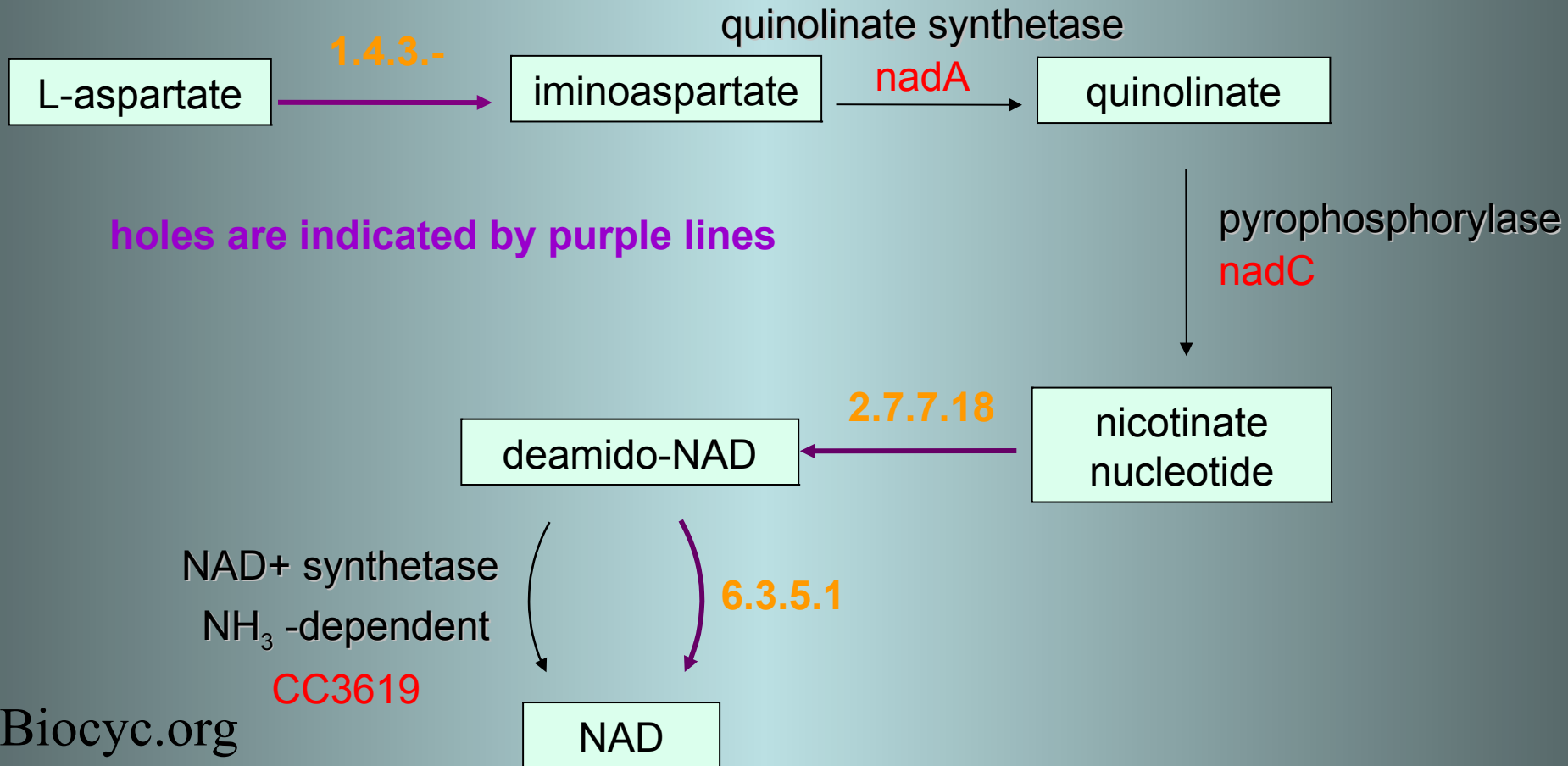




# Identify Pathway Hole Fillers

Definition: Pathway Holes are reactions in metabolic pathways for which no enzyme is identified in the PGDB.



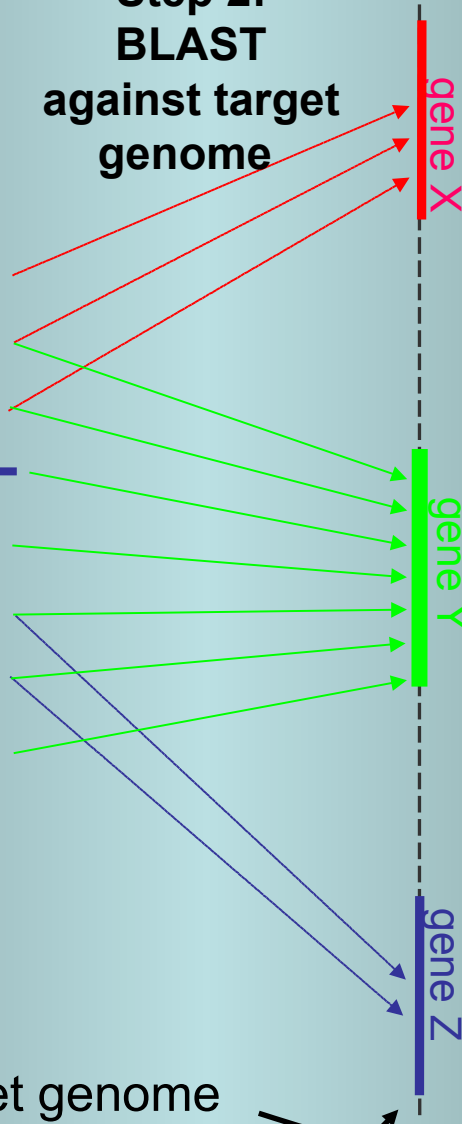
# Algorithm for identifying candidates and consolidating data



Step 1: collect query  
isozymes of function  
A based on EC#

organism 1 enzyme A   
organism 2 enzyme A   
organism 3 enzyme A   
organism 4 enzyme A   
organism 5 enzyme A   
organism 6 enzyme A   
organism 7 enzyme A   
organism 8 enzyme A 

Step 2:  
BLAST  
against target  
genome



Step 3 & 4: Consolidate  
hits and evaluate  
evidence

Candidates

**Gene X**

**Gene Y**

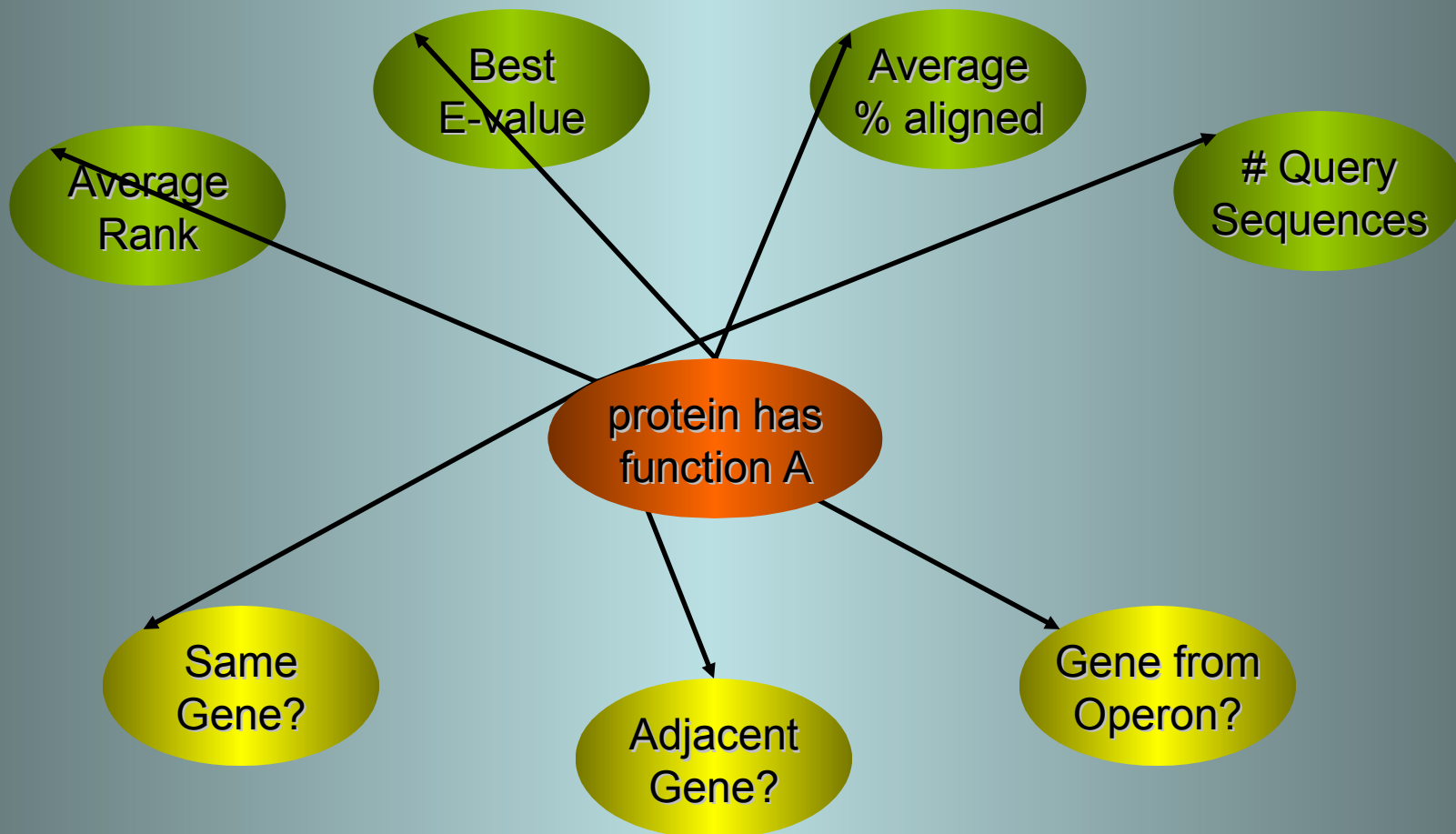
**Gene Z**



# Features used to calculate the probability that a protein has the desired function

- Best E-value
- Avg. rank of candidate sequence in BLAST output
- Avg. length % aligned
- Number of query sequences aligned
- Candidate in same direction as another pathway gene?
- Candidate is adjacent to a gene that catalyzes an adjacent reaction?
- Candidate catalyzes another pathway reaction?

# Use Bayesian classifier to evaluate candidates

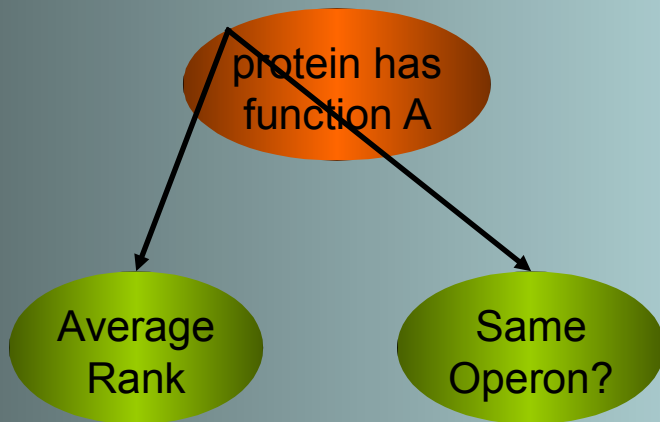




# Computing P(has function)

Apply Bayes' rule:

$$P(\text{true} | \text{evidence}) = \frac{P(\text{true}) P(\text{evidence} | \text{true})}{P(\text{true}) P(\text{evidence} | \text{true}) + P(\text{false}) P(\text{evidence} | \text{false})}$$



Compute probability distributions, i.e.,  $P(\text{evidence} | \text{true})$  and  $P(\text{evidence} | \text{false})$ , from the “known” reactions in the database.

e.g., Same operon?

In operon?	True hit Has-Fn(A)	False hit $\sim$ Has-Fn(A)
yes	0.24 (TP)	0.04 (FP)
no	0.76 (FN)	0.96 (TN)



# Computing P(has function)

Example:

Candidate X has avg-rank 1.5 and is in a directon with another pathway gene.

From training data:

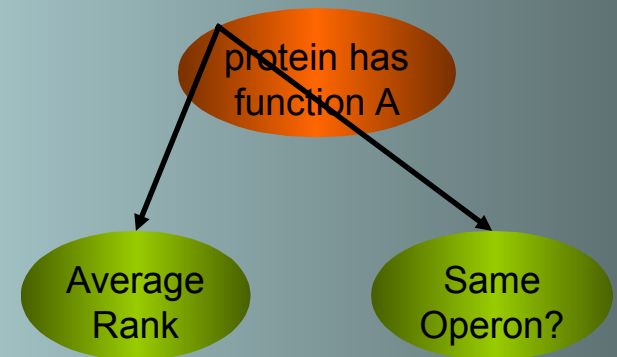
$$P(\text{average-rank} = 1.5 \mid \text{has-function}) = 0.40$$

$$P(\text{average-rank} = 1.5 \mid \neg \text{has-function}) = 0.03$$

$$P(\text{pathway-directon} = \text{true} \mid \text{has-function}) = 0.24$$

$$P(\text{pathway-directon} = \text{true} \mid \neg \text{has-function}) = 0.04$$

$P(\text{has-function}) = 0.041$  (4.1% of candidates in training data are true hits)



$$P(\text{has}_{\text{function}_A}) = \frac{0.041 * 0.40 * 0.24}{0.041 * 0.40 * 0.24 + 0.959 * 0.03 * 0.04}$$



# Steps that must be completed before running the Pathway Hole Filler

- Install BLAST executable (see Installation instructions)
- Prepare BLAST protein db for the PGDB
  - Need FASTA format genome nucleotide sequence. (If only ESTs are available, see User Guide regarding Prepare BLAST Reference Data->Protein from ESTs)
- In general, the more pathways in your PGDB, the more candidates the pathway hole filler will have to find



## ***Conceptual stages of the pathway hole filler***

### 1. Prepare training data for Bayes classifier

- Collect feature data for known rxns in PGDB
- Calculate probability distributions for classifier

### 2. Identify and evaluate candidates

- Collect feature data for each candidate
- Use classifier to determine  $P(\text{has-function})$

### 3. Choose holes to fill in KB

- Either select all above a cut-off or manually review candidates



# Navigating to the Pathway Hole Filler

The screenshot shows the PathoLogic software interface. The 'Refine' menu is open, and 'Pathway Hole Filler' is selected. The interface includes a sidebar with organism names, a main panel with organism details, and a right-hand panel with a file path and a list of numbers.

**Organism:** ID: MTBRV  
**Name:** M. tb.  
**Strain:** H37Rv  
**Status:** Built

**Genetic Elements:** H37Rv Chromosome

**Refine Menu:**

- Resolve Ambiguous Name Matches
- Assign Probable Enzymes
- Assign Modified Proteins
- Create Protein Complexes
- Re-Run Name Matcher
- Rescore Pathways
- Predict transcription units
- Run Consistency Checker
- Update Overview
- Pathway Hole Filler**

**Right Panel:**

/hapuna4/aic/ecocyc/mtbrvcyc/beta/kb/mtbrvbase.occ |  
ECOCYC  
7000 8000 9000 10000 11000 12000 13000 14000

**Pathway Hole Filler Options:**

- Fully-Automatic
- Wizard
- Expert Mode, Step 1: Prepare Training Data
- Expert Mode, Step 2: Identify and Evaluate Candidates
- Expert Mode, Step 3: Choose Holes to Fill in KB

**Footer:** L: Copy Region To Clipboard; R: Menu.



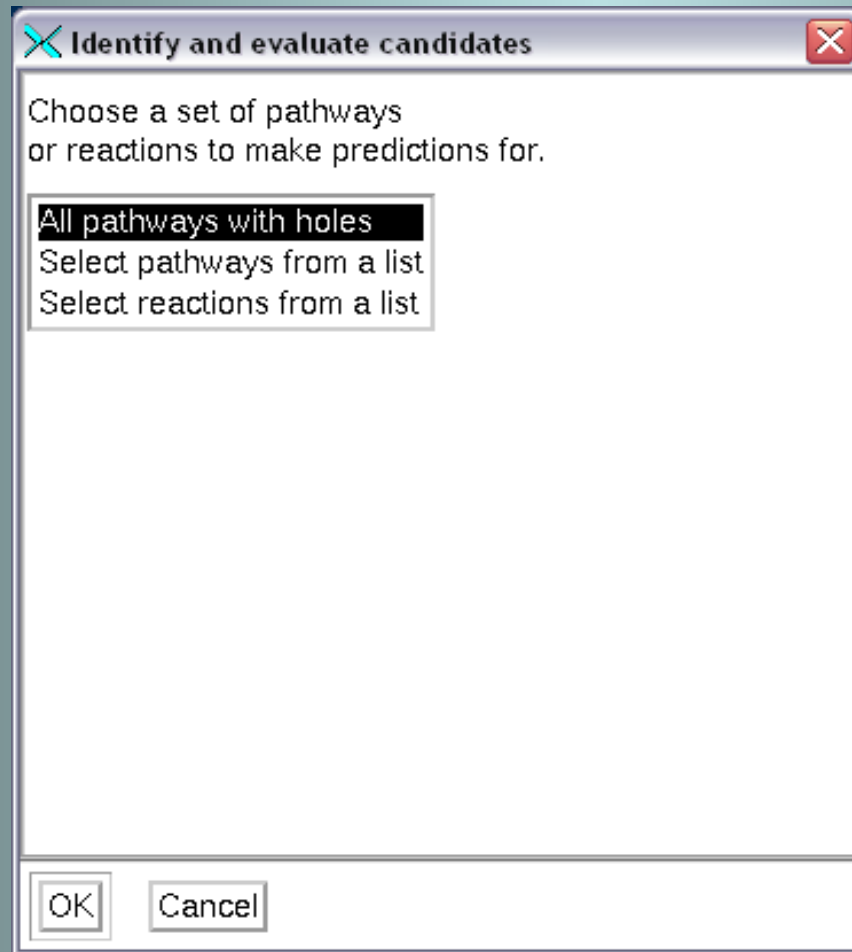
# Step 1: Prepare Training Data

Calculate training data from your organism or use existing training data.

- Once Step 1 has been completed, the training data are saved and can be reused (even in another Pathway Tools session).
- If using existing data from *E. coli* the training data are based on data from the literature.



# Step 2: Identify & Evaluate Candidates...





# Step 2: Identify & Evaluate Candidates

## Select reactions from a list

Identify and evaluate candidates

Choose a set of pathways or reactions to make predictions for.

All pathways with holes  
Select pathways from a list  
Select reactions from a list

- Malate-Dehydrogenase-(Acceptor)-Rxn
- Pyruvate-Carboxylase-Rxn
- $\beta$ -Oxoacid-Coa-Transferase-Rxn
- Methglysyn-Rxn
- Glyoxiii-Rxn
- Methylglyreduct-Rxn
- 2transketo-Rxn
- Pyruvdeh-Rxn
- Pyruvformly-Rxn
- Hydrog-Rxn

Select All      Deselect All

OK      Cancel

## Select pathways from a list

Identify and evaluate candidates

Choose a set of pathways or reactions to make predictions for.

All pathways with holes  
Select pathways from a list  
Select reactions from a list

- Entbacsyn-Pwy
- Udpnagsyn-Pwy
- Ursin-Pwy
- Udpnacetylalsyn-Pwy
- Betsyn-Pwy
- Pwy0-163
- Pwy0-181
- Salvadehypox-Pwy
- Pwy0-166
- Denovopurine2-Pwy

Select All      Deselect All

OK      Cancel



# Modes of operation...

## *Fully automatic*

- No interaction required from user
- All default values used
  - Prepare training data – all known rxns in KB
  - Identify and evaluate candidates – all pathways with pathway holes
  - Choose holes to fill in KB – all holes with  $P > 0.9$  filled
- Evidence code: “Automatic inference from sequence similarity”



# Modes of operation

## *Wizard*

Wizard prompts user for training data source and for which holes to make predictions. Wizard runs Steps 1 & 2, then prompts user to complete Step 3.

## *Power-user mode*

User must proceed through each step in order. Program still prompts user for required parameters, but each step must be completed before advancing to next step.

# Step 3: Choose Holes to Fill in KB

Choose Holes to Fill in KB

**Instructions:**

If you click on the name or description of any biological object in the table below, it will be displayed in the Navigator window.

To consider only high-probability hole-filling candidates, please specify the minimum probability that you would accept.

Minimum probability cutoff (Range: 0.0000000 to 1.0000000):

**Fill hole with top candidate?**

Holes/Reactions	Top candidate	
EC# 6.2.1.9: <b>coenzyme A + malate + ATP = phosphate + malyl-CoA + ADP</b>	<b>sucD/CC0338</b> P = 0.9884	<input checked="" type="radio"/> No <input type="radio"/> Yes, by adding function <input type="radio"/> Yes, by replacing function



Candidates to fill pathway hole: EC# 6.2.1.9

**Hole in pathway** *serine-isocitrate lyase pathway* [ Of 14 steps in this pathway, 4 are holes and 9 are present in other pathways in addition to this one. ]  
 EC# 6.2.1.9: **coenzyme A + malate + ATP = phosphate + malyl-CoA + ADP**

Show definitions

<b>Candidate hole filler</b>	CC0338-MONOMER <b>succinyl-CoA synthetase, alpha subunit</b> <input type="button" value="Move candidate to last column"/>	CC0337-MONOMER <b>succinyl-CoA synthetase, beta subunit</b> <input type="button" value="Move candidate to last column"/>
<b>Fill hole?</b>	<input checked="" type="radio"/> No <input type="radio"/> Yes, by adding function <input type="radio"/> Yes, by replacing function	<input checked="" type="radio"/> No <input type="radio"/> Yes, by adding function <input type="radio"/> Yes, by replacing function
<b>Gene</b>	CC0338 <b>sucD</b>	CC0337 <b>sucC</b>
<b>Probability</b>	0.9884	0.9748
<b>Current reactions catalyzed</b>	(none)	(none)
<b>Associated MetaCyc reactions</b>	In pathways <b>TCA cycle -- aerobic respiration,</b> <b>TCA cycle variation VIII:</b> EC# 6.2.1.5: <b>succinyl-CoA + ADP + phosphate = succinate + coenzyme A + ATP</b>	In pathways <b>TCA cycle -- aerobic respiration,</b> <b>TCA cycle variation VIII:</b> EC# 6.2.1.5: <b>succinyl-CoA + ADP + phosphate = succinate + coenzyme A + ATP</b>
<b>Average rank</b>	1.0000	1.0000
<b>Best E-value</b>	1E-180	1E-180
<b>Shotgun score</b>	1 of 3	2 of 3
<b>Average fraction aligned</b>	0.9846	0.9988
<b>Adjacent reactions?</b>	(none)	(none)
<b>Pathway direction?</b>	no	no
<b>History note</b>	<input type="text"/>	<input type="text"/>

OK





Candidates to fill pathway hole: EC# 6.2.1.9

Hide definitions

<b>Candidate hole filler</b> An enzyme that may have the function needed to catalyze the missing reaction.	CC0338-MONOMER <b>succinyl-CoA synthetase, alpha subunit</b> <input type="button" value="Move candidate to last column"/>
<b>Fill hole?</b> Should this enzyme be assigned to the missing reaction?	<input checked="" type="radio"/> No <input type="radio"/> Yes, by adding function <input type="radio"/> Yes, by replacing function
<b>Gene</b> The gene that codes the candidate enzyme.	CC0338 <b>sucD</b>
<b>Probability</b> Probability that the candidate really catalyzes the reaction.	0.9864
<b>Current reactions catalyzed</b> A list of reactions catalyzed by the candidate enzyme in this organism.	(none)
<b>Associated MetaCyc reactions</b> A list of reactions from MetaCyc that are catalyzed by the same enzyme that catalyzes the missing reaction.	In pathways <b>TCA cycle -- aerobic respiration,</b> <b>TCA cycle variation VIII:</b> EC# 6.2.1.5: <b>succinyl-CoA + ADP + phosphate =</b> <b>succinate + coenzyme A + ATP</b>
<b>Average rank</b> The average rank of the candidate enzyme sequence in the BLAST output lists (e.g., if a candidate is the best hit in each search, the average rank for the candidate is 1).	1.0000
<b>Best E-value</b> The negative log of the E-value for the best alignment of the candidate with a query sequence.	1E-180
<b>Shotgun score</b> The number of query sequences whose BLAST output included the candidate sequence.	1 of 3
<b>Average fraction aligned</b> The average of each alignment length normalized by the length of the query sequence.	0.9846
<b>Adjacent reactions?</b> Is the gene coding the candidate enzyme adjacent in the genome to one of the genes coding the enzyme for an adjacent reaction in the pathway?	(none)
<b>Pathway direction?</b> Is the candidate gene in the same direction as another gene in the same pathway; a direction is a contiguous series of genes transcribed in the same direction.	no
<b>History note</b> If desired, you may associate a history note with this enzyme. If no history note is entered, the Pathway Hole Filler will generate a note describing why this enzyme was associated with this pathway hole.	<input type="text"/>



Candidates to fill pathway hole: EC# 6.2.1.9  
 Hole in pathway **serine-isocitrate lyase pathway** [ Of 14 steps in this pathway, 4 are holes and 9 are present in other pathways in addition  
 EC# 6.2.1.9: **coenzyme A + malate + ATP = phosphate + malyl-CoA + ADP**

### Choose Holes to Fill in KB

**Instructions:**  
 If you click on the name or description of any biological object in the table below, it will be displayed in the Navigator window.  
 To consider only high-probability hole-filling candidates, please specify the minimum probability that you would accept.

Minimum probability cutoff (Range: 0.0000000 to 1.0000000):

**Fill hole with top candidate?**

Holes/Reactions	Top candidate	
EC# 6.2.1.9: <b>coenzyme A + malate + ATP = phosphate + malyl-CoA + ADP</b>	<b>sucD</b> /CC0338 P = 0.9884 <input type="button" value="Show all 7 candidates"/> Other candidates already selected: <b>sucC</b> /CC0337	<input type="radio"/> No <input checked="" type="radio"/> Yes, by adding function <input type="radio"/> Yes, by replacing function



# Output from Pathway Hole Filler - from “Identify and Evaluate Candidates” step

ROOT/ptools-local/pgdbs/user/ORGIDcyc/VERSION/reports/  
(e.g., ROOT/ptools-local/pgdbs/user/caulocyc/1.0/reports/)

- ***ORGID\_filled-holes.html*** = the list of holes that user selected to fill in the KB in Step 3.
- ***ORGIDholesX-Y.html*** (e.g., ***CAULOholes0-10.html***)
- **blasterrors.log** = log of each rxn describing whether or not any candidates were found
- **hole-data** = file containing data found for each rxn, used to generate list in “Choose holes to fill in KB” dialogue. If this file is available, step 3 can be initiated without repeating Step 2.

\* Each file is overwritten each time you run this step.



# Reference for the Pathway Hole Filler

**Green, ML and Karp, PD.**

**A Bayesian method for identifying missing enzymes in  
predicted metabolic pathway databases.**

***BMC Bioinformatics 2004, 5:76.***



# Pathway Hole Filler Demo (2)

- **once more:**
  - Refine->PHF->Step 1: Prepare Training Data**
- **In popup, select EcoCyc and say Yes to use existing Training Data**
- **Refine->PHF->Step 2: Identify Candidates**
  - **In popup, select Pathways from a List**
  - **Select Pyridnucsyn-Pwy**
- **Refine->PHF->Step 3: Choose Holes to Fill in KB**



# Pathway Hole Filler Demo (1)

## Prerequisites:

- HpyCyc installed
- BLAST installed and working
- For EcoCyc, the data/priors/ directory needed

## Demo:

- Using Power User mode, to save time
- Select HpyCyc
- Refine->PHF->Step 1: Prepare Training Data
- In popup, select HpyCyc and 2-3 reactions



# Pathway Hole Filler Demo (2)

- **once more:**
  - Refine->PHF->Step 1: Prepare Training Data**
- **In popup, select EcoCyc and say Yes to use existing Training Data**
- **Refine->PHF->Step 2: Identify Candidates**
  - **In popup, select Pathways from a List**
  - **Select Pyridnucsyn-Pwy**
- **Refine->PHF->Step 3: Choose Holes to Fill in KB**