# Pathway Bioinformatics

**Peter D. Karp, PhD**

**Bioinformatics Research Group**

**SRI International**

**Menlo Park, CA**

**pkarp@ai.sri.com**

**BioCyc.org**

# *Overview*

- **Definitions**

- **BioCyc collection of Pathway/Genome Databases**

- **Algorithms for pathway bioinformatics**

- **Pathway Tools software**
  - Navigation and analysis
  - Infer metabolic pathways from genomes
- **Pathway Tools ontology**

# Pathway Bioinformatics

- **The subfield of bioinformatics concerned with ontologies, algorithms, databases and visualizations of pathways**

- **Examples:**
  - Inference of metabolic pathways from genomes
  - Schemas for pathway DBs
  - Exchange formats for pathway data
  - Classification systems for pathway data
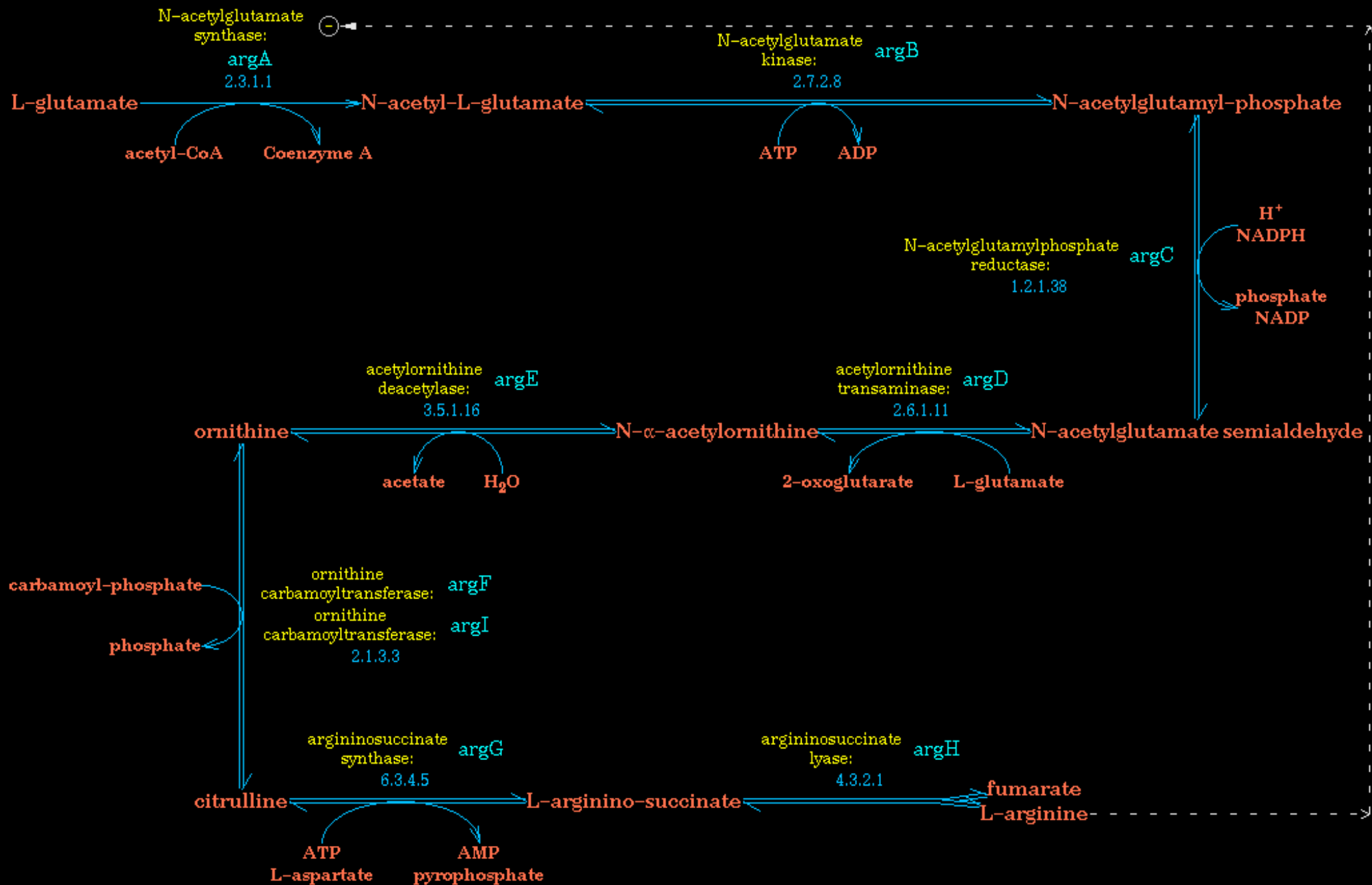  - Pathway diagram layout algorithms

# Definition of Metabolic Pathways

- **A chemical <u>reaction</u> interconverts chemical compounds (analogous to a production rule)**

$$A + B = C + D$$

- **An <u>enzyme</u> is a protein that accelerates chemical reactions.  Each enzyme is encoded by one or more genes.**

- **A <u>pathway</u> is a linked set of reactions (analogous to a chain of rules)**
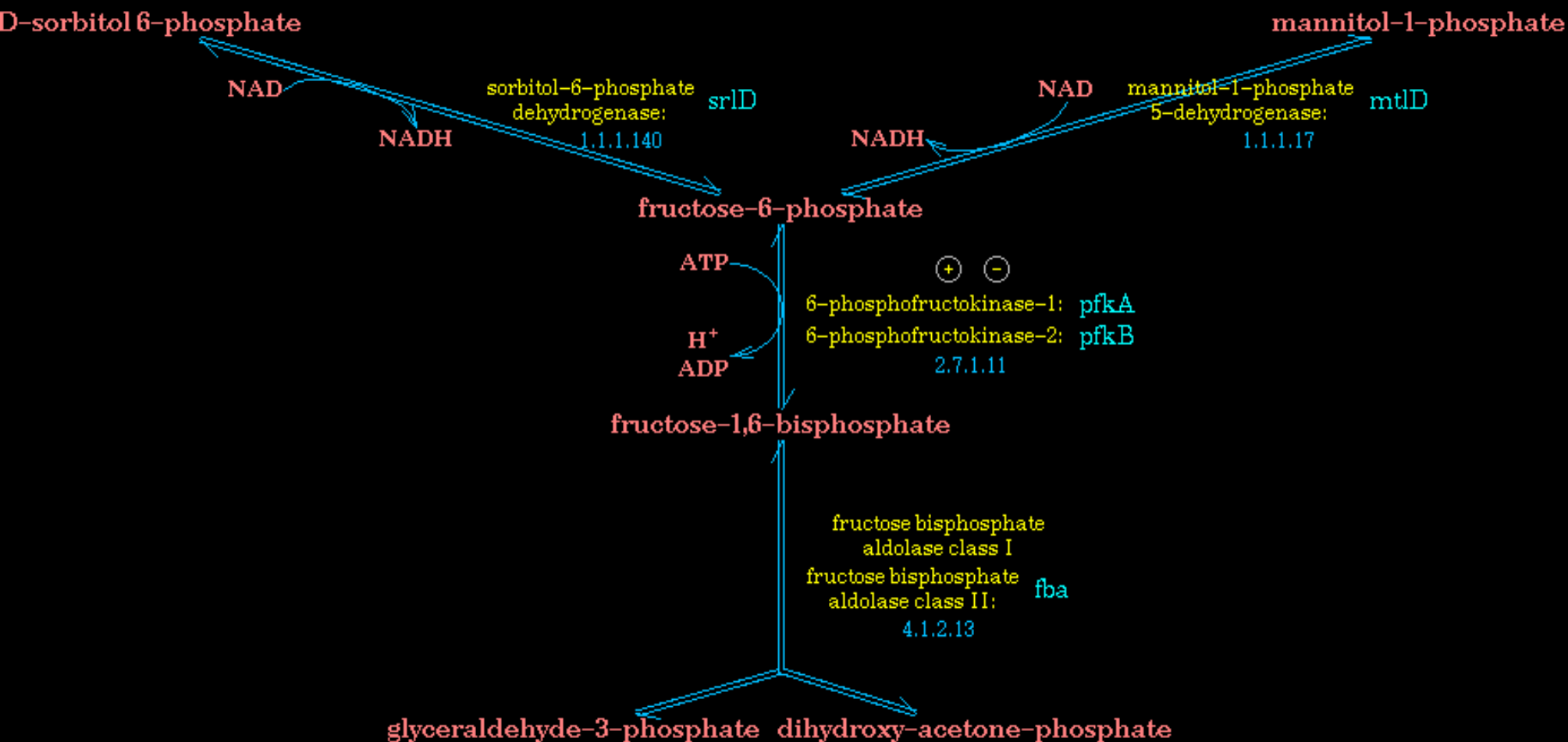
$$A \longrightarrow C \longrightarrow E$$

N-acetylglutamate
synthase:

argA
2.3.1.1

L-glutamate

acetyl-CoA    Coenzyme A

N-acetyl-L-glutamate

N-acetylglutamate
kinase:               argB

2.7.2.8

ATP    ADP

N-acetylglutamyl-phosphate

$H^+$
NADPH

N-acetylglutamylphosphate
reductase:           argC

1.2.1.38

phosphate
NADP

acetylornithine
deacetylase:       argE

3.5.1.16

ornithine

acetate    $H_2O$

N-α-acetylornithine

acetylornithine
transaminase:       argD

2.6.1.11

2-oxoglutarate    L-glutamate

N-acetylglutamate semialdehyde

carbamoyl-phosphate

ornithine
carbamoyltransferase:    argF

ornithine
carbamoyltransferase:    argI

2.1.3.3

phosphate

argininosuccinate
synthase:           argG

6.3.4.5

citrulline

ATP
L-aspartate

AMP
pyrophosphate

L-arginino-succinate

argininosuccinate
lyase:             argH

4.3.2.1

fumarate
L-arginine

# E. coli Pathway: hexitol degradation super-pathway

More Detail   Less Detail

D-sorbitol 6-phosphate

mannitol-1-phosphate

NAD

sorbitol-6-phosphate
dehydrogenase:   srlD
1.1.1.140

NADH

NAD

mannitol-1-phosphate
5-dehydrogenase:   mtlD
1.1.1.17

NADH

fructose-6-phosphate

ATP

$+$  $-$

6-phosphofructokinase-1:  pfkA
6-phosphofructokinase-2:  pfkB
2.7.1.11

$H^+$
ADP

fructose-1,6-bisphosphate

fructose bisphosphate
aldolase class I
fructose bisphosphate
aldolase class II:   fba
4.1.2.13

glyceraldehyde-3-phosphate   dihydroxy-acetone-phosphate

Superclasses: Carbon compounds, Super-Pathways

Subpathways: mannitol degradation, sorbitol degradation

Locations of Mapped Genes:

# *Mp. pneumoniae Pathway: pyrimidine ribonucleotide/ribonucleoside metabolism*

More Detail          Less Detail

cytosine

orotidine-5'-phosphate

$H_2O$

$CO_2$          4.1.1.23

$NH_3$          3.5.4.1

UMP

ATP

uridylate kinase: pyrH
2.7.4..

uracil
phosphoribosyltransferase:          upp
2.4.2.9

pyrophosphate

ADP

PRPP

GDP

uracil          uridine kinase:          udk
2.7.1.48

GTP

UDP

ribose-1-phosphate

ATP

Uridine phosphorylase
2.4.2.3

ADP

phosphate

uridine

UTP

$NH_3$

ATP
$H_2O$
L-glutamine

cytidine deaminase:          cdd
3.5.4.5

$H_2O$

L-glutamate
ADP
phosphate

6.

cytidine

CTP

GTP

ADP

GDP

ATP

CMP

CDP

ATP          ADP

cytidylate kinase:          cmk

# Definition of Small-Molecule Metabolism

- **Small-molecule metabolism**
  - Biochemical factory within the cell
  - Hundreds of enzyme-catalyzed reactions operating principally on small-molecule substrates

# *Small Molecule Metabolism*

# *What is a Metabolic Pathway?*

- **A pathway is a conceptual unit of the metabolism**
- **An ordered set of interconnected, directed biochemical reactions**
- **A pathway forms a coherent unit:**
  - Boundaries defined at high-connectivity substrates
  - Regulated as a single unit
  - Evolutionarily conserved across organisms as a single unit
  - Performs a single cellular function
  - Historically grouped together as a unit
  - All reactions in a single organism

# *BioCyc Collection of 507 Pathway/ Genome Databases*

- **Pathway/Genome Database (PGDB) – combines information about**
  - Pathways, reactions, substrates
  - Enzymes, transporters
  - Genes, replicons
  - Transcription factors/sites, promoters, operons

- **Tier 1: Literature-Derived PGDBs**
  - MetaCyc
  - EcoCyc -- *Escherichia coli* K-12

- **Tier 2: Computationally-derived DBs, Some Curation -- 24 PGDBs**
  - HumanCyc
  - Mycobacterium tuberculosis

- **Tier 3: Computationally-derived DBs, No Curation -- 481 DBs**

# Family of Pathway/Genome Databases

# *Pathway Tools Software: PathoLogic*

- **Computational creation of new Pathway/Genome Databases**

- **Transforms genome into Pathway Tools schema and layers inferred information above the genome**

- **Predicts operons**
- **Predicts metabolic network**
- **Predicts pathway hole fillers**
- **Infers transport reactions**

# *Pathway Tools Software: Pathway/Genome Editors*

- **Interactively update PGDBs with graphical editors**

- **Support geographically distributed teams of curators with object database system**

- **Gene editor**
- **Protein editor**
- **Reaction editor**
- **Compound editor**
- **Pathway editor**
- **Operon editor**
- **Publication editor**

# *Pathway Tools Software: Pathway/Genome Navigat...*

- **Querying, visualization of pathways, chromosomes, operons**

- **Analysis operations**
  - Pathway visualization of gene-expression data
  - Global comparisons of metabolic networks
  - Comparative genomics

- **WWW publishing of PGDBs**
- **Desktop operation**

# MetaCyc Data  -- Version 13.6

| | |
|---|---|
| Pathways | 1,436 |
| Reactions | 8,200 |
| Enzymes | 6,060 |
| Small Molecules | 8,400 |
| Organisms | 1,800 |
| Citations | 21,700 |

# Taxonomic Distribution of MetaCyc Pathways – version 13.1

| Bacteria | 883 |
|---|---|
| Green Plants | 607 |
| Fungi | 199 |
| Mammals | 159 |
| Archaea | 112 |

# MetaCyc Enzyme Data

- Reaction(s) catalyzed
- Alternative substrates
- Cofactors / prosthetic groups
- Activators and inhibitors
- Subunit structure
- Molecular weight, pI
- Comment, literature citations
- Species

# HumanCyc -- HumanCyc.org

- **Derived from Ensembl and LocusLink**

- **Tier 2 PGDB**
- **Curation has just resumed**

- **235 metabolic pathways**
- **1,523 small-molecule reactions**
- **1,188 substrates**

- *Genome Biology* 6:1-17 2004.

# *EcoCyc Project – EcoCyc.org*

- ***E. coli* Encyclopedia**
  - Review-level Model-Organism Database for *E. coli*
  - Tracks evolving annotation of the *E. coli* genome and cellular networks
  - The two paradigms of EcoCyc

- **Collaborative development via Internet**
  - Paulsen (TIGR) – Transport, flagella, DNA repair
  - Collado (UNAM) -- Regulation of gene expression
  - Keseler, Shearer (SRI) -- Metabolic pathways, cell division, proteases
  - Karp (SRI) -- Bioinformatics

# Paradigm 1: EcoCyc as Textual Review Article

- **All gene products for which experimental literature exists are curated with a minireview summary**
  - Found on protein and RNA pages, not gene pages!
  - 3257 gene products contain summaries
- **Summaries cover function, interactions, mutant phenotypes, crystal structures, regulation, and more**

- **Additional summaries found in pages for operons, pathways**

- **EcoCyc cites 14,269 publications**

# *Paradigm 2:*
# *EcoCyc as Computational Symbolic Theory*

- **Highly structured, high-fidelity knowledge representation provides computable information**
- **Each molecular species defined as a DB object**
  - Genes, proteins, small molecules
- **Each molecular interaction defined as a DB object**
  - Metabolic reactions
  - Transport reactions
  - Transcriptional regulation of gene expression
- **220 database fields capture extensive properties and relationships**

# *Demonstration*

# Pathway Tools Schema and Semantic Inference Layer

# *Guiding Principles for the Pathway Tools Ontology of Biological Function*

- **Encode distinct molecular species as separate objects**
- **Describe all molecular interactions as reactions**

- **Layered approach:**
  - Molecular species form the base
  - Reactions built from molecular species
  - Pathways built from reactions

- **Link catalyst to reaction via Enzymatic-Reaction**

| Reaction | ↔ | Enzymatic Reaction | ↔ | Enzyme |
|---|---|---|---|---|

# *Pathway Tools Ontology / Schema*

- **Ontology classes: 1621**
  - Datatype classes: Define objects from genomes to pathways
  - Classification systems / controlled vocabularies
    - Pathways, chemical compounds, enzymatic reactions (EC system)
    - Protein Feature ontology
    - Cell Component Ontology
    - Evidence Ontology

- **Comprehensive set of 279 attributes and relationships**

# *Overview of Schema Presentation*

- **Survey of important classes**

- **What slots are present within these classes**

- **How objects are linked together to form a network**

# Use GKB Editor to Inspect the Pathway Tools Ontology

- GKB Editor = Generic Knowledge Base Editor
- Type in Navigator window:   (GKB)        or
- [Right-Click]  Edit->Ontology Editor

- View->Browse Class Hierarchy
- [Middle-Click] to expand hierarchy
- To view classes or instances, select them and:
    - Frame -> List Frame Contents
    - Frame -> Edit Frame

# Root Classes in the Pathway Tools Ontology

- **Chemicals**    -- **All molecules**
- **Polymer-Segments**    -- **Regions of polymers**
- **Protein-Features**    -- **Features on proteins**
- **Paralogous-Gene-Groups**

- **Organisms**

- **Generalized-Reactions**    -- **Reactions and pathways**
- **Enzymatic-Reactions**    -- **Link enzymes to reactions they catalyze**
- **Regulation**    -- **Regulatory interactions**

- **CCO**    -- **Cell Component Ontology**
- **Evidence**    -- **Evidence ontology**

- **Notes**    -- **Timestamped, person-stamped notes**
- **Organizations**
- **People**
- **Publications**

# *Principal Classes*

- **Class names are usually capitalized, plural, separated by dashes**

- **Genetic-Elements, with subclasses:**
  - Chromosomes
  - Plasmids
- **Genes**
- **Transcription-Units**
- **RNAs**
  - rRNAs, snRNAs, tRNAs, Charged-tRNAs
- **Proteins, with subclasses:**
  - Polypeptides
  - Protein-Complexes

# *Principal Classes*

- **Reactions**

- **Enzymatic-Reactions**

- **Pathways**

- **Compounds-And-Elements**

- **Regulation**

# *Semantic Network Diagrams*

TCA Cycle

↑ in-pathway

Succinate + FAD = fumarate + FADH$_2$

↑ reaction

Enzymatic-reaction

↑ catalyzes

Succinate dehydrogenase

component-of

Sdh-flavo  Sdh-Fe-S  Sdh-membrane-1  Sdh-membrane-2

product

sdhA  sdhB  sdhC  sdhD

# Pathway Tools Schema and Semantic Inference Layer

# Genes, Operons, and Replicons

# *Representing a Genome*



- **Classes:**
  - ORG is of class Organisms
  - CHROM1 is of class Chromosomes
  - PLASMID1 is of class Plasmids
  - Gene1 is of class Genes
  - Product1 is of class Polypeptides or RNA

# *Polynucleotides*



Review slots of COLI and of COLI-K12

# *Polymer-Segments*



Review slots of Genes

# *Proteins*

# *Proteins and Protein Complexes*

- **Polypeptide: the monomer protein product of a gene (may have multiple isoforms, as indicated at gene level)**

- **Protein complex: proteins consisting of multiple polypeptides or protein complexes**

- **Example: DNA pol III**
  - DnaE is a polypeptide
  - pol III core enzyme contains DnaE, DnaQ, HolE
  - pol III holoenzyme contains pol III core enzyme plus three other complexes

# Slots of Proteins (DnaE)

- comments, citations
- pI, molecular-weight
- features
- component-of
- gene
- catalyzes   [link to Enzymatic-Reaction]
- dblinks

# Semantic Network Diagrams

TCA Cycle

↑ in-pathway

Succinate + FAD = fumarate + FADH$_2$

↑ reaction

Enzymatic-reaction

↑ catalyzes

Succinate dehydrogenase

component-of

Sdh-flavo    Sdh-Fe-S    Sdh-membrane-1    Sdh-membrane-2

↑    ↑    ↑ product    ↑

sdhA    sdhB    sdhC    sdhD

# *Semantic Inference Layer*

- **Reactions-of-protein (prot)**
  - Returns a list of rxns this protein catalyzes
- **Transcription-units-of-proteins(prot)**
  - Returns a list of TU's activated/inhibited by the given protein
- **Transporter? (prot)**
  - Is this protein a transporter?
- **Polypeptide-or-homomultimer?(prot)**
- **Transcription-factor? (prot)**
- **Obtain-protein-stats**
  - Returns 5 values
    - Length of : all-polypeptides, complexes, transporters, enzymes, etc…

# Compounds / Reactions / Pathways

# *Compounds / Reactions / Pathways*

- **Think of a three tiered structure:**
  - Compounds at the bottom
  - Reactions built on top of compounds
  - Pathways built on top of reactions
- **Metabolic network can be defined by reactions alone**
- **Pathways are an additional "optional" structure**
- **Some reactions not part of a pathway**
- **Some reactions have no attached enzyme**
- **Some enzymes have no attached gene**

# *Compounds*

# *Slots of Compounds*

- **common-name, abbrev-name, synonyms**
- **comment, citations**
- **charge, gibbs-0, molecular-weight**
- **empirical-formula**
- **structure-atoms, structure-bonds**
- **appears-in-left-side-of, appears-in-right-side-of**

# *Semantic Inference Layer*

- **Reactions-of-compound (cpd)**
- **Pathways-of-compound (cpd)**
- **Activated/inhibited-by? (cpds slots)**
  - Returns a list of enzrxns for which a cpd in cpds is a modulator (example slots: activators-all, activators-allosteric)
- **All-substrates (rxns)**
  - All unique substrates specified in the given rxns
- **Has-structure-p (cpd)**

# *Reactions*

# *Reactions*

- **Represent information about a reaction that is independent of enzymes that catalyze the reaction**

- **Connected to enzyme(s) via enzymatic reaction frames**

- **Classified with EC system when possible**

- **Example: 2.7.7.7 – DNA-directed DNA polymerization**
  - Carried out by five enzymes in *E. coli*

# *Slots of Reaction Frames*

- **Keq**
- **Left and Right  (reactants / products)**
  - Can include modified forms of proteins, RNAs, etc here
- **Enzymatic-reaction**
- **In-pathway**

# *Semantic Inference Layer*

- **Genes-of-reaction (rxn)**
- **Substrates-of-reaction (rxn)**
- **Enzymes-of-reaction (rxn)**
- **Lacking-ec-number (organism)**
  - Returns list of rxns with no ec numbers in that database
- **Get-reaction-direction-in-pathway (pwy rxn)**
- **Reaction-type(rxn)**
  - Indicates types of Rxn as: Small molecule rxn, transport rxn, protein-small-molecule rxn (one substrate is protein and one is a small molecule), protein rxn (all substrates are proteins), etc.
- **All-rxns(type)**
  - Specify the type of reaction (see above for type)
- **Obtain-rxn-stats**
  - Returns six values
    - Length of : all-rxns, transport, non-transport, etc…

# *Enzymatic Reactions (DnaE and 2.7.7.7)*

- **A necessary bridge between enzymes and "generic" versions of reactions**
- **Carry information specific to an enzyme/reaction combination:**
  - Cofactors and prosthetic groups
  - Alternative substrates
  - Links to regulatory interactions

- **Frame is generated when protein is associated with reaction (via protein or reaction editor)**

# *Regulation of Enzyme Activity*

**Frame Editor**

Application  Knowledge Base  Frame  Value  Slot  Preferences

Editing Instance REG0-6616, instance of Regulation-of-Enzyme-Activity

CITATIONS —— "11515538"

COMMENT

COMMENT-INTERNAL

COMMON-NAME

CREATION-DATE —— 16-Jun-2007 01:08:08

CREATOR —— paley

CREDITS

DATA-SOURCE

HISTORY

MECHANISM —— NONCOMPETITIVE

MODE —— "-"

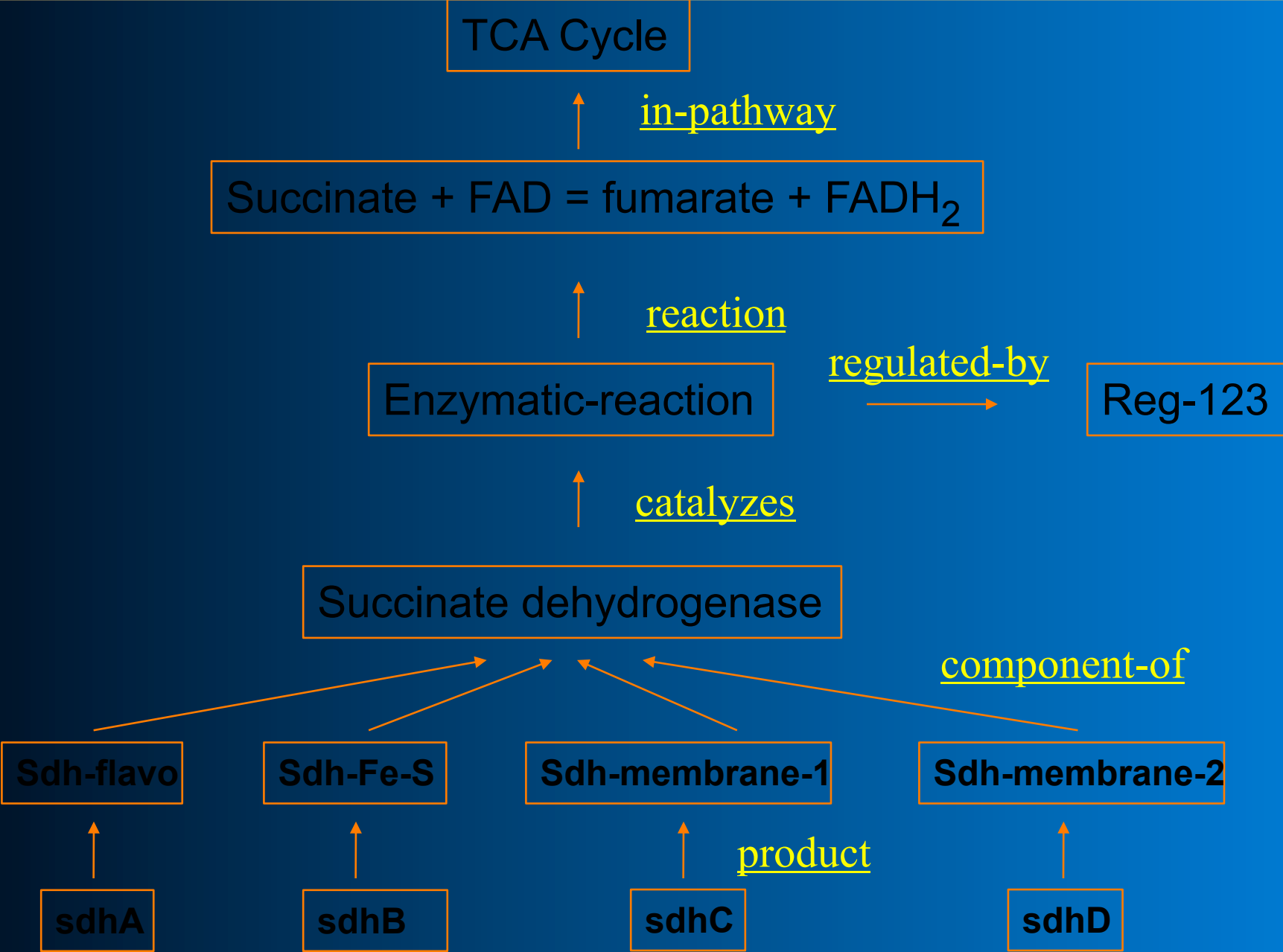PHYSIOLOGICALLY-RELEVANT?

REGULATED-BY

REGULATED-ENTITY —— aminopeptidase

REGULATOR —— EDTA

SCHEMA? ——  T
             inherited from Regulation

SYNONYMS

# Semantic Network Diagrams

TCA Cycle

↑ in-pathway

Succinate + FAD = fumarate + FADH$_2$

↑ reaction

Enzymatic-reaction → regulated-by → Reg-123

↑ catalyzes

Succinate dehydrogenase

component-of

Sdh-flavo     Sdh-Fe-S     Sdh-membrane-1     Sdh-membrane-2

product

sdhA     sdhB     sdhC     sdhD

# *Pathway Tools Schema and Semantic Inference Layer: Pathways*

# *Pathway Ontology*

- **Slots in pathway:**
  - Reaction-List,  Predecessor-List

A $\xrightarrow{R1}$ B $\xrightarrow{R2}$ C $\xrightarrow{R3}$ D

A $\xrightarrow{R1}$ B $\nearrow^{R2}$ C
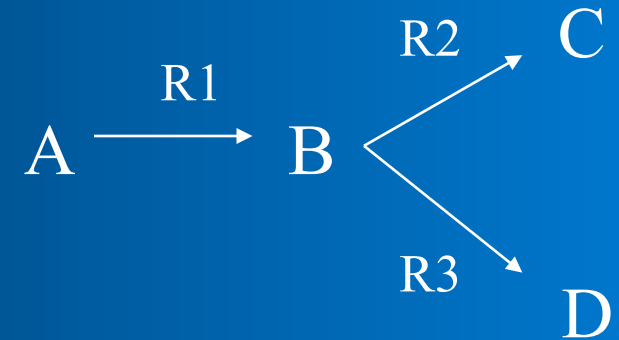
B $\searrow_{R3}$ D

R1:    Left = A, Right = B
R2:    Left = B, Right = C
R3:    Left = C, Right = D

Predecessor list:
   (R1 R2) (R2 R3)

R1: Left = A, Right = B
R2: Left = B, Right = C
R3: Left = B, Right = D

Predecessor list:
   (R1 R2) (R1 R3)

# *Super-Pathways*

- **Collection of pathways that connect to each other via common substrates or reactions, or as part of some larger logical unit**
- **Can contain both sub-pathways and additional connecting reactions**
- **Can be nested arbitrarily**
- **REACTION-LIST: a pathway ID instead of a reaction ID in this slot means include all reactions from the specified pathway**
- **PREDECESSORS: a pathway ID instead of a tuple in this slot means include all predecessor tuples from the specified pathway**

# *Querying Pathways Programmatically*

- **See http://bioinformatics.ai.sri.com/ptools/ptools-resources.html**
- **(all-pathways)**
- **(base-pathways)**
  - Returns list of all pathways that are not super-pathways
- **(genes-of-pathway pwy)**
- **(unique-genes-of-pathway pwy)**
  - Returns list of all genes of a pathway that are not also part of other pathways
- **(enzymes-of-pathway pwy)**
- **(substrates-of-pathway pwy)**
- **(variants-of-pathway pwy)**
  - Returns all pathways in the same variant class as a pathway
- **(get-predecessors rxn pwy), (get-successors rxn pwy)**
- **(get-rxn-direction-in-pathway pwy rxn)**
- **(pathway-inputs pwy), (pathway-outputs pwy)**
  - Returns all compounds consumed (produced) but not produced (consumed) by pathway (ignores stoichiometry)

# *Regulation*

# Encoding Cellular Regulation in Pathway Tools -- Goals

- **Facilitate curation of wide range of regulatory information within a formal ontology**
- **Compute with regulatory mechanisms and pathways**
  - Summary statistics, complex queries
  - Pattern discovery
  - Visualization of network components

- **Provide training sets for inference of regulatory networks**
- **Interpret gene-expression datasets in the context of known regulatory mechanisms**

# *Regulation in Pathway Tools*

- **Substrate-level regulation of enzyme activity**
- **Binding to proteins or small molecules (phosphorylation)**
- **Regulation of transcription initiation**
- **Attenuation of transcription**
- **Regulation of translation by proteins and by small RNAs**

# *Regulation*

- **Class Regulation with subclasses that describe different biochemical mechanisms of regulation**
- **Slots:**
  - Regulator
  - Regulated-Entity
  - Mode
  - Mechanism

# *Regulation of Enzyme Activity*

- **Class Regulation-of-Enzyme-Activity**
- **Each instance of the class describes one regulatory interaction**

- **Slots:**
  - Regulator  --  usually a small molecule
  - Regulated-Entity  --  an Enzymatic-Reaction
  - Mechanism  -- One of:
    - Competitive, Uncompetitive, Noncompetitive, Irreversible, Allosteric, Unkmech, Other
  - Mode  --  One of:  + , -

# *Transcription Initiation*

- **Class Regulation-of-Transcription-Initiation**

- **Slots:**
  - Regulator  --  instance of Proteins or Complexes (a transcription-factor)
  - Regulated-Entity  --  instance of Promoters or Transcription-Units or Genes
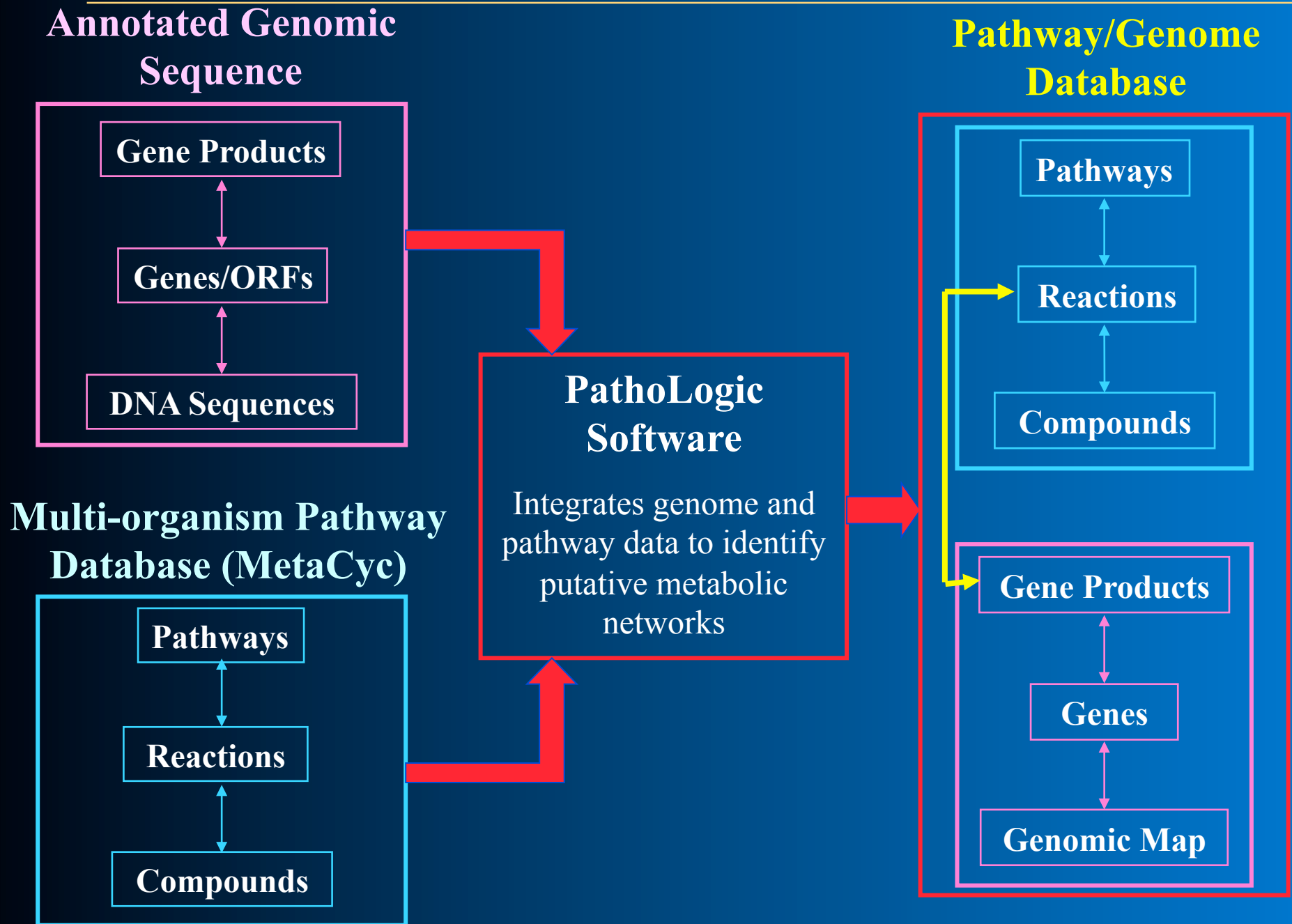  - Mode  --  One of:  +  , -

# *Other Features of Ontology*

- **Evidence codes**

- **Curator crediting system**

# *Inference Algorithms*

# PathoLogic: Inference of Pathway Complement

- An additional level of inference after genome annotation

- Place predicted genes in their biochemical context
- Information reduction device
- Assess coherence of the set of genes in a genome
- Identify pathway holes and singleton enzymes
- Provides a framework for analysis of functional-genomics data

```
Genbank Format:

    gene                   422054..423490
                           /gene="aroE"
    CDS                    422054..423490
                           /gene="aroE"
                           /label="CT370"
                           /product="Shikimate 5-Dehydrogenase"
                           /db_xref="PID:g3328794"

PathoLogic Format:

        ID             CT370
        NAME           aroE
        STARTBASE      422054
        ENDBASE        423490
        PRODUCT        Shikimate 5-Dehydrogenase
        DBLINK         PID:g3328794
        PRODUCT-TYPE   P
        EC             1.1.1.25
```

# *Pathway Prediction*

- **Step 1: Infer reactome**

- **Step 2: Infer metabolic pathways from reactome**
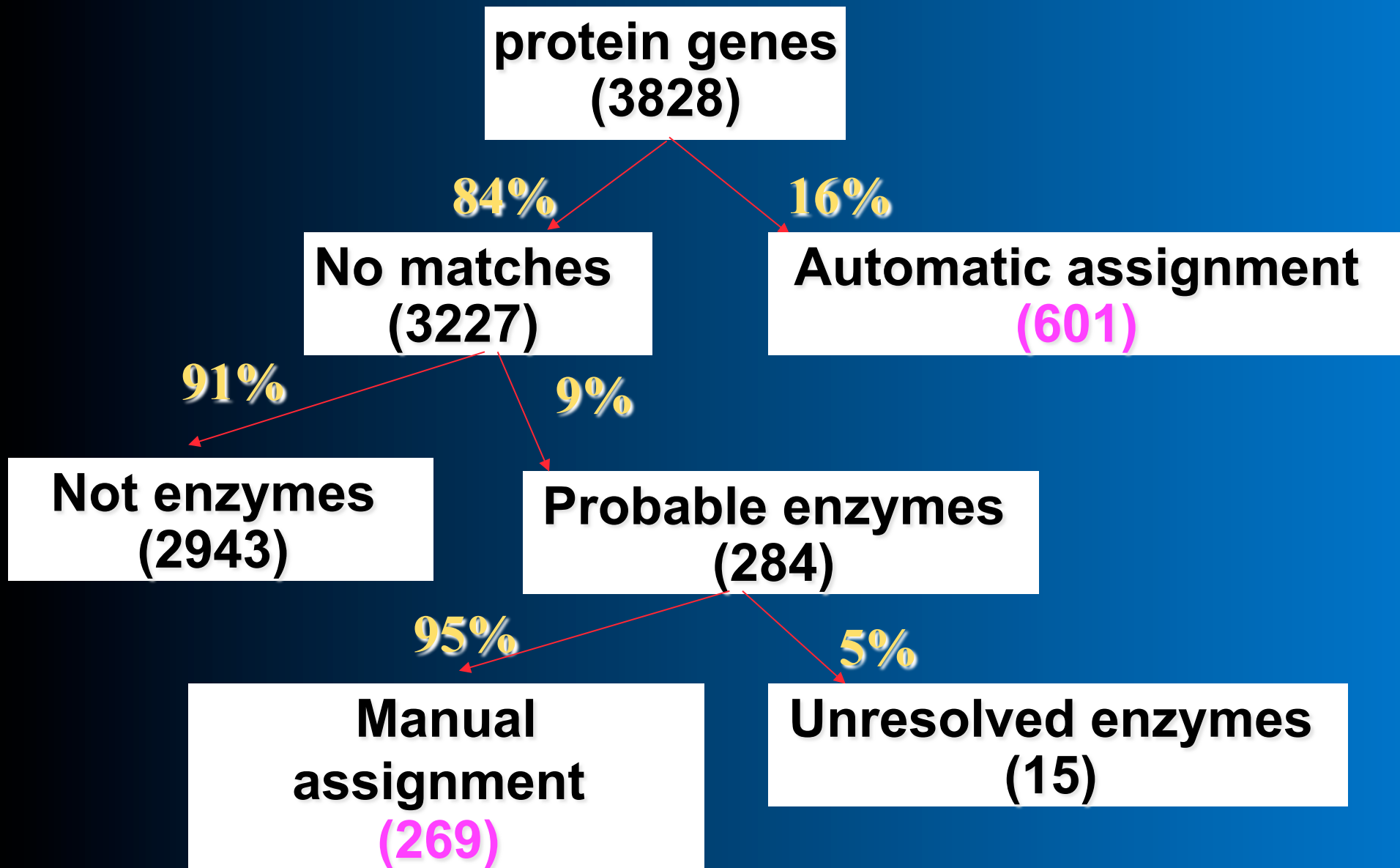
# *Inference of Reactome*

- **Given genome annotation, infer metabolic reactions that can be catalyzed by the genome**
  - EC numbers
  - Enzyme names
  - Gene Ontology annotations

- **Complications:**
  - Most genomes contain a subset of above annotations
  - Enzyme names sometimes ambiguous
  - Some reactions occur in multiple pathways
    - 99 of 744 reactions in *E. coli*
  - Pathway variants

# *Match Enzymes to Reactions*

Gene product

5.1.3.2

UDP-glucose-4-epimerase

MetaCyc

Match

no — Probable enzyme-ase

yes — Assign

UDP-D-glucose → UDP-galactose

no — Not a metabolic enzyme

yes — Manually search

no — Can't Assign

yes — Assign

# *Pathway Prediction Algorithm*

- **Two pathway lists:**
  - U: Undecided status
  - K: Keep

- **Initialize U to contain all MetaCyc pathways for which at least one reaction has an enzyme**

# *Pathway Prediction Algorithm*

- **For each P in U:**
  - If current organism is outside taxonomic range of P AND at least one reaction in P lacks an enzyme, delete P from U
  - If all reactions of P designated as key reactions have no enzyme, delete P from U

# *Pathway Prediction Algorithm*

- **Iterate through P in U until U is unchanged:**
  - If P should be kept, move P to K
    - A reaction in P is unique to P and has an enzyme
    - At most one reaction in P has no enzyme
    - The enzymes present for P are not a subset of the enzymes present for a variant pathway of P
  - If P should be deleted, delete P from U
    - At most one reaction R in P has an enzyme, and R is not unique to P
    - The pathway is a biosynthetic pathway missing its final steps
    - The pathway is a catabolic pathway missing its initial steps

- **Accuracy: 91%**

# *Pathway Evidence Report*

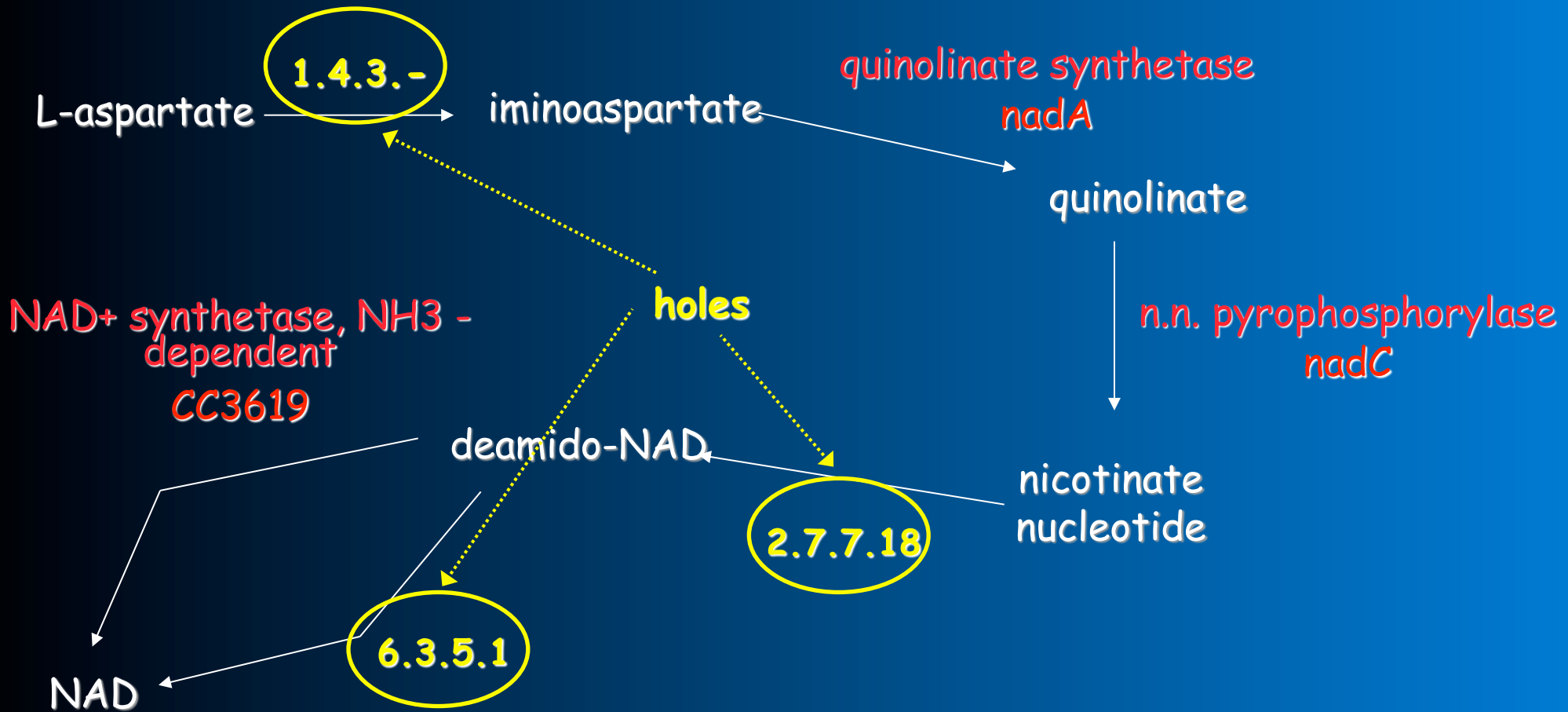## Biosynthesis:Cofactors, Prosthetic Groups, Electron Carriers

| Pathway | Pathway Glyph | Total Rxns | Rxns Present in V. cholerae | Rxns Present in Other Pwys | Other Pwys |
|---|---|---|---|---|---|
| biosynthesis of proto- and siroheme | | 15 | 12 | 2 | tRNA charging pathway cobalamin biosynthesis II (aerobic pathway) |
| biotin biosynthesis I | | 4 | 4 | 0 | (none) |
| cobalamin biosynthesis I | | 7 | 6 | 2 | cobalamin biosynthesis II (aerobic pathway) |
| cobalamin biosynthesis II (aerobic pathway) | | 18 | 5 | 3 | cobalamin biosynthesis I biosynthesis of proto- and siroheme |

# *Limitations of Pathway Inference*

- **Can be misled by missing or incorrect functional assignments**

- **No sequences known for many enzymes**

- **Uncertainty for short pathways**

# *Pathway Hole Filler*

- **Why should hole filler find things beyond the original genome annotation?**

- **Reverse BLAST searches more sensitive**
- **Reverse BLAST searches find second domains**
- **Integration of multiple evidence types**

- **In 90 minutes, I only got to here**
- **Included a 10-15 min demo**
- **3/8/2007 Brutlag class lecture**

# PathoLogic Step 6:
# Build Cellular Overview Diagram

- **Diagram encompassing metabolic, transport, and other cellular networks**

- **Automatically generated for every BioCyc DB using advanced graph layout algorithm**

- **Harness the power of the human visual system to interpret patterns in a mechanistic context**

- **Can be zoomed, interrogated, and painted with experimental or comparative data**

File  Overviews  Pathway  Reaction  Protein  RNA  Gene  Compound  Chromosome  Tools  Help

Escherichia coli  Home  Back  Forward  History  Next Answer  Clone  Save DB

E. coli Chromosome Global Overview

Escherichia coli K-12 Chromosome:

| | |
|---|---|
| 190 | 146,694 |
| 146,968 | 279,099 |
| 278,402 | 416,176 |
| 416,366 | 562,542 |
| 562,553 | 695,727 |
| 695,765 | 837,148 |
| 837,413 | 997,630 |
| 997,713 | 1,130,657 |
| 1,130,661 | 1,258,292 |
| 1,258,347 | 1,395,116 |
| 1,395,389 | 1,525,176 |
| 1,525,914 | 1,661,631 |
| 1,661,633 | 1,805,323 |
| 1,805,424 | 1,940,607 |
| 1,940,686 | 2,062,491 |
| 2,062,503 | 2,203,706 |
| 2,203,717 | 2,365,089 |
| 2,365,093 | 2,506,264 |
| 2,506,483 | 2,651,560 |
| 2,651,537 | 2,786,260 |
| 2,786,399 | 2,923,302 |
| 2,923,370 | 3,071,713 |
| 3,071,998 | 3,207,509 |
| 3,207,552 | 3,344,172 |
| 3,344,195 | 3,469,351 |
| 3,469,422 | 3,614,208 |
| 3,614,205 | 3,780,585 |
| 3,780,665 | 3,920,074 |
| 3,920,083 | 4,062,318 |
| 4,062,386 | 4,212,146 |
| 4,212,303 | 4,374,465 |
| 4,374,576 | 4,505,486 |
| 4,505,489 | 4,639,651 |

Command: :Cmd Menu Node Network "chbR"
Command: :Cmd Menu Node Network "prpR"
Command: Genbro Show Chromosome Global View
Command: []

File   Overviews   Pathway   Reaction   Protein   RNA   Gene   Compound   Chromosome   Tools   Help

Escherichia coli   | Home | Back | Forward | History | Next Answer | Clone |                    Save DB



Command: :Cmd Menu Node Network "uhpA"
Command: :Cmd Menu Node Network "chbR"
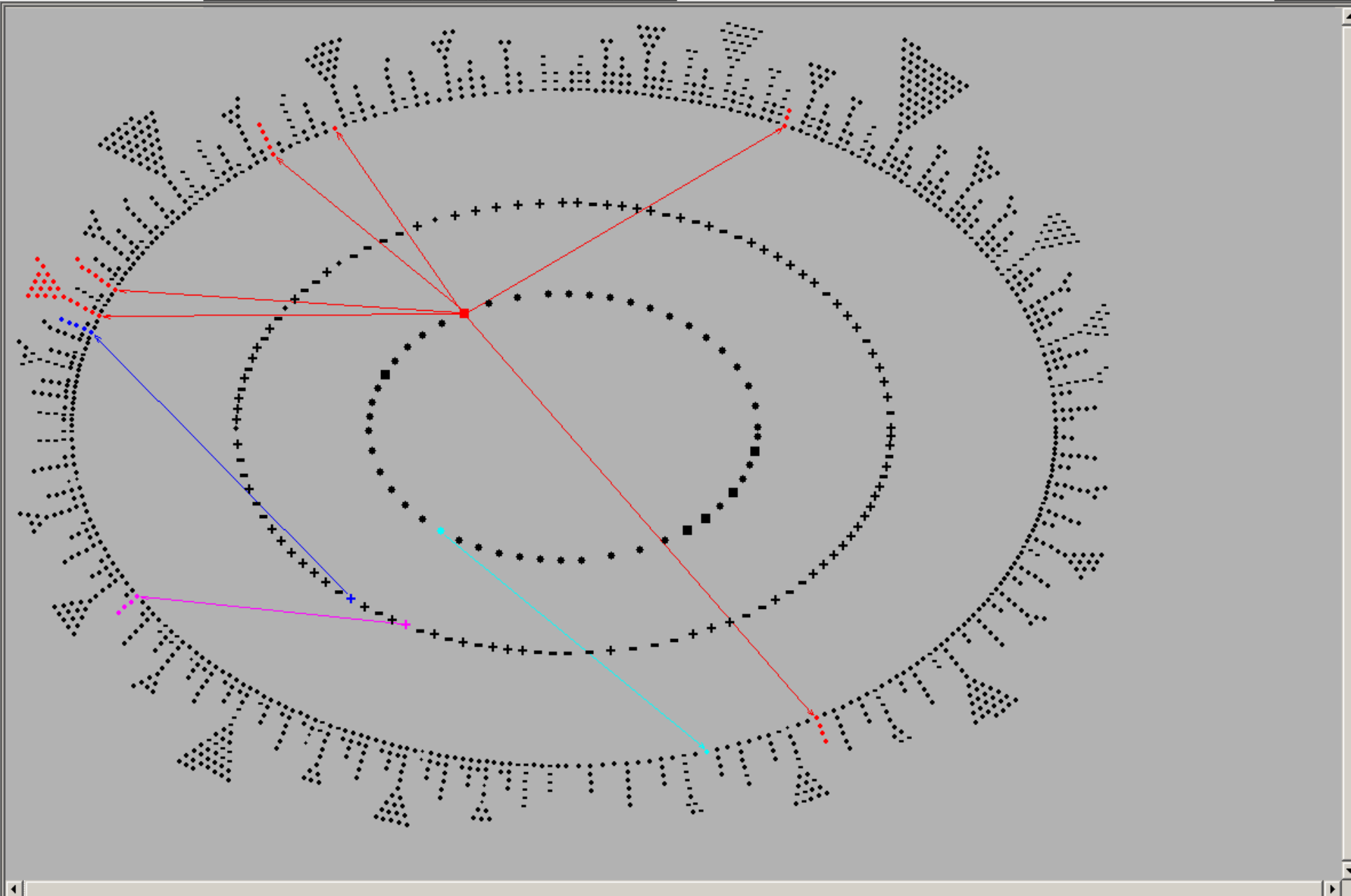Command: :Cmd Menu Node Network "prpR"
Command:

# *Pathway Algorithms*

- **Automated layout of metabolic pathways**
  - *Bioinformatics and Genome Research Conference* 1994 p225
- **Automated layout of cellular overview diagram**
- **Automated generation of metabolic map poster**

- **Forward propagation of metabolites through the metabolic network**
  - Consistency of a PGDB with respect to known growth-media requirements
  - *Pacific Symp Biocomputing* 2001:471
- **Identify dead-end metabolites**

- **Infer drug targets as choke points in metabolic network**
  - *Genome Research* 14:917 2004

# Dead End Metabolites

- **Clues to extra/missing reactions**
- **A small molecule C is a dead-end if:**
  - (Def 1 easier to compute; Def 2 more accurate)
- **Definition 1:**
  - C is a substrate in only one reaction of the set of SMM reactions occurring in Compartment  AND
  - No transporter acts on C in Compartment, nor on parent classes of C
- **Definition 2:**
  - C is produced only by SMM reactions in Compartment, and no transporter acts on C in Compartment  OR
  - C is consumed only by SMM reactions in Compartment, and no transporter acts on C in Compartment

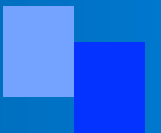# Global Consistency Checking of Biochemical Network

- **Given:**
  - A PGDB for an organism
  - A set of initial metabolites

- **Infer:**
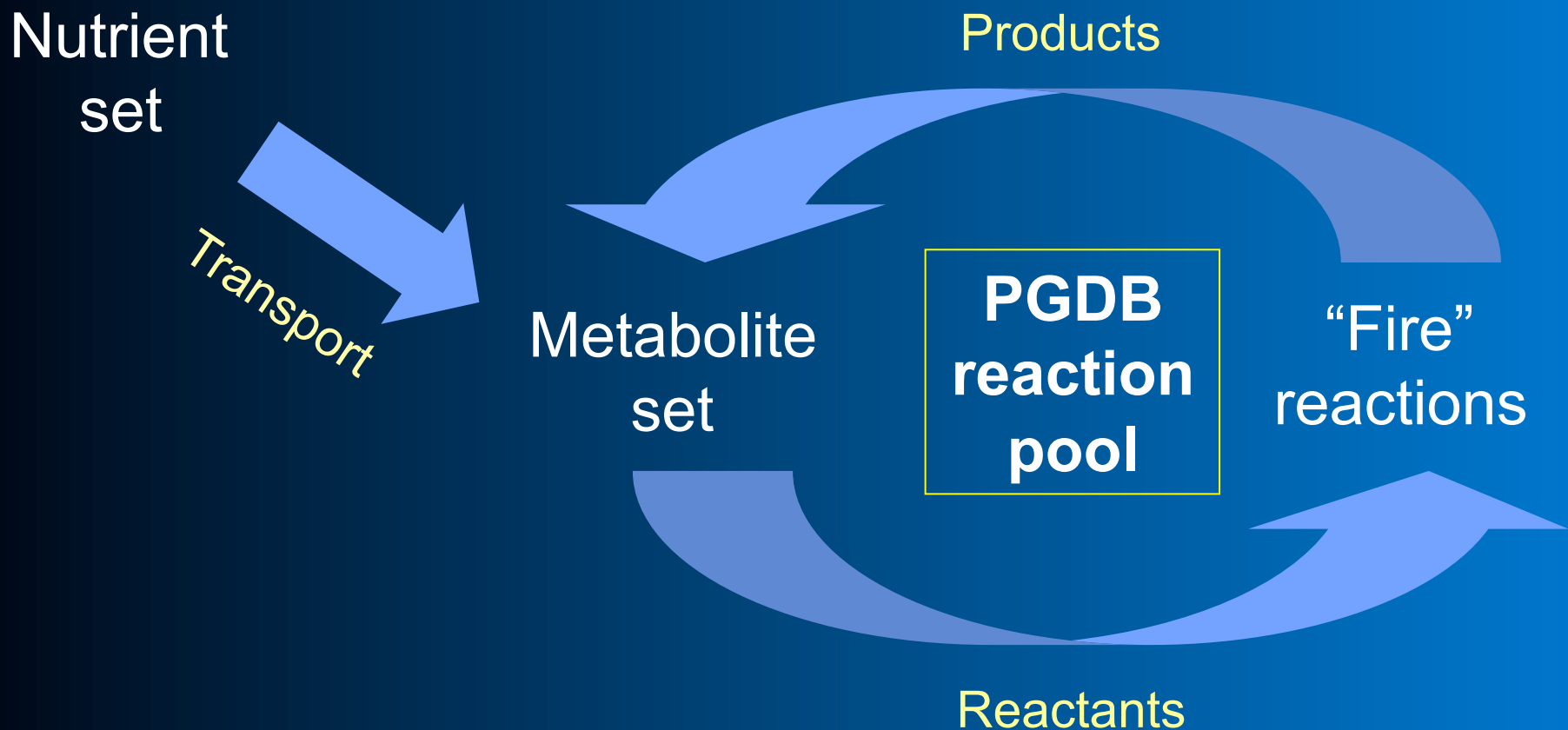  - What set of products can be synthesized by the small-molecule metabolism of the organism

- **Can known growth medium yield known essential compounds?**

# Algorithm: Forward Propagation Through Production System

- Each reaction becomes a production rule
- Each metabolite in nutrient set becomes an axiom

Nutrient set

Transport

Metabolite set

Products

PGDB reaction pool

"Fire" reactions

Reactants

# Initial Metabolite Nutrient Set (Total: 21 compounds)

| | |
|---|---|
| Nutrients (8) (M61 Minimal growth medium) | $H^+$, $Fe^{2+}$, $Mg^{2+}$, $K^+$, $NH_3$, $SO_4^{2-}$, $PO_4^{2-}$, Glucose |
| Nutrients (10) (Environment) | Water, Oxygen, Trace elements ($Mn^{2+}$, $Co^{2+}$, $Mo^{2+}$, $Ca^{2+}$, $Zn^{2+}$, $Cd^{2+}$, $Ni^{2+}$, $Cu^{2+}$) |
| Bootstrap Compounds (3) | ATP, NADP, CoA |

# Essential Compounds
## E. coli Total: 41 compounds

- **Proteins (20)**
  - Amino acids
- **Nucleic acids (DNA & RNA) (8)**
  - Nucleosides
- **Cell membrane (3)**
  - Phospholipids
- **Cell wall (10)**
  - Peptidoglycan precursors
  - Outer cell wall precursors (Lipid-A, oligosaccharides)

Nutrients: A, B, C, E, F

A + B → W

C + D → X

E + F → Y

W + Y → Z

Produced Compounds: W, Y, Z

# *Results*

- **Phase I: Forward propagation**
  - 21 initial compounds yielded only half of the 41 essential compounds for *E. coli*

- **Phase II: Manually identify**
  - Bugs in EcoCyc (e.g., two objects for tryptophan)
    - A $\rightarrow$ B      B' $\rightarrow$ C
  - Incomplete knowledge of *E. coli* metabolic network
    - A + **B** $\rightarrow$ C + D
  - "Bootstrap compounds"
  - Missing initial protein substrates (e.g., ACP)
    - Protein synthesis not represented

- **Phase III: Forward propagation with 11 more initial metabolites**
  - Yielded all 41 essential compounds

# *Summary*

- **Pathway/Genome Databases**
  - MetaCyc non-redundant DB of literature-derived pathways
  - Additional organism-specific PGDBs available through SRI at BioCyc.org
  - Computational theories of biochemical machinery

- **Pathway Tools software**
  - Extract pathways from genomes
  - Morph annotated genome into structured ontology
  - Distributed curation tools for MODs
  - Query, visualization, WWW publishing

# *How to Learn More*

- **BioCyc Webinars**
  - See BioCyc.org

- **BioCyc publications page**
  - BioCyc.org

- **Pathway Tools training course**
- **Pathway Tools feedback sessions**
  - ptools-support@ai.sri.com

- **Try out Pathway Tools**

# Additional Pathway Tools Algorithms

- **Predict metabolic pathway complement**
- **Automatic layout of Cellular Overview diagram**
- **Paint Omics datasets onto Cellular Overview**
- **Compare metabolic networks**
- **Reaction balance checker**
- **Chemical substructure search**
- **Predict operons**
- **Predict pathway hole fillers**
- **Qualitative path tracing from network inputs to network outputs**