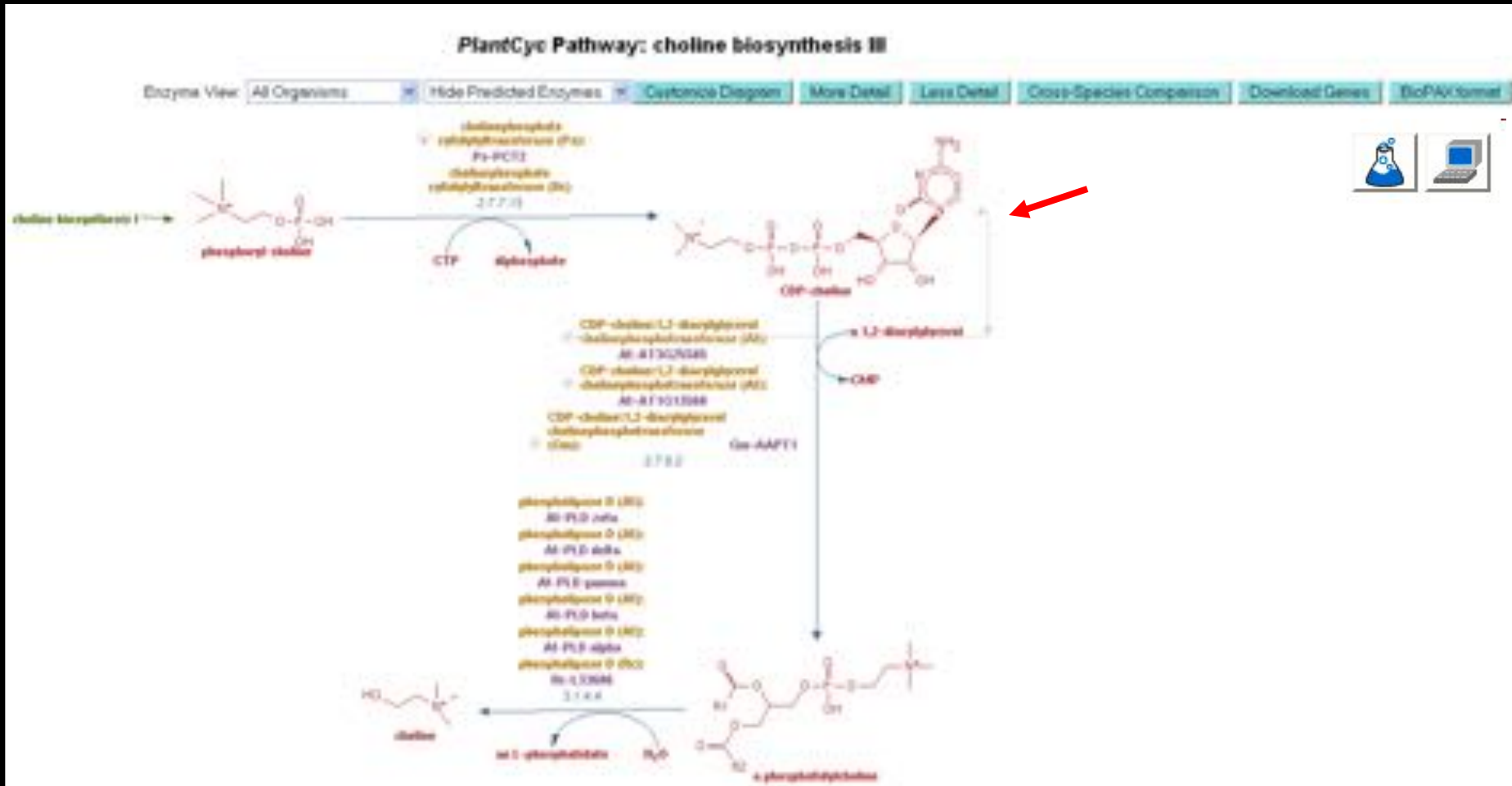sanguinarine and macarpine biosynthesis

# Searching Databases

# A Typical Pathway Detail Page

# Linking to Other Data Detail Pages

# Compound Detail Pages

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: $C_{14}H_{27}N_4O_{11}P_2$

Molecular Weight: 489.34 daltons

**Molecular Weight / Formula**

Smiles: C(OP(O)(=O)OP(O)(=O)OCC[N+](C)(C)C)C1OC(C(O)C1O)n2ccc(=O)nc(N)cc2))

**Smiles / InChI**

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

**Appears as Reactant**

In Pathway Reactions as a Product:

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

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

**Appears as Product**

# Enzyme Detail Pages

*Arabidopsis* **Enzyme: phosphatidyltransferase**
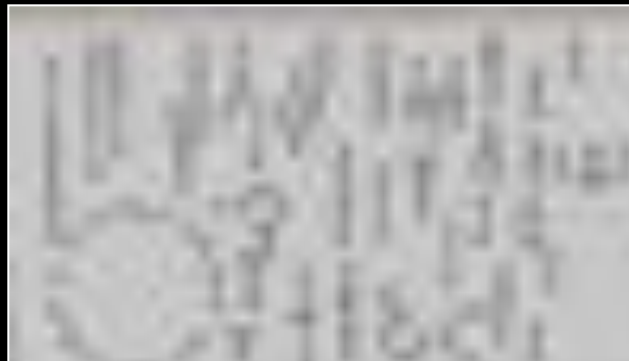


Evidence

Summary

Inhibitors, Kinetic Parameters, etc.

# Visualizing and Interpreting Omics Data in a Metabolic Context



- Gene expression data
- Proteomic data
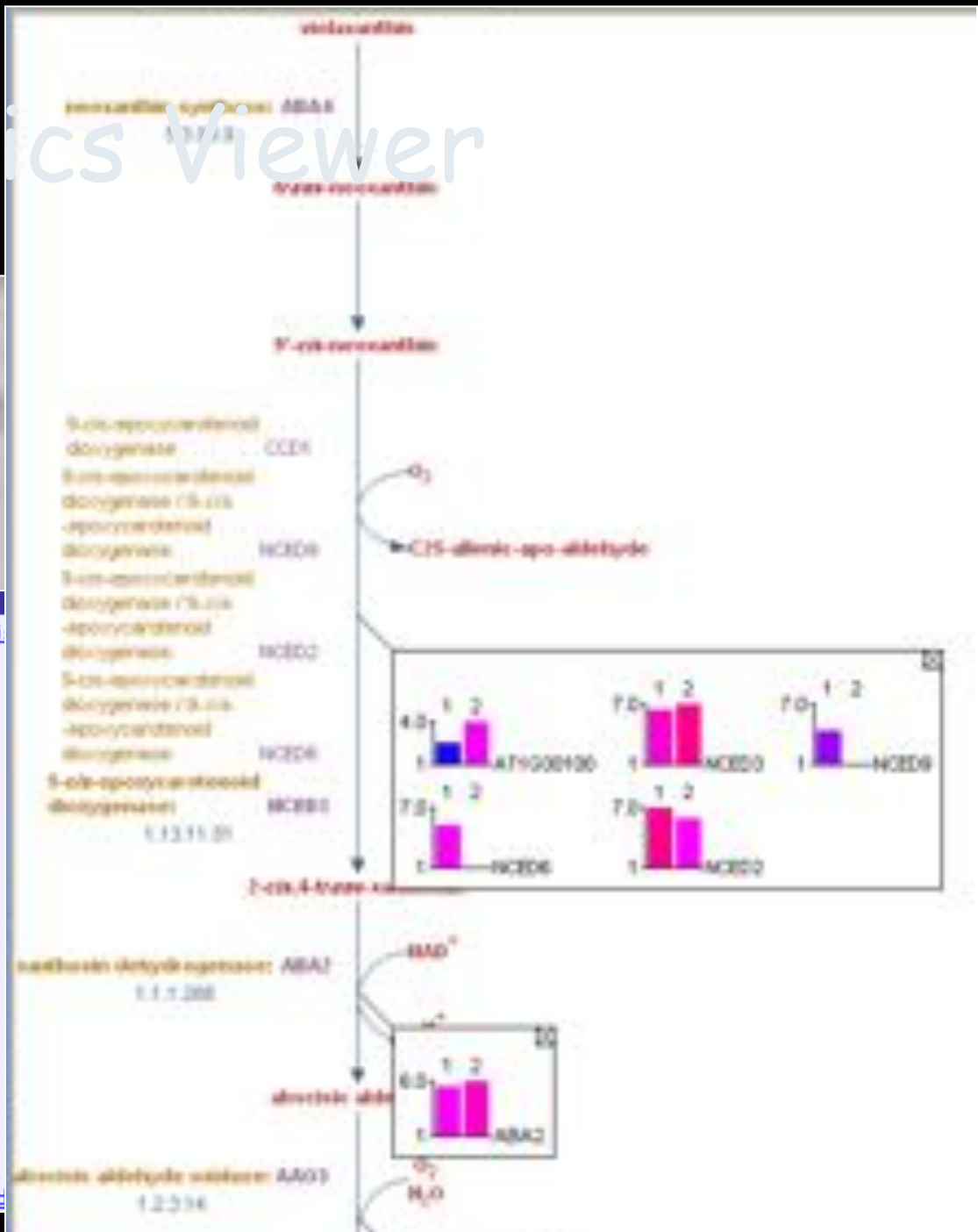- Metabolic profiling data
- Reaction flux data

# Omics Viewer

# Comparing Across Species



Compound: 3-phosphoglycerate
Pathway: serine biosynthesis

# Pathways

## Table 1: Breakdown of Pathways by Pathway Class

This table present:
down further to sh
assigned to more t
will see only those

Biosynthesis
- Amines and Poly
- Amino acids Bios
- Aminoacyl-tRNA
- Aromatic Compo
- Carbohydrates B
- Cell structures B
- Cofactors, Prost
- Fatty Acids and
- Hormones Biosy
- Metabolic Regul.
- Nucleosides and
- Other Biosynthe
- Polysaccharides
- Secondary Metal
- Secondary Metal
- Siderophore Bios

**Pathway Class:** Biosynthesis - Hormones Biosynthesis | AraCyc col | P. trichocarpa

cis-zeatin biosynth
ent-kaurene biosyr
trans-zeatin biosyr
abscisic acid biosyr
abscisic acid glucos
aldehyde oxidation
brassinosteroid bio
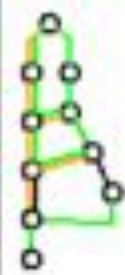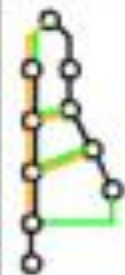brassinosteroid bio
brassinosteroid bio
cytokinins 7-N-gluc
cytokinins 9-N-gluc
cytokinins-O-glucos
ethylene biosynthe
GA$_{12}$ biosynthesis
gibberellin biosynth
gibberellin biosynth
gibberellin biosynth
hydroxyjasmonate
IAA biosynthesis I
IAA biosynthesis II
IAA biosynthesis V
indole-3-acetyl-ami
jasmonic acid biosy
jasmonoyl-amino a

| Organism | Evidence Glyph | Enzymes and Genes for brassinosteroid biosynthesis II | |
|---|---|---|---|
| AraCyc col | | EC# 5.3.3.1 | $\Delta^{5-3}$-ketosteroid isomerase: |
| | | EC# 1.1.1.145 | 3β-hydroxylsteroid dehydrogenase: |
| | | EC# 1.3.99.- | sterol 5-alpha reductase: DET2 |
| | | EC# 1.14.13.- | steroid 22α-hydroxylase: DWF4 |
| | | RXN-712 | None |
| | | EC# 1.14.13.- | steroid 22α-hydroxylase: DWF4 |
| | | EC# 1.14.13.- | steroid 22α-hydroxylase: DWF4 |
| | | EC# 1.14.13.- | steroid 22α-hydroxylase: DWF4 |
| | | RXN-4226 | SAX1: SAX1 |
| | | EC# 1.3.99.- | sterol 5-alpha reductase: DET2 |
| | | RXN-4228 | None |
| | | EC# 1.14.-.- | 23alpha hydroxylase / cathasterone 23α-hydroxylase: CBB3 |
| P. trichocarpa | | EC# 5.3.3.1 | None |
| | | EC# 1.1.1.145 | None |
| | | EC# 1.3.99.- | None |
| | | EC# 1.14.13.- | monooxygenase: JGI-820024 monooxygenase: JGI-796803 |
| | | RXN-712 | None |
| | | EC# 1.14.13.- | monooxygenase: JGI-820024 monooxygenase: JGI-796803 |
| | | EC# 1.14.13.- | monooxygenase: JGI-820024 monooxygenase: JGI-796803 |
| | | EC# 1.14.13.- | monooxygenase: JGI-820024 monooxygenase: JGI-796803 |

# Introduction (cont)

- Background and rationale
  - Plants (food, feed, forest, medicine, biofuel...)
  - An ocean of sequences
    - More than 60 species in genome sequencing projects, hundreds in EST projects
  - Putting individual genes onto a network of metabolic reactions and pathways
    - Annotating, visualizing and analyzing at system level
  - AraCyc (Arabidopsis thaliana, TAIR/PMN)
    - predicted by using the Pathway Tools software, followed by manual curation
  - Other plant pathway databases predicted by using the Pathway Tools
    - RiceCyc (Oryza sativa, Gramene)
    - MedicCyc (Medicago truncatula, Noble Foundation)
    - LycoCyc (Solanum lycopersicum, SGN), ...

# Pathway Prediction and Pathway Database Creation

- Infer the reactome of an organism from the enzymes present in its annotated genome
  - Mapping annotated enzyme sequences to reactions
- Infer the metabolic pathways of the organism from its reactome
  - Pathway-calling based on supporting evidence of reactions

# Limitations

- Creating pathway databases includes three major components, and is resource-intensive
  - Sequence annotation
  - Reference pathway database
  - Pathway prediction, validation, refinement
- Heterogeneous sequence annotation protocols and varying levels of pathway validation impact quality and hinder meaningful cross-species comparison
  - Using a non-plant reference database causes many false-positive and false-negative pathway predictions

# Introducing the PMN

- Scope
  - A platform for plant metabolic pathway database creation
  - A community for data curation
    - Curators, editorial board, ally in other databases, researchers

- Major goals
  - Create a plant-specific reference pathway database (PlantCyc)
  - Create an enzyme sequence annotation pipeline
  - Enhance pathway prediction by using PlantCyc, and including an automated initial validation step
  - Create metabolic pathway databases for plant species
    - e.g. PoplarCyc (Populus trichocarpa), SoyCyc (soybean)

# PlantCyc Creation

- Nature
  - Multiple-species, plants-only, curator-reviewed pathways, primary and secondary metabolism

- Major Source
  - All AraCyc pathways and enzymes
  - Plant pathways and enzymes from MetaCyc
  - Additional pathways and enzymes manually curated and added
  - Enzymes from RiceCyc, LycoCyc and MedicCyc

# PMN Database Content Statistics

| | PlantCyc 4.0 | AraCyc 7.0 | PoplarCyc 2.0 |
|---|---|---|---|
| Pathways | 685 | 369 | 288 |
| Enzymes | 11058 | 5506 | 3420 |
| Reactions | 2929 | 2418 | 1707 |
| Compounds | 2966 | 2719 | 1397 |
| Organisms | 343 | 1 | 1* |

Valuable <u>plant natural products</u>, many are specialized metabolites that are limited to a few species or genus.
- medicinal: e.g. artemisinin and quinine (treatment of malaria),
    codeine and morphine (pain-killer),
    ginsenosides (cardio-protectant),
    lupenol (antiinflammatory),
    taxol and vinblastine (anti-cancer)
- industrial materials: e.g. resin and rubber
- food flavor and scents: e.g. capsaicin and piperine (chili and pepper flavor), geranyl acetate (aroma of rose) and menthol (mint).

# Enzyme Sequence Annotation (version 1.0)

- Reference sequences, <u>enzymes</u> with <u>known</u> functions
  - 14,187 enzyme sequences compiled from UniProt, Brenda, MetaCyc, and TAIR
  - 3805 functional identifiers (full EC number, MetaCyc reaction id, GO id)
- Annotation methods
  - BLASTP
- Cut-off
  - unique e-value threshold for each functional identifier

# Enzyme Sequence Annotation (version 2.0)

- Reference sequences, <u>proteins</u> with <u>known</u> functions (ERL)
  - SwissProt
    - 117,000 proteins, 26,000 enzymes, 2,400 full EC numbers
  - Additional enzymes from MetaCyc, TAIR, Brenda and UniProt
  - Functional identifiers, full EC number, MetaCyc reaction id, GO id,
- Annotation methods
  - BLASTP
  - Priam (enzyme-specific, motif-based)
  - CatFam (enzyme-specific, motif-based)
- Function calling
  - Ensemble and voting

# Enzyme Sequence Annotation (version 2.0)



A) Individual classification of query sequence

B) Ensemble classification outputs final prediction

Lee Chae (unpublished)

# Automated Initial Pathway Validation

– Remove non-plant pathways
  • A list of 132 MetaCyc pathways

– Add universal plant pathways
  • A list of 115 PlantCyc pathways

# Manual Curation

- Who
  - Curators identify, read and enter information from published journal articles

- What
  - Remove false-positive pathway predictions
  - Remove false-positive enzyme annotations
  - Add missing pathways (pathway diagrams)
  - Add missing enzymes
  - Curate enzyme properties, kinetic data
  - Update existing pathways (pathway diagrams)
  - Add new reactions
  - Add new compounds and curate compound structures

# Conventions Used in Curation and Data Presentation

- A pathway, as drawn in the text books, is a functional unit, regulated as a unit
- Pathway displayed is expected to operate as such in the individual species listed

# *PlantCyc* Pathway: caffeine biosynthesis I

Enzyme View: No Enzymes ▾  [More Detail] [Less Detail] [Species Comparison]



If an enzyme name is shown in bold, there is experimental evidence for this enzymatic activity.

Superclasses: Biosynthesis -> Secondary Metabolites Biosynthesis -> Nitrogen-Containing Secondary Compounds Biosynthesis -> Alkaloids Biosynthesis -> Purine alkaloids -> Caffeine Biosynthesis

Species Data Available for: Camellia sinensis , Camellia sinensis assamica , Camellia taliensis , Coffea arabica , Coffea canephora

# Conventions Used in Curation and Data Presentation

- Pathway, as drawn in the text books, is a functional unit, regulated as a unit
- Pathway displayed is expected to operate as such in the individual species shown
- Alternative routes that have been observed in different organisms are curated separately as pathway <u>variants</u>

*PlantCyc* Pathway: caffeine biosynthesis I

Enzyme View: No Enzymes | More Detail | Less Detail | Species Comparison

2.1.1.158

salvage pathways of purine nucleosides II (plant) ——▶ xanthosine ——————————————————————————————————▶ 7-methylxanthosine

S-adenosyl-L-methionine    S-adenosyl-L-homocysteine
H⁺

3.2.2.25        H₂O

▶ D-ribose

2.1.1.160                                                    2.1.1.159

caffeine ◀——————————————— theobromine ◀——————————————————— 7-methylxanthine

H⁺        S-adenosyl-L-methionine     S-adenosyl-L-homocysteine     S-adenosyl-L-methionine
S-adenosyl-L-homocysteine                                      H⁺



*PlantCyc* Pathway: caffeine biosynthesis II (via paraxanthine)

Enzyme View: No Enzymes | More Detail | Less Detail | Species Comparison

2.1.1.158                                            3.2.2.25

xanthosine ————————————————————▶ 7-methylxanthosine —————————————————▶ 7-methylxanthine

S-adenosyl-L-methionine     S-adenosyl-L-homocysteine       H₂O        D-ribose
H⁺

S-adenosyl-L-methionine

S-adenosyl-L-homocysteine
H⁺

2.1.1.160

caffeine ◀——————————————————— paraxanthine

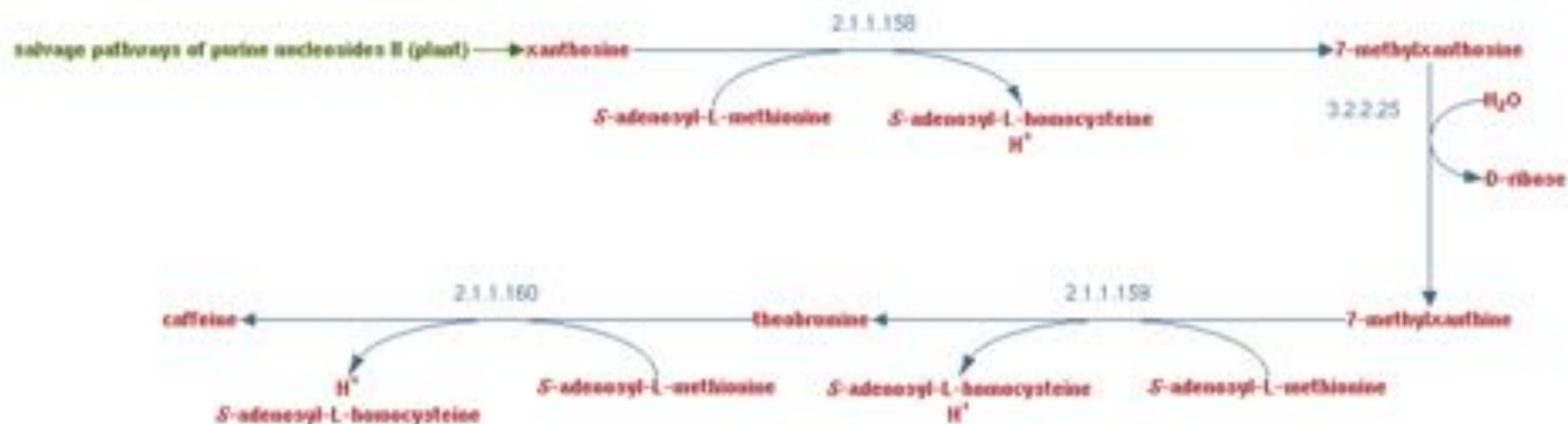H⁺        S-adenosyl-L-methionine
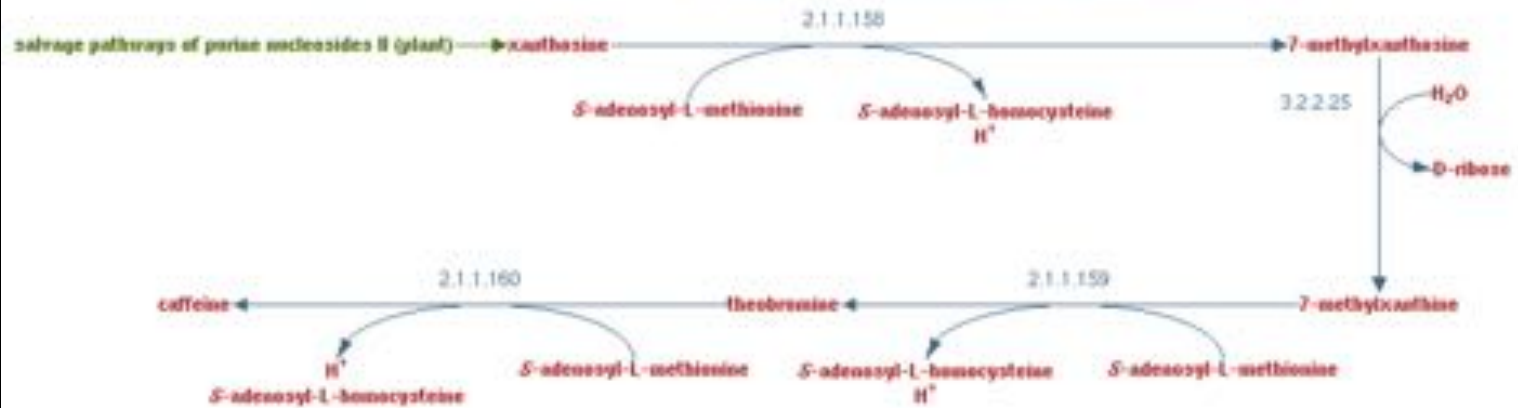S-adenosyl-L-homocysteine

# Conventions Used in Curation and Data Presentation

- Pathway, as drawn in the text books, is a functional unit, regulated as a unit
- Pathway displayed is expected to operate as such in the individual species shown
- Alternative routes that have been observed in different organisms are curated separately as pathway variants
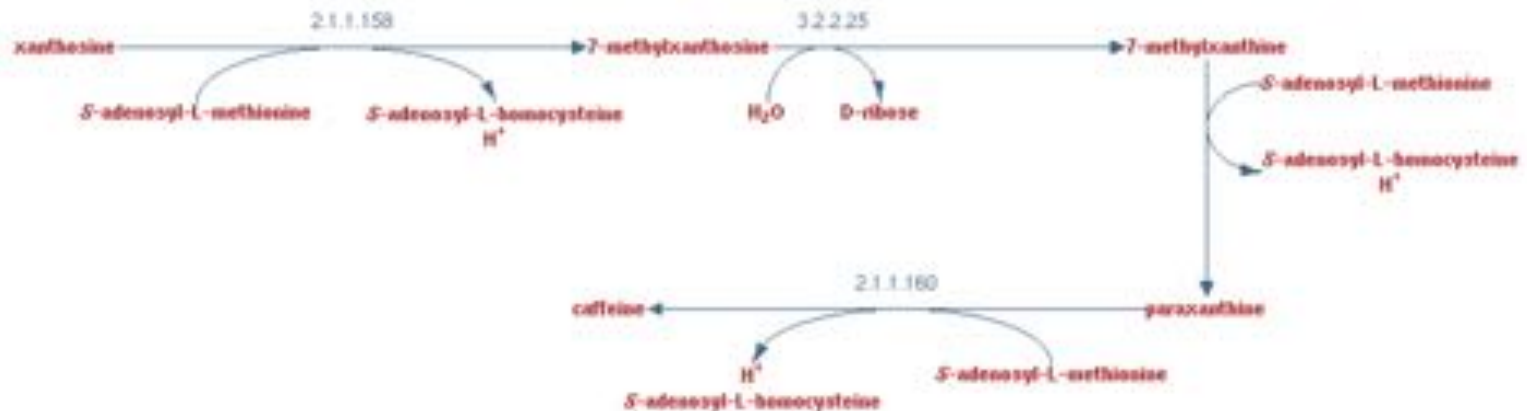- Mosaics combined of alternative routes from several different species are curated as <u>superpathways</u>
- Connected pathways, extended networks, are curated as <u>superpathways</u>

[ More Detail ]  [ Less Detail ]  [ Species Comparison ]

# Future Work

- Enhance pathway prediction and validation
  - Using additional evidence, such as presence of compounds, weighted confidence of enzyme annotations

- Refine pathways, hole-filling
  - Including non-sequence homology based information in enzyme function prediction, such as phylogenetic profiles, co-expression

- Add new data types, critical for strategic planning of metabolic engineering
  - Rate-limiting step
  - Transcriptional regulator

- Create new pathway databases
  - moss (P. patens), Selaginella, maize, cassava, wine grape …

Thank you!