

MetaPathway Hunter (MPH)

Alignment of Metabolic Pathways

Ron Y. Pinter
Oleg Rokhlenko
Esti Yeger-Lotem
Michal Ziv-Ukelson

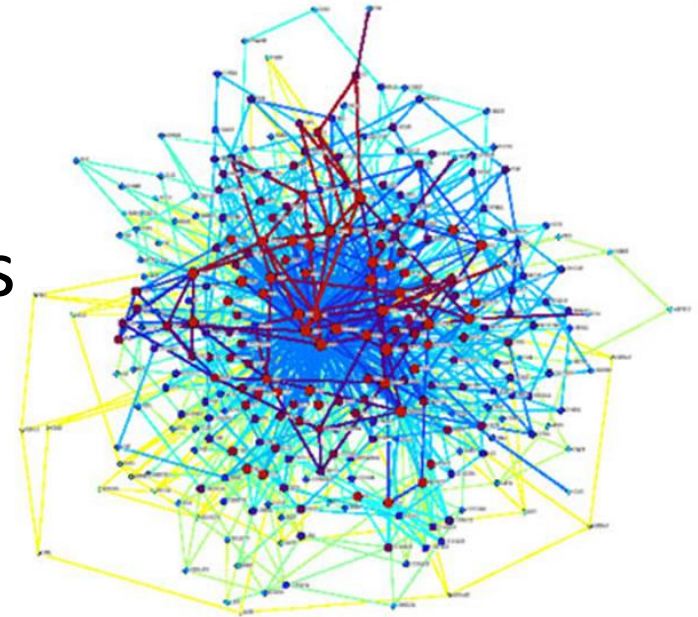


Content

- Metabolic Pathways
- Target
- Representation
- Evaluation Problem
- Approximate Labeled Subtree
Homeomorphism Problem (ALSH)
- Meta Pathway Hunter (MPH) As A Tool

Metabolic Pathways and Network

- Current efforts to reconstruct genome-scale metabolic networks
 - *E. coli*, *S. cerevisiae*, human
- Metabolic pathways databases
 - BioCyc (MetaCyc), KEGG, SGD
 - tools for pathway visualization, queries on pathway components
- Dearth of tools for comparison between pathways



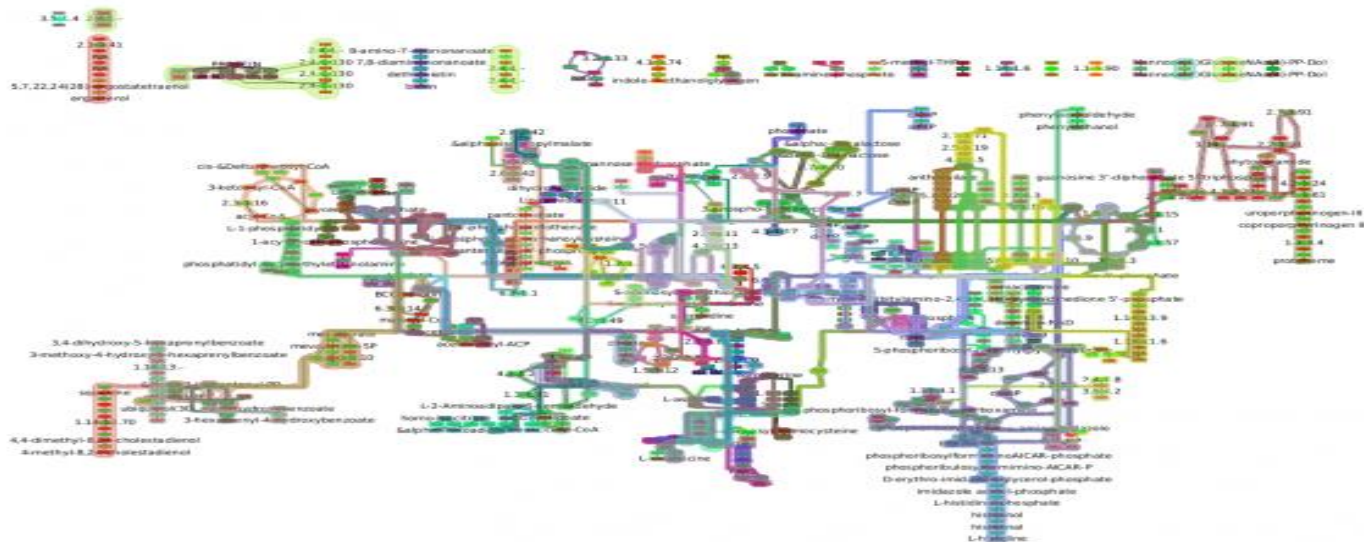
Targets

- Develop a pathway alignment tool for the detection
 - Novel pathway against a database of pathways.
 - Pathway fragment in a network.
 - Pathway conservation between species.
 - Pathway evolution within a species.
- Pathway similarity (**not necessarily identity!**) should rely on
 - pathway structure
 - enzyme similarity



Representation

- A mathematical representation is needed
 - We will use a labeled graph.
 - Nodes are **enzymes**.
 - Labels are the Enzyme Commission (**EC**).
 - Edges are **corresponding enzymes**.



Evaluation Problem

- Search for repeating patterns in a labeled graph.
- Labeled graphs pattern matching
 - Yet labeled subgraphs isomorphism is a NP-Complete problem (Garey and Johnson, 1979).
- Solution – ALSH approximate labeled subtree homeomorphism.



Convert Graphs to Trees

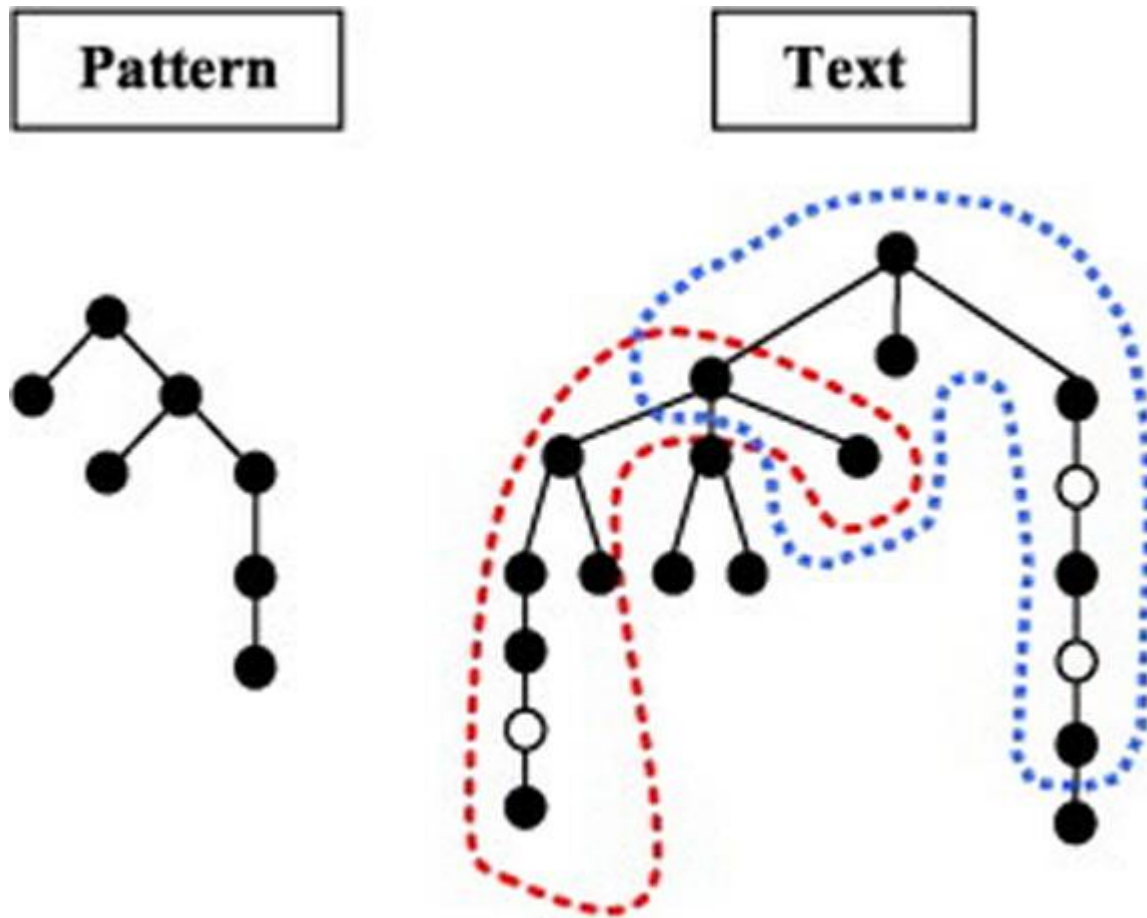
- Most metabolic pathways can be easily cast as a multi source tree or transformed to them without much loss of generality.
- Breaking directed cycles
 - generate alternative multi source trees that cover all possible cycle splitting variations.
 - surprisingly rare (less than $\frac{10}{\text{organism}}$).
- Handling DAGs (Directed acyclic graphs)
 - duplicate and split merge nodes.
 - fits well with biology: *alternative pathways*.



ALSH Content

- Subtree homeomorphism
- Definitions
- Approach
- Weighted matching (min-cost max-flows)
- Dynamic programming algorithms unordered trees
- ALSH on ordered trees

Subtree Homeomorphism

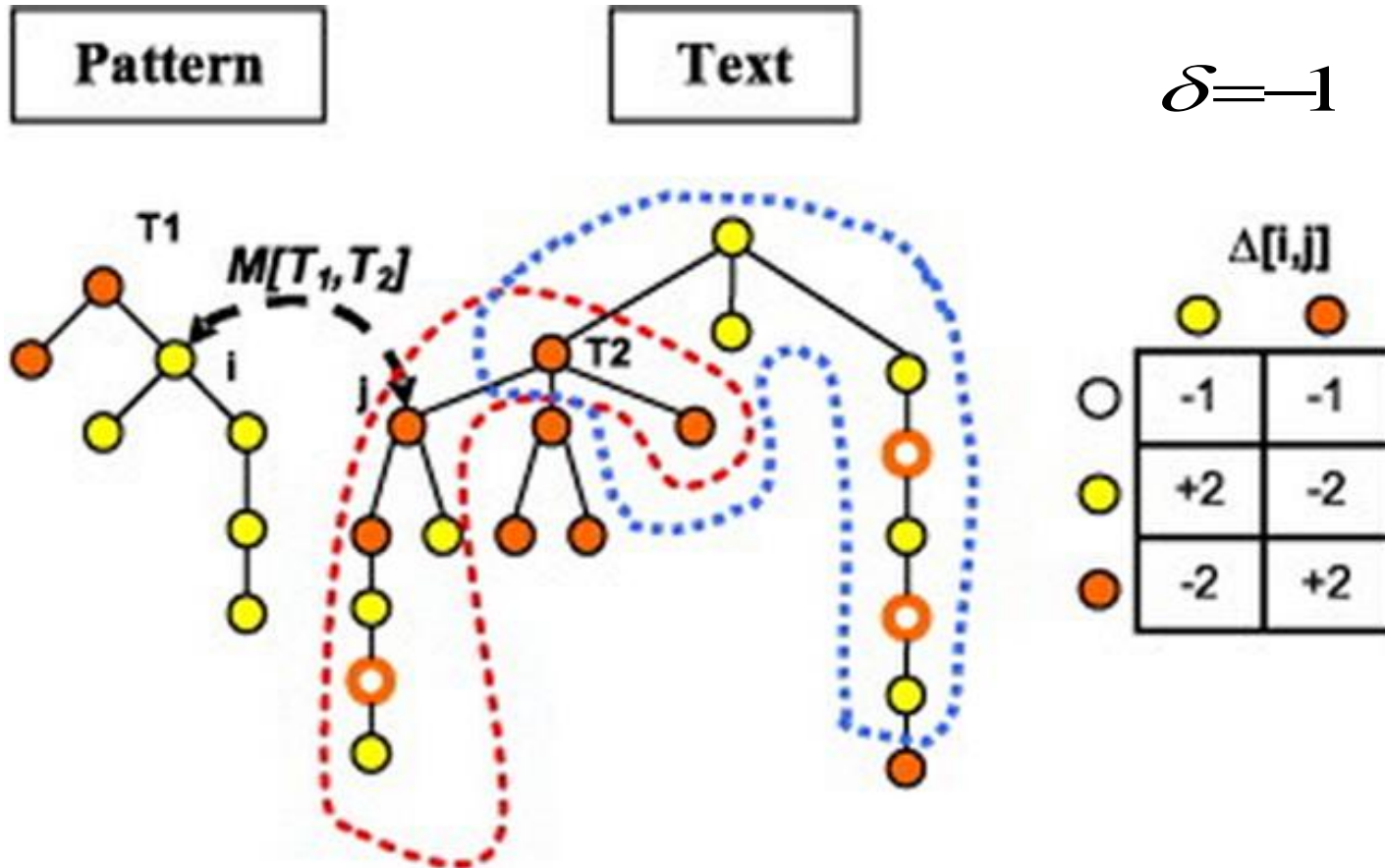


Definition I - LSH similarity score

- Let T_1, T_2 be 2 labeled trees and T_1 is homeomorphic to T_2 .
- Let δ be a predefined node deletion penalty.
- Let M be a homeomorphism from T_1 to T_2 .
- Let Δ be scoring table containing a similarity score between any labeled node in T_1 and any labeled node in T_2 .
- Let the LSH similarity score of M be

$$\sum_{u \in T_1} \Delta(u, M(u)) - \delta \cdot |T_1 - T_2|$$

Best Subtree Homeomorphism



$$LSH(M) = \delta \cdot (|T_2| - |T_1|) + \sum_u \Delta[u, M(u)]$$

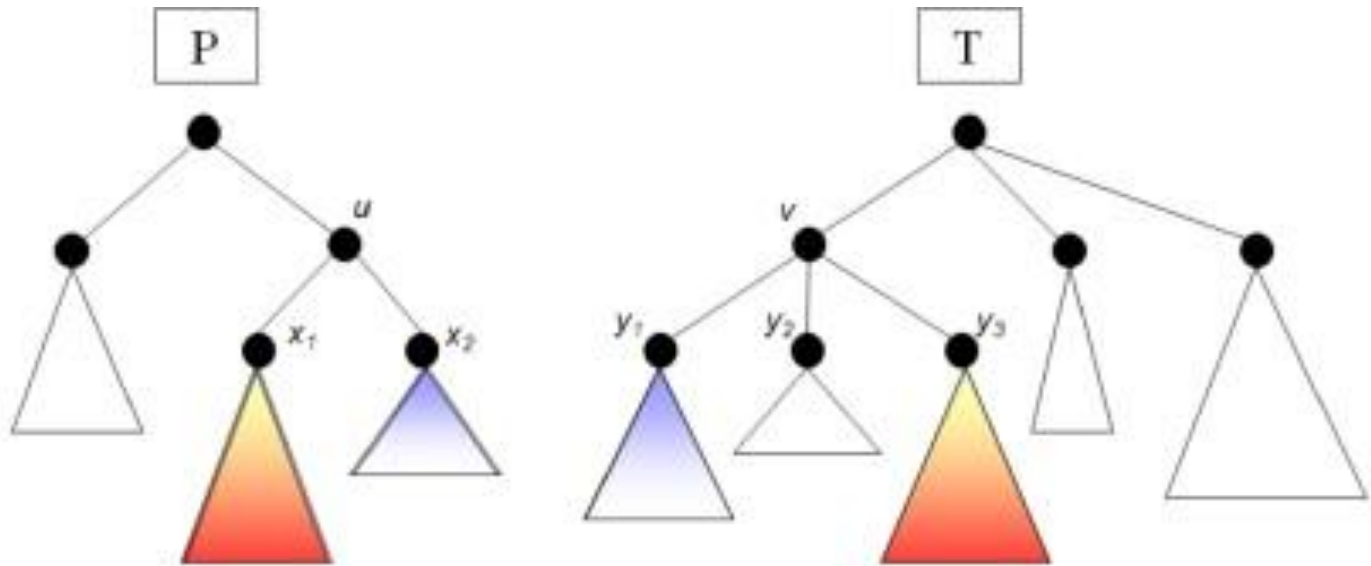
Definitions II - ALSH

- Let P, T be 2 undirected labeled trees.
- ALSH is finding the homeomorphism-preserving mapping M from P to some subtree t of T such that $\text{LSH}(M)$ is maximal.

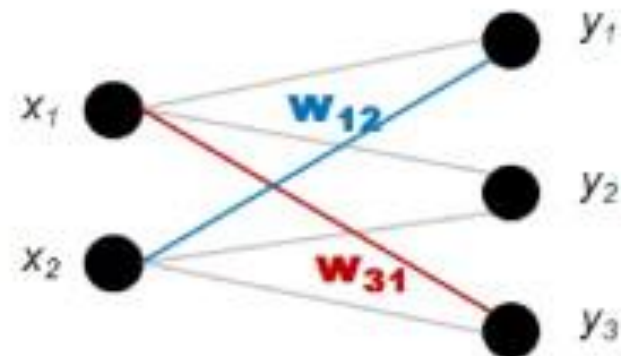
Solution Approach

- Using a recursive structure to separate the ALSH into smaller ALSH problems.
 - We will use postorder.
 - Starting from the leaves according to Δ .
- Each pair of internal nodes (u,v) induces a max weight matching problem (Min cost – Max flow).
 - Using bipartite graph ~~G~~ $G_{X,Y}$
 - X is u's children and Y is v's Children.
- The results are combined using dynamic programming.

Weight Matching



	x_1	x_2	...	u	...
y_1	w_{11}	w_{12}		w_{u1}	
y_2	w_{21}	w_{22}		w_{u2}	
y_3	w_{31}	w_{32}		w_{u3}	
...					
v					
...					



Max Weight Matching Algorithm

- Fredman & Tarjan $O(\sqrt{V}E \lg V)$ (using Fibonacci heaps).
- Under the assumption that the input costs are integers in the range $[-C, \dots, 0, \dots, C]$, Gabow & Tarjan (in “Faster scaling algorithms for network problems”) use cost scaling and blocking flow techniques to obtain an $O(\sqrt{V}E \lg(V/C))$ time algorithm. **This algorithm is conditioned by the similarity assumption** (i.e. $C = n^k$ for some constant k).
- **Lemma I** - *A flow f is minimum cost iff its residual graph R has no negative cost cycle.*

Dynamic Programming

Basic Algorithm for Rooted Unordered Trees - RScore

- Let T^r and $P^{r'}$ be 2 text trees which are rooted in r & r' .
- Let t_v^r be a subtree of T^r which is rooted in v and contains all of v 's descendants. Similarly we will use $p_u^{r'}$ a subtree of $P^{r'}$.
- $\forall v, \forall u$, the maximal LSH similarity score between $p_u^{r'}$ and some homeomorphic subtree t_v^r will be defined as $RScore[v, u]$ if such subtree does not exist $RScore[v, u] = \emptyset$.
- $RScore[v, u]$ is calculated recursively in postorder over T^r .
- $\forall v \in V_T, \forall u \in V_P$ all of u 's children will be $x_1, \dots, x_{c(u)}$ and all of v 's children will be $y_1, \dots, y_{c(v)}$.
- let $RScore[v] = \max_{1 \leq i \leq c(v)} \{ RScore[v, x_i] \}$

Require: Rooted trees $T = (V_T, E_T, r)$ and $P = (V_P, E_P, r')$.

Ensure: The root of the subtree t of T that has the highest similarity score to P , if T has a subtree which is homeomorphic to P .

```
1: for each node  $u$  of  $P$  in postorder do
2:   for each node  $v$  of  $T$  in postorder do
3:     if  $u$  is leaf then
4:       if  $v$  is leaf then
5:          $\text{RScores}(v, u) \leftarrow \Delta[v, u]$ 
6:       else
7:          $\text{RScores}(v, u) \leftarrow \text{ComputeScoresForTextNode}(v, u)$ 
8:       end if
9:     else
10:      if  $\text{Level}(u) > \text{Level}(v)$  then
11:         $\text{RScores}(v, u) \leftarrow -\infty$ 
12:      else
13:         $\text{RScores}(v, u) \leftarrow \text{ComputeScoresForTextNode}(v, u)$ 
14:      end if
15:    end if
16:  end for
17: end for
```



level > level

Procedure $\text{ComputeScoresForTextNode}(v, u)$

```
1: Let  $k$  denote the out-degree of node  $u$  and  $\ell$  denote the out-degree of node  $v$ .
```

```
2: if  $k > \ell$  then
```

```
3:   AssignmentScore  $\leftarrow -\infty$ 
```

```
4: else
```

```
5:   Construct a bipartite graph  $G$  with node bipartition  $X$  and  $Y$  such that  $X = \{x_1, \dots, x_k\}$  is the set of children of  $u$ ,  $Y = \{y_1, \dots, y_\ell\}$  is the set of children of  $v$ , and every node  $u_i \in X$  is connected to every node  $v_j \in Y$  via an edge whose weight is  $\text{RScores}(v_j, u_i)$ .
```

```
6:   Set AssignmentScore to the maximum weight of a matching in  $G$ .
```

```
7: end if
```

```
8: BestChild  $\leftarrow \max_{j=1}^{\ell} \text{RScores}(y_j, u)$ 
```

```
9: return  $\max\{\Delta[v, u] + \text{AssignmentScore}, \text{BestChild} + \delta\}$ 
```



k > l

Dynamic Programming

Basic algorithm for rooted unordered trees – best score

- Let *best_score* be the optimal LSH similarity score, it's defined ~~$best_score = \max_{v_j \in V_T} \text{sim}(P, T_{v_j})$~~ .
- $\forall v_j \in V_T$ where ~~$\text{sim}(P, T_{v_j}) = best_score$~~ v_j will be reported as a possible root of a subtree of T^r which has the maximal similarity to P under the LSH similarity score measure.



Dynamic Programming

Basic algorithm for rooted unordered trees - complexity

- Observation 1 $\sum_{u \in V} \sum_{v \in V} m(u, v)$
- For each pair $u, v \in V$ our algorithm calls a function (ComputeScoresForTextNode)
 - The greatest computational complexity comes from calculating a min cost max flow on a bipartite graph.
 - The graphs has $\alpha(v) + \alpha(u)$ nodes.
 - The graphs has $\alpha(v) \cdot \alpha(u)$ edges.
- Using Fredman and Tarjan's algorithm in the bipartite graph we compute a match for each one of u 's children invoking Dijkstra's shortest path algorithm.

$$\begin{aligned}
 & \left(\sum_{u \in V} \sum_{v \in V} (\alpha(v) \cdot \alpha(u) + \alpha(v) \cdot \alpha(u) \lg(\alpha(v))) \right) = \\
 & \left(\sum_{u \in V} (n \cdot \alpha(u) + n \cdot \alpha(u) \lg(u)) \right) = (n^2 + n \lg(n))
 \end{aligned}$$

Dynamic Programing

Expansion for unrooted unordered trees - intuition

- Why unrooted is more complicated?
 - It seems like $|T| \cdot |P|$ times the previous (rooted) solution.
- How to solve this?
 - We can choose an arbitrary root for one of the trees preferably the bigger one T.
 - We will cover all the options of roots in the pattern tree.

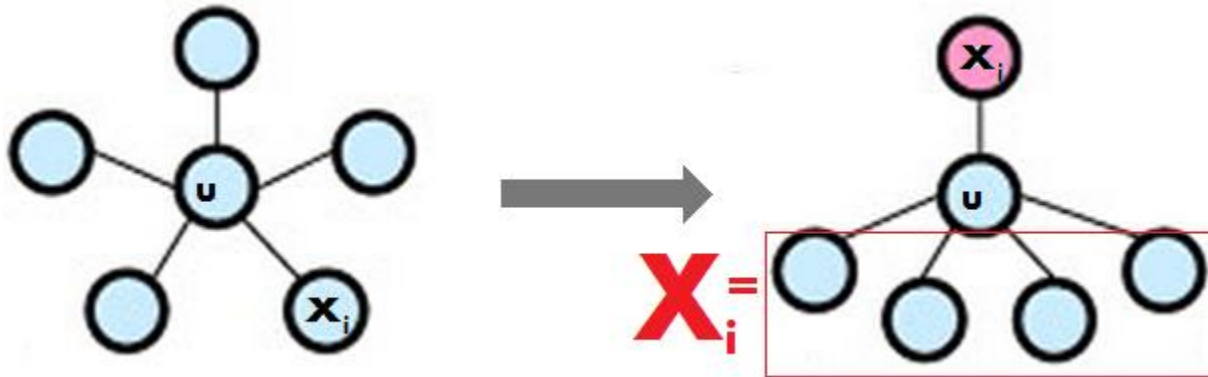


How to calculate it
with the same
complexity???

Dynamic Programming

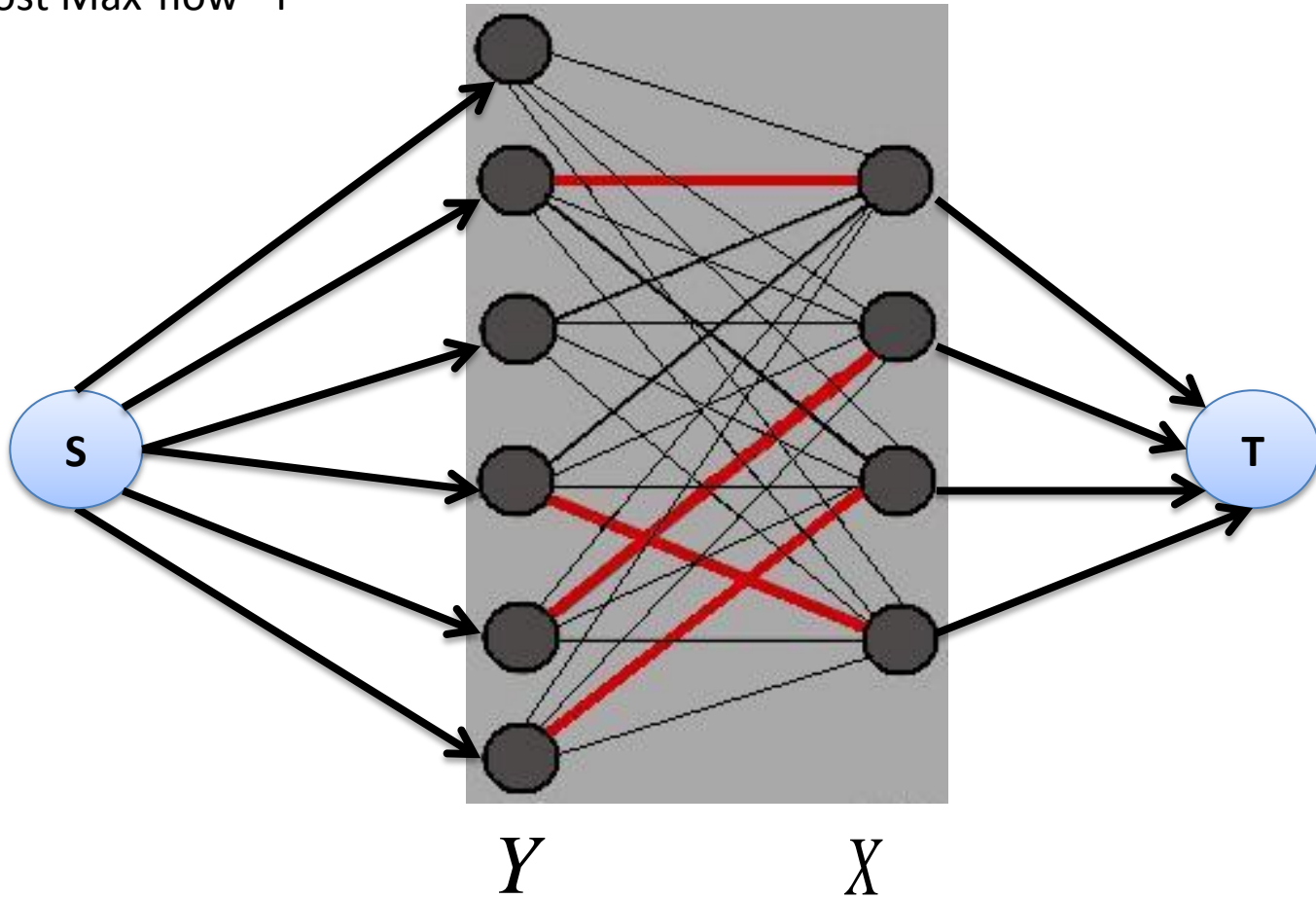
Expansion for unrooted unordered trees - G_i 's

- In unrooted trees there are only neighbors so let $d(u)$ be the number of neighbors of a node $u \in P$ and let $\{x_1, \dots, x_{d(u)}\}$ be the neighbors.
- Let X_i be ~~X_i~~
- Let G_i be ~~G_i~~ .



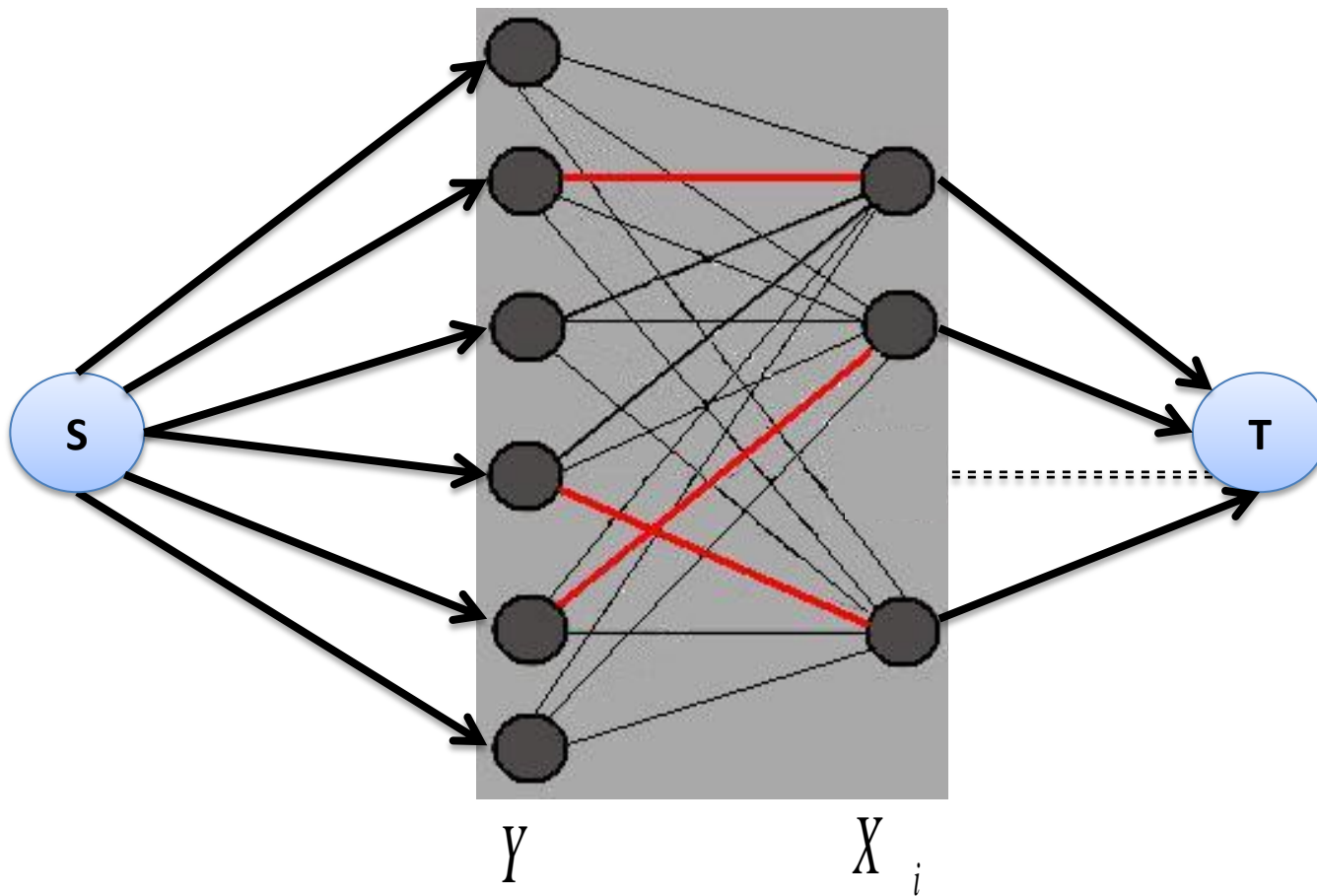
G'

— Min-cost Max-flow - f



G_i

— Flow — f'

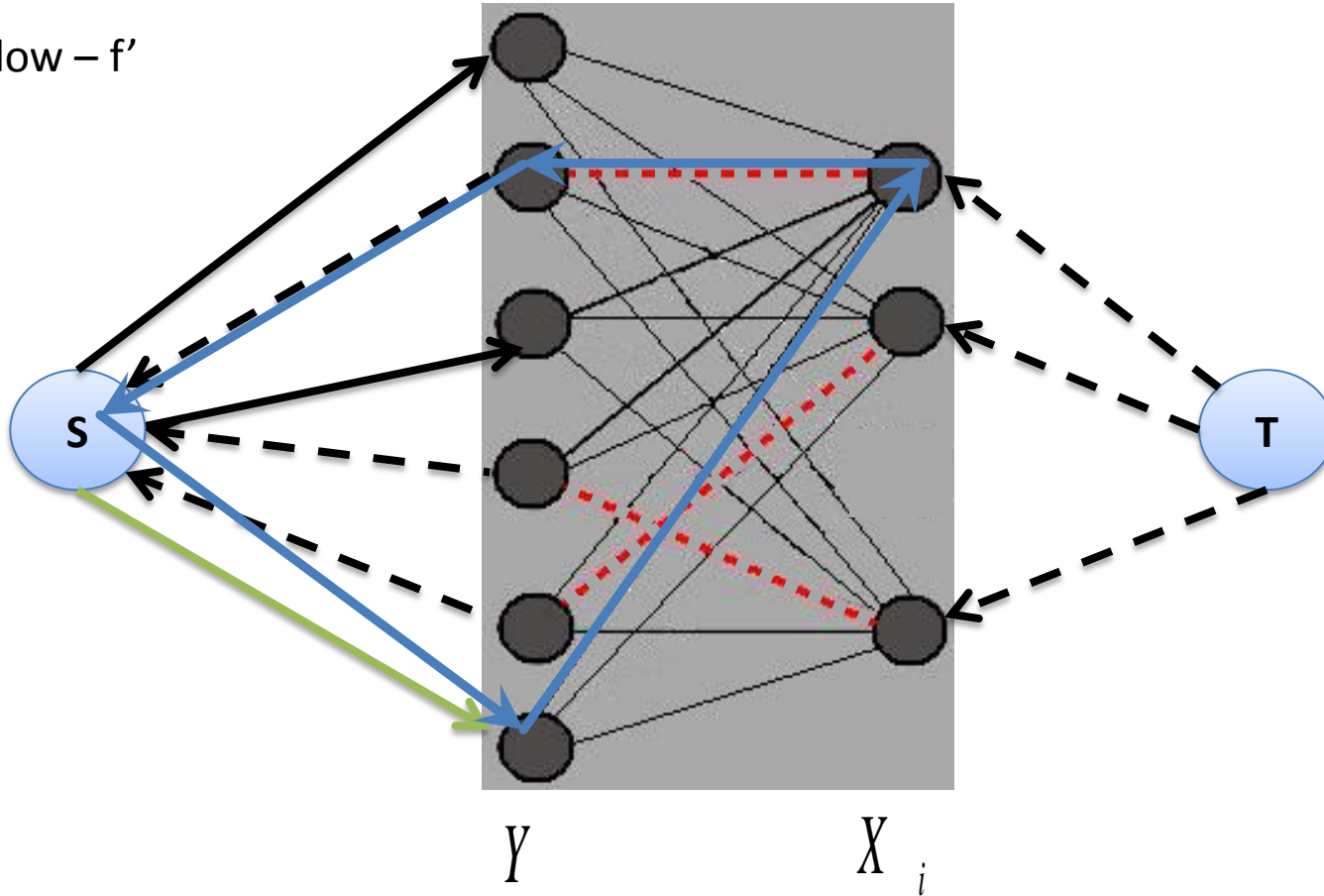


R'

→ *Correction path p*


- minimal cost cycle which contains →.

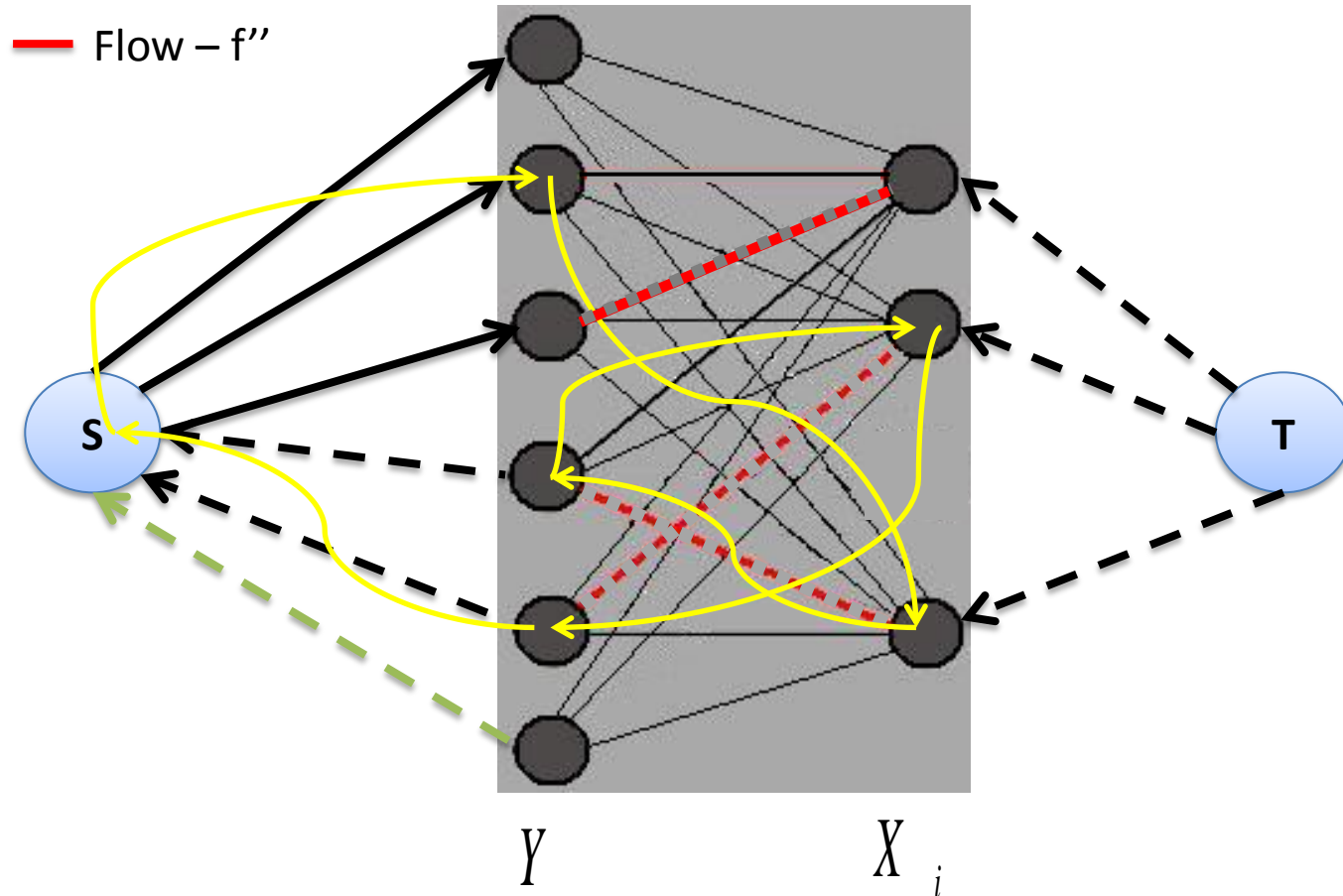
— Flow — f'



$$R''$$

Claim - the corrected residual graph (flow) has no negative cost cycles (lemma I)

 **Negative cost cycle c** - Conversely assume that c exists.



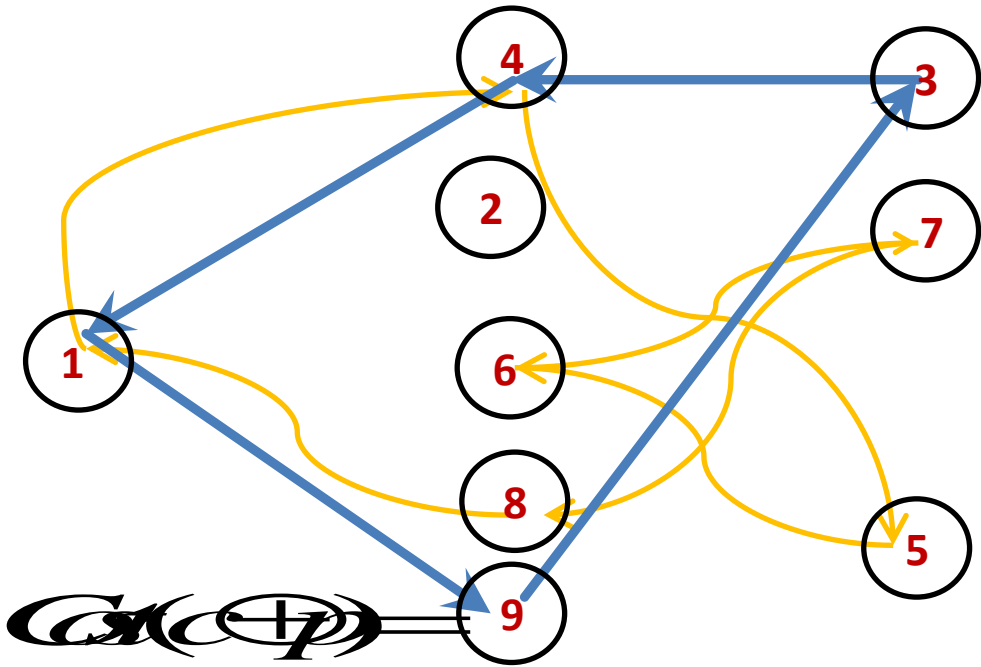
Claim - the corrected residual graph has no negative cost cycles, by **lemma I**: min-cost max-flow

→
Correction path *p*

- minimal cost cycle which contains →.

→
Negative cost cycle *c*

- Conversely assumed.



~~$$C(c \oplus p) =$$~~

~~$$C(c) - C(4) + C(p) + C(4) =$$~~

~~$$C(c) + C(p) < C(p)$$~~

Contradicts *p*'s minimalism!

Dynamic Programming

Expansion for unrooted unordered trees - UScore

- Let T, P be 2 labeled trees which are unrooted.
- Let r be an arbitrary root of T .
- $\forall v \in V_T$, the maximal LSH similarity score between $p_u^{x_i}$ and some homeomorphic subtree t_v^r will be defined as $UScore[v, u]$.
- If such subtree does not exist $UScore[v, u] = \epsilon$.
- $UScore[v, u]$ is calculated recursively in postorder over T^r .
- $\forall v \in V_T, \forall u \in V_P$ all of u 's neighbors will be $x_1, \dots, x_{d(u)}$ and all of v 's children will be $y_1, \dots, y_{c(v)}$.
- Let $UScore[v, u] = \max_{1 \leq i \leq d(u)} \{ UScore[v, x_i] \}$
- Let $UScore[v, u] = \max_{1 \leq i \leq c(v)} \{ UScore[y_i, u] \}$

Require: Unrooted trees $T = (V_T, E_T)$ and $P = (V_P, E_P)$.

Ensure: The root of the subtree t of T which has the highest similarity score to P , if T has a subtree which is homeomorphic to P .

```

1: Pick a vertex  $r$  of  $T$  to be the root of  $T$ .
2: for all  $u \in P, v \in T, x_i \in P$  do
3:   UScores[ $v, u, x_i$ ]  $\leftarrow -\infty$ .
4: end for
5: for each leaf  $v$  of  $T_r$  do
6:   for each leaf  $u$  of  $P$  do
7:     UScores[ $v, u, \text{parent}(u)$ ]  $\leftarrow \Delta[v, u]$ .
8:   end for
9: end for
10: for each internal node  $v$  of  $T$  in postorder do
11:   ComputeScoresForTextNode( $v$ )
12: end for
13:  $\text{best\_score} \leftarrow \max_{i=1, \dots, m, j=1, \dots, n} \text{UScores}[v_j, u_i, \phi]$ 

```



left ob

Procedure ComputeScoresForTextNode(v)

```

1: for each node  $u$  of  $P$  do
2:    $\text{BestChild}(v, u, x_i) \leftarrow \max_{j=1, \dots, \ell} \text{UScores}[y_j, u, x_i]$ 
3:   Construct a bipartite graph  $G$  with node bipartition  $X$  and  $Y$  such that  $X = \{x_1, \dots, x_k\}$  is the set of neighbors of  $u$ ,  $Y = \{y_1, \dots, y_\ell\}$  is the set of children of  $v$ , and every node  $x_i \in X$  is connected to every node  $y_j \in Y$  via an edge whose weight is  $\text{UScores}[y_j, x_i, u]$ .
4:   Let  $X_0 = X$  and  $X_i = X - \{x_i\}$ .
5:   for all  $1 \leq i \leq k$  do
6:     if  $\text{Level}(v) < \text{Level}(u, x_i)$  then
7:       UScores[ $v, u, x_i$ ]  $\leftarrow -\infty$ 
8:     else
9:       if  $k > \ell + 1$  then
10:        AssignmentScore( $X_i, Y$ )  $\leftarrow -\infty$ .
11:       else
12:        Compute the score AssignmentScore( $X_i, Y$ ) of the maximum weight matching in  $G_i$ .
13:       end if
14:       UScores[ $v, u, x_i$ ]  $\leftarrow \max\{\Delta[v, u] + \text{AssignmentScore}(X_i, Y), \text{BestChild}(v, u, x_i) + \delta\}$ 
15:     end if
16:   end for
17:   if  $k > \ell$  then
18:     UScores[ $v, u, \phi$ ]  $\leftarrow -\infty$ 
19:   else
20:     UScores[ $v, u, \phi$ ]  $\leftarrow \Delta[v, u] + \text{maximum weight matching in } G$ .
21:   end if
22: end for

```



Induced sub



$k > l + 1$



$k > l$

Dynamic Programming

Expansion for unrooted unordered trees – best score

- Let $best_score$ be the optimal LSH similarity score, it's defined ~~$best_score = \max_{v \in V_T} LSH(T^r, P^u, v)$~~ .
- $\forall v \in V_T, \forall u \in V_P$ where ~~$best_score = LSH(T^r, P^u, v)$~~ , v will be reported as a possible root of a subtree of T^r which has the maximal similarity to P^u under the LSH measure.



Dynamic Programming

Expansion for unrooted unordered trees – Complexity I

- Bottleneck: G_1, G_2, \dots, G_m seems $d(u)$ times more complicated than solving one G .
- The solution is that from G 's min-cost max-flow we can compute in $O(E + V \log V)$ all of the G_i 's flow solutions.
- Observation II $\sum_{u=1}^m d(u) = 2m - 2$.
- For each pair $u \in P, v \in T$ our algorithm calls a function (ComputeScoresForTextNode)
 - The greatest computational complexity comes from calculating a min cost max flow on $d(u)$ bipartite graphs.
 - The graphs has $(v) + d(u)$ nodes.
 - The graphs has $(v) \cdot d(u)$ edges.

Dynamic Programming

Expansion for unrooted unordered trees – Complexity II

- Our bottleneck can be solved in $O(E \log V)$
 - Using Fredman and Tarjan's algorithm in the bipartite graph G we compute a match for each one of u 's neighbors invoking Dijkstra's shortest path algorithm $O(E \log V)$ run time complexity.
 - For each $x_i \in X$ we create a bipartite graph G_i a match can be calculated from G 's matching using a one run of a single source shortest path algorithm $O(d(u) \cdot (E + V \log(V)))$ run time complexity.
- Therefore the total run time upper bound is

$$O\left(\sum_{u=1}^m \sum_{v=1}^n (d(u)^2 + d(u) \log(d(u)))\right) = O\left(\sum_{u=1}^m (nd(u)^2 + nd(u) \log(n))\right) = O(n^2 \sum_{u=1}^m d(u) \log(n))$$

Dynamic Programming

Expansion for unrooted unordered trees – Complexity II

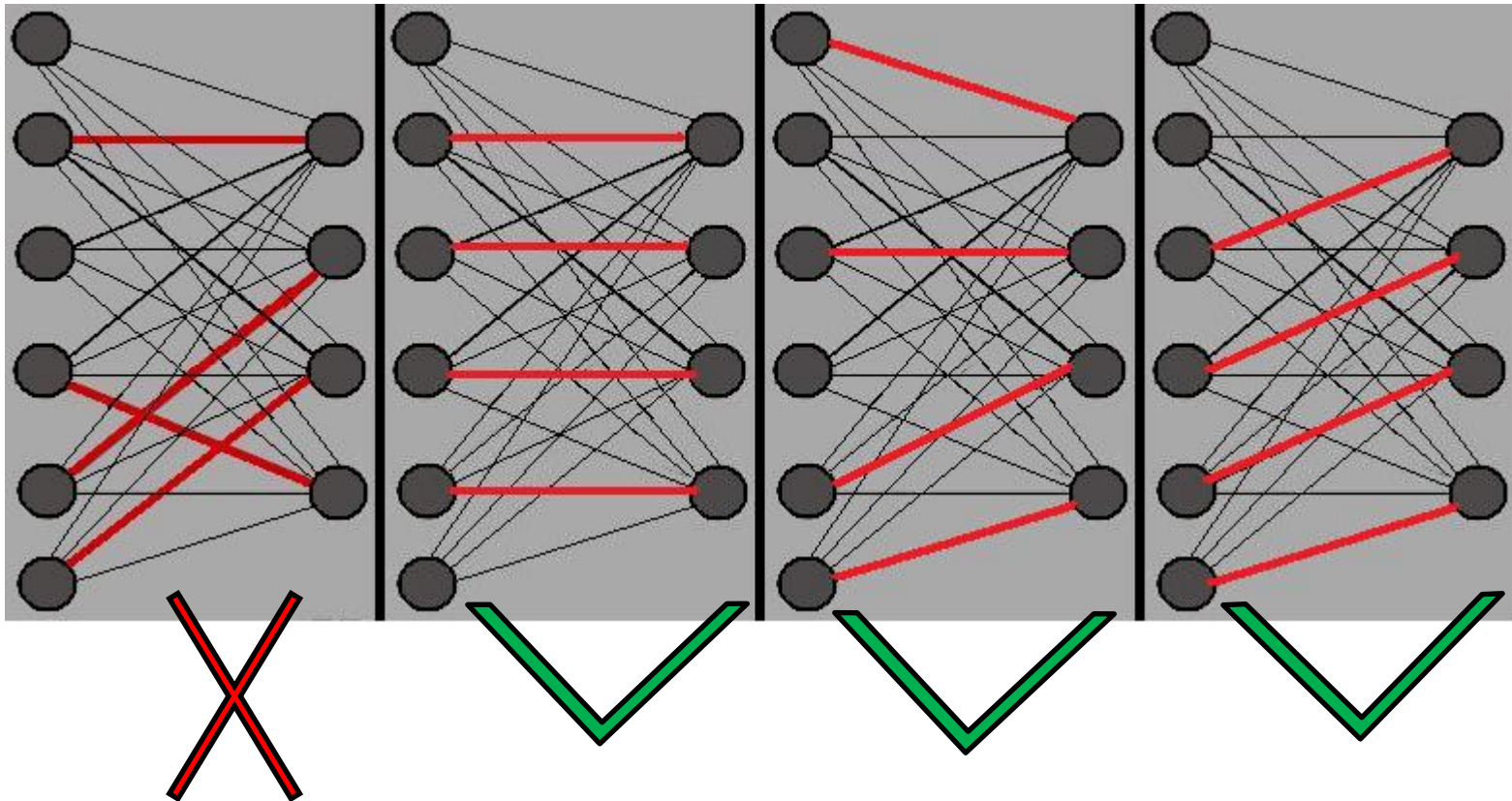
- Our bottleneck can be solved in $O(E \log V)$
 - Using Fredman and Tarjan's algorithm in the bipartite graph G we compute a match for each one of u 's neighbors invoking Dijkstra's shortest path algorithm $O(E \log V)$ run time complexity.
 - For each $x_i \in X$ we create a bipartite graph G_i a match can be calculated from G 's matching using a one run of a single source shortest path algorithm $O(d(u) \cdot (E + V \log(V)))$ run time complexity.
- Therefore the total run time upper bound is

$$O\left(\sum_{u=1}^m \sum_{v=1}^n (d(u)^2 + d(u) \log(d(v)))\right) = O(n^2 + n \log(n))$$

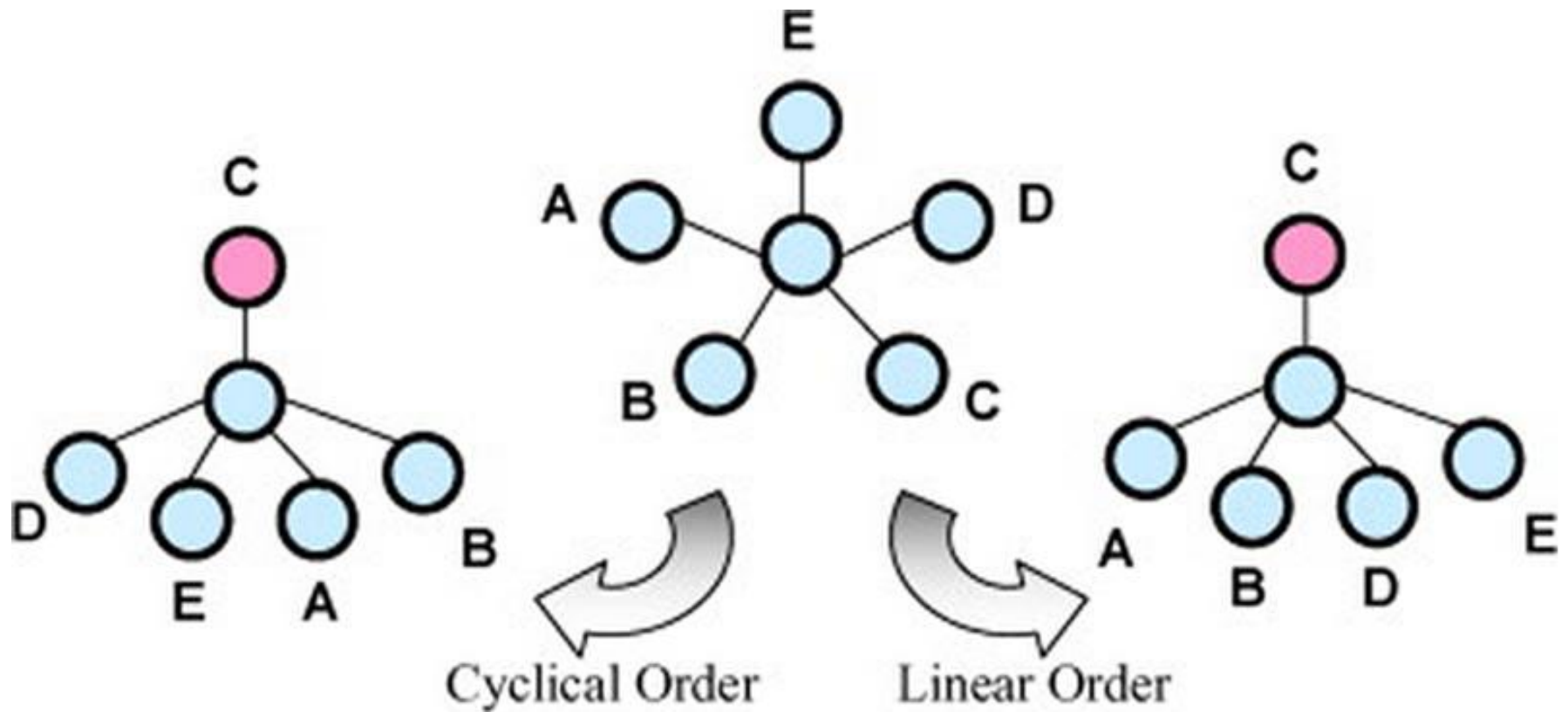
$$O\left(\sum_{u=1}^m (nd(u)^2 + nd(u) \log(n))\right) = O(n^2 + n \log(n))$$

Ordered rooted trees

- It exactly like comparing 2 strings in the matching.
- it's complexity is $O(n^2)$



Ordered unrooted trees



ALSH Summary

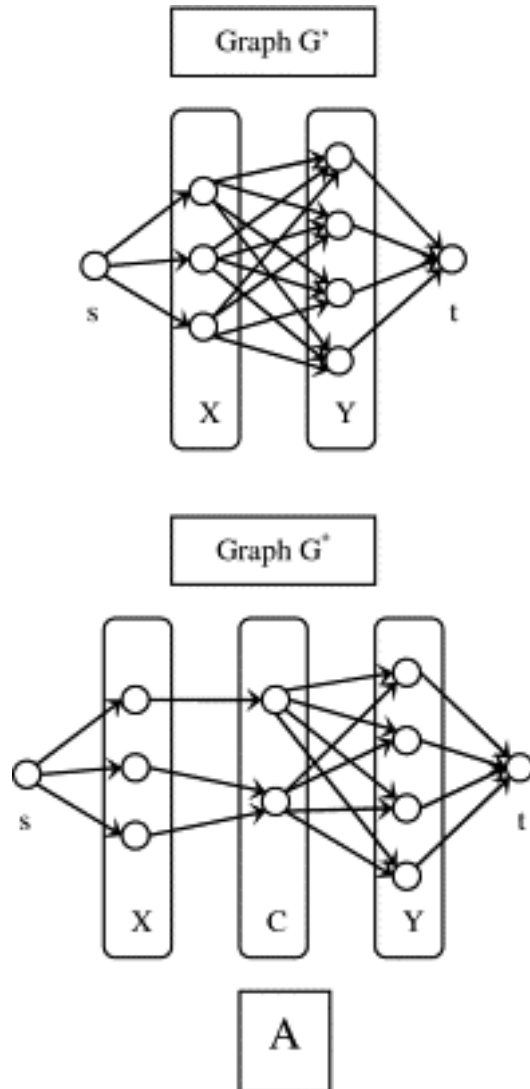
ALSH running time	ordered	unordered
rooted	$O(n \cdot m)$	$O(n \cdot m^2)$
unrooted	<div> cyclic $O(n \cdot m)$ </div> <div> linear $O(n \cdot m)$ </div>	$O(n \cdot m^2 + n \cdot m \cdot \log n)$

m - # of vertices in P

n - # of vertices in T

Dynamic Programming

more efficient ALSH algorithm for unordered trees



	Graph G^*	Graph G'
1		
2		
3		
4		
5		
6		

B

Dynamic Programming

more efficient ALSH algorithm for unordered trees

- Lemma:
 - The matching between $u \in V_p$ and $v \in V_T$ can be computed in time:

$$E + V \lg V$$

- where
 - $d(u)$ is the number of neighbors of u .
 - $D(u)$ is the number of distinct trees in the forest of trees rooted at neighbors of u .
 - $c(v)$ is the number of children of v .

Dynamic Programing

more efficient ALSH algorithm for unordered trees

- For each pair $u \in P, v \in T$ our new algorithm will work in:



- The total sum for all the pairs is :

$$O\left(\sum_{u=1}^m \sum_{v=1}^n (d(u) \cdot (clusters_u \cdot c(v) + c(v) \cdot \log(c(v))))\right)^{ds.1} =$$

$$O\left(\sum_{u=1}^m (d(u) \cdot clusters_u \cdot n + d(u) \cdot n \cdot \log(n))\right)^{ds.2} =$$

$$O\left(n \cdot \sum_{u=1}^m (d(u) \cdot D(u)) + m \cdot n \cdot \log(n)\right)$$

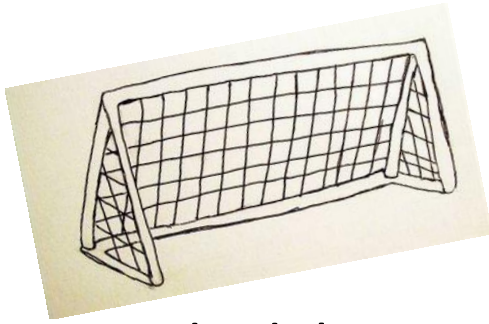
Dynamic Programming

more efficient ALSH algorithm for unordered trees

- Lemma: $\sum_{u=1}^m d(u, T(u)) = O\left(\frac{m^2}{\log m}\right)$

(similar to [Feder and Motwani 1991, Shamir and Tsur 1999])

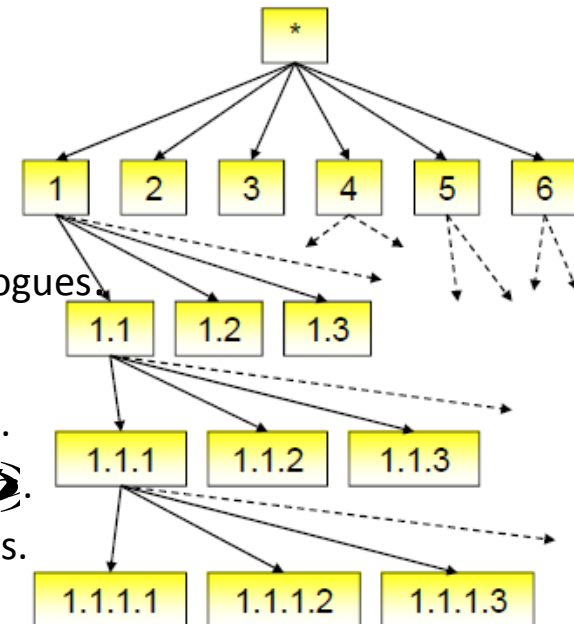
- Therefore the optimal *ALSH* solution for two rooted unordered trees run-time complexity is $O\left(\frac{n^2}{\log n} + m \log n\right)$



Alignment scoring



- Node deletion score, modeling gaps in the pattern, which entails a fixed penalty δ , a **parameter for MPH**.
 - A low value encourages alignments of evolutionary remote pathways, where only bits of the pathway are conserved.
 - A higher value encourages matching further enzymes to avoid gaps.
- Node substitution scores that are rated by a label substitution table.
 - **EC** (Enzyme Commission) :
 - Numbering system of enzymes.
 - Each enzyme is represented by a set of 4 numbers.
 - A functional classification of enzymes.
 - Enzymes with similar classification are functions homologues.
 - For an enzyme class h .
 - $C(h)$ denotes the number of enzymes included under h .
 - $I(h)$ the information content of h is $I(h) = \log_2(C(h))$.
 - For enzymes e_i, e_j , h_{ij} is the lowest common upper class.
 - The similarity score of e_i, e_j is $I(h_{ij})$.



Statistical Significance

- The p-value cutoff used in the analysis is 0.01.
- Pathway pairs with at least one statistically significant alignment between them as significantly aligned.
- To assess whether the number of significantly aligned pathway pairs in the comparisons deviate significantly from the number expected by pure chance at a cutoff of 0.01.
- The binomial test (k, n, p) is used per comparison to ensure. This test computes the probability of having at least k (successes) significantly aligned pairs in n (experiments) total number of aligned pathway pairs with probability $p(=0.01)$ for success.
- This test was performed using the R project for Statistical Computation.

MetaPathway Hunter As A Tool

- Searches a pattern **P** against a database of labeled trees F .
- Each alignment is given a score **s** and statistical significance **p**.
- Statistical significance **p** computation:
 - Executes **P** against a database of 100 randomized F .
 - $p(P)$: the percentage of randomized F which scored more than **s**.
 - $p = \frac{\text{Number of randomized } F \text{ with score} \geq s}{100}$
- Enables visualization of the alignment, which is color-coded by similarity.

Each random graph has the same nodes with the same degrees and the same number of edges



Implementation Details

- Implementation

- Code in C++
- Java-based GUI to allow usage as a website.

- System requirements

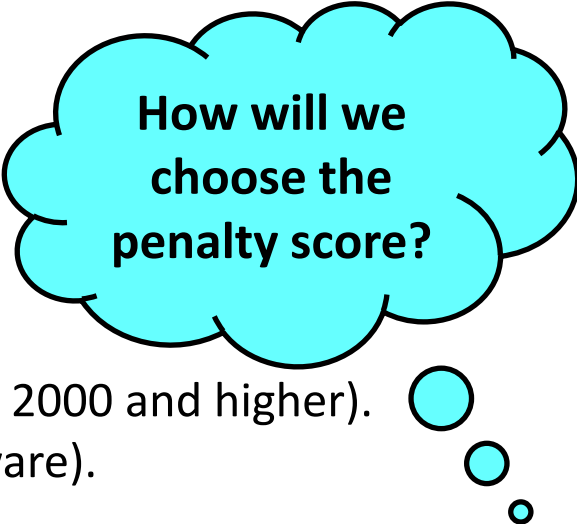
- Any Intel Pentium-based computers.
- Microsoft Windows operating system (version 2000 and higher).
- No further requirements (hardware and software).

- Query response

- The 5 best matches per pathway sorted by score & statistical sig.
- An HTML file as a visual aid with the query drawn on the aligned MPs.

- Query input

- The penalty score δ .
- Pattern ***P*** in the processed tree form.
- A number of metabolic pathways (MPs) ***T*** (it can accept more than one ***T***) in the processed tree form.



How will we
choose the
penalty score?

Performance

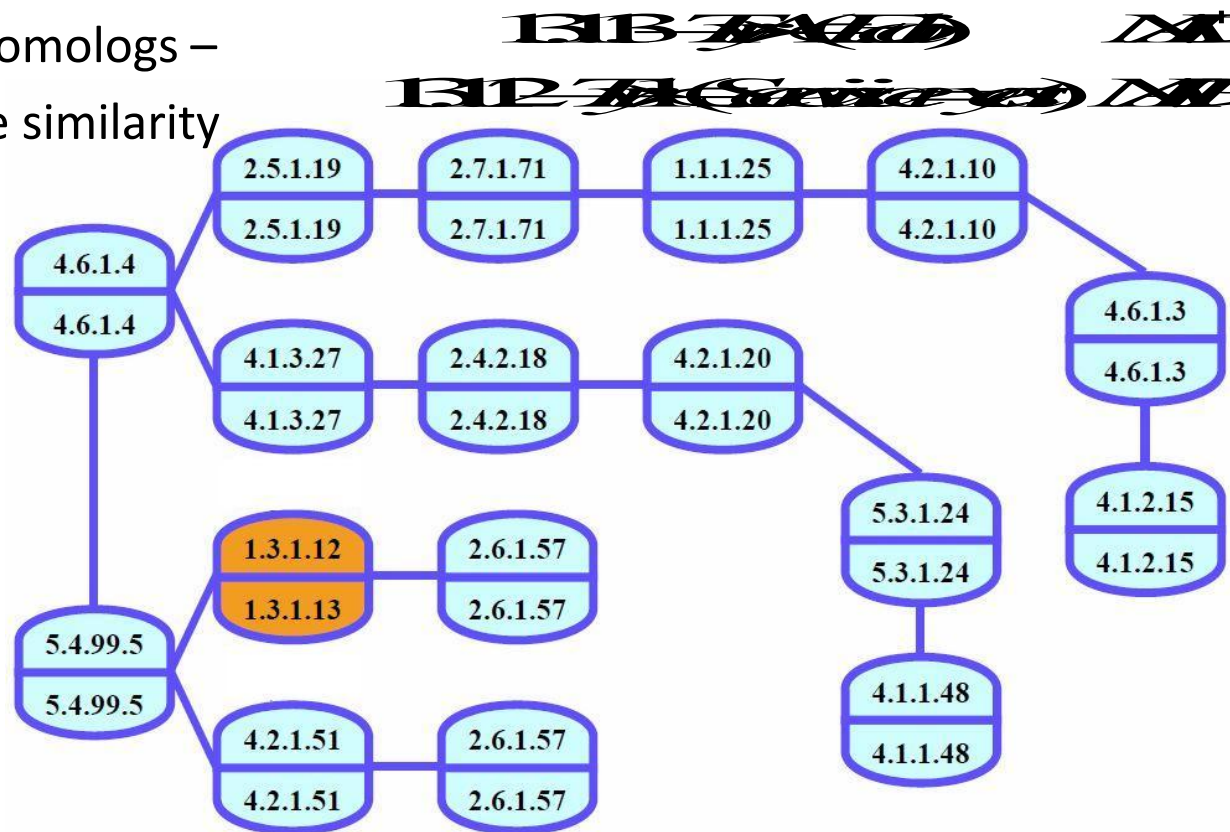
- All against all alignments were run between *E.coli* (prokaryotic) and *S.cerevisiae* (eukaryotic) .
- Took 3.66 hours on a Pentium 4, 2.6 GHz clock, 512 MB RAM computer.
- An average of 47 s per query.

Inter-Species Alignments I

- Conservation of metabolic pathways:
 - The alignments between 62 out of 80 pairs of analogous pathways were significant.
 - Conservation is not limited to small pathways, *e.g.*

functional homologs –
no sequence similarity

Biosynthesis of
phenylalanine,
tyrosine, and
tryptophan,
 $s=-4.28$; $p<0.01$

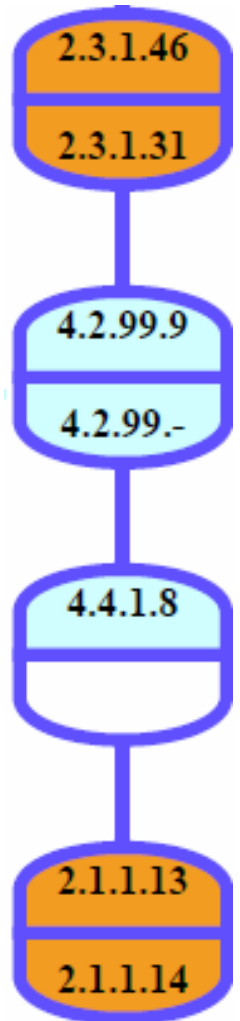


Inter-Species Alignments II

- Alignment with gaps

- A gap found
- Both sequences do the same thing.
 - May be gene fusion in yeast.
 - Or gene duplication in E.coli.

~~4299~~ ~~*AMB*~~ (~~*E. coli*~~)
~~4418~~ ~~*MC*~~ (~~*E. coli*~~)
~~4292~~ ~~*MT7*~~ (~~*Saccharomyces cerevisiae*~~)

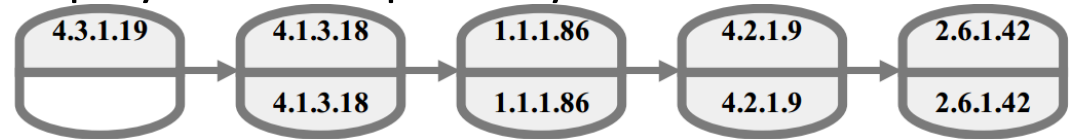


Homoserine and methionine biosynthesis, *E. coli* vs. yeast:
 $s = -13.15$, $p < 0.01$

Intra-Species Alignments

- Hints to evolution of paralogous pathways
- Amino acid biosynthesis pathways are conserved
- Valine, Leucine, Isoleucine have similar biosynthesis pathways.
- Valine vs. isoleucine employ the same pathway but on different substrates

$s = 0, p < 0.01$

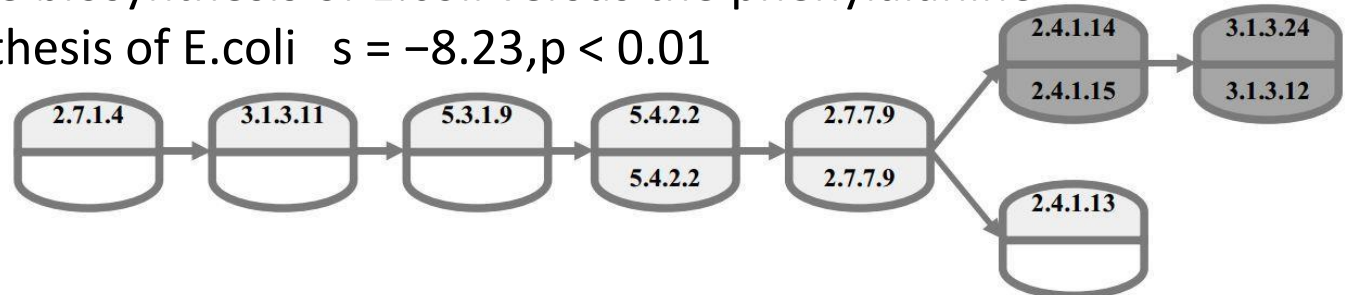


- The trehalose anabolism pathways of yeast versus the sucrose biosynthesis pathway of yeast

$s = -9.58, p < 0.01$



- tyrosine biosynthesis of E.coli versus the phenylalanine biosynthesis of E.coli $s = -8.23, p < 0.01$



MetaPathway Queries

Flexible searching – possible queries

- Search for a metabolic connection between enzymes (source and product query).
- Search for a metabolic pathway using only a part of it (common pathway core query)

