

Graemlin: General and robust alignment of multiple large interaction networks

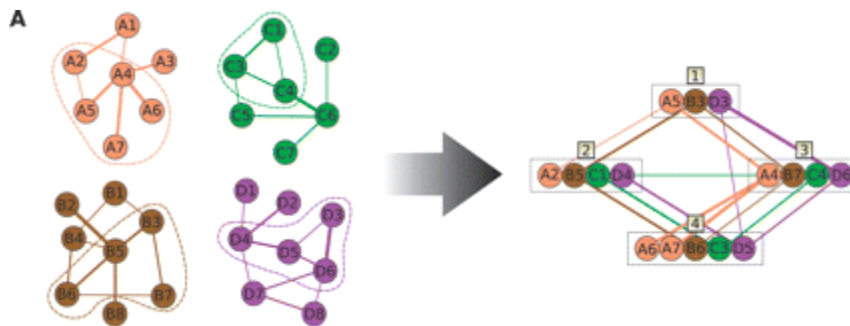
Flannick et al.

Introduction

- Graemlin: BLAST for ppi networks
 - Fast
 - Heuristic
 - Dual or multi alignments
- First quantitative comparison of network aligners

Methods

- What is a graph alignment algorithm?
 - Grouping proteins into equivalence classes
 - Determining an alignment's score
 - Choosing what to align



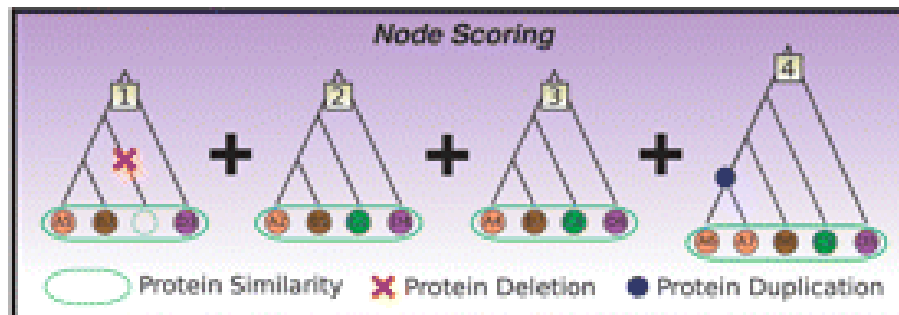
Methods: Scoring

- Log ratio of model versus background
- Score the equivalence class
- Score the edges

Methods: Node Scoring

- Model trained on pairs from COGs (Clusters of Orthologous Groups)
- Random pairs as background
- Score based on mutations/duplications/insertions etc.
- Node score is sum-of-pairs

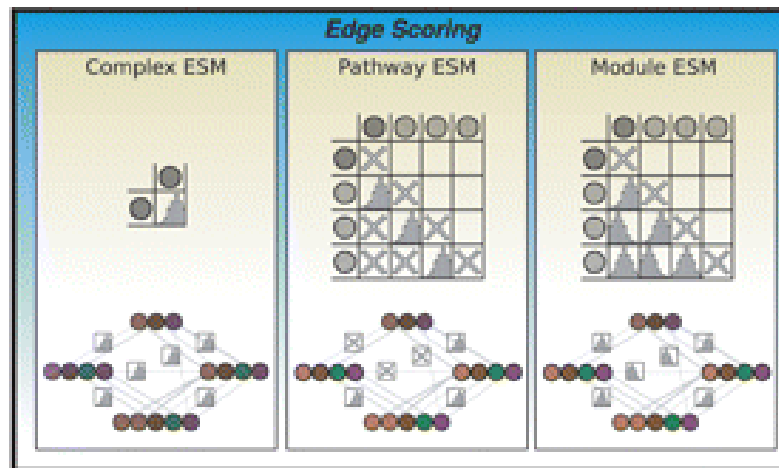
B



Methods: Edge Scoring

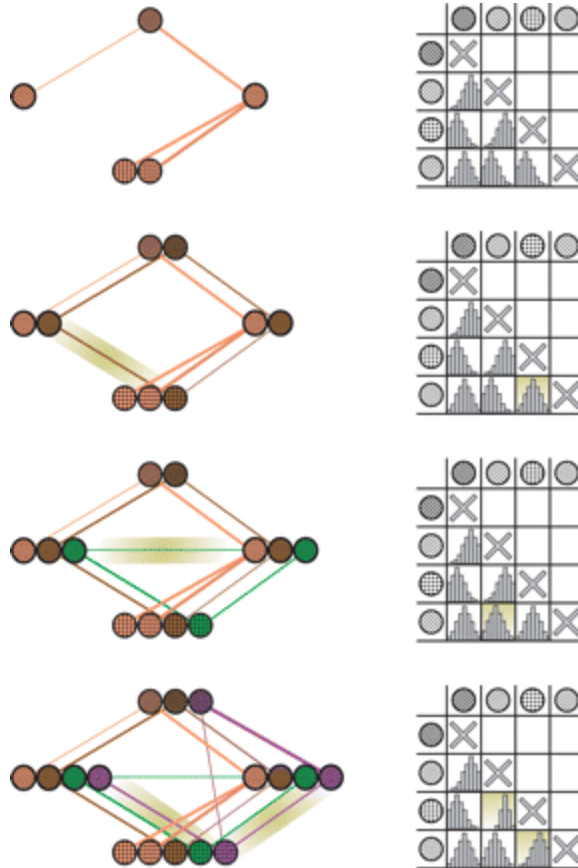
- Edge scoring matrix determines edge weights
- Different scoring methods prioritize different kinds of clusters

C



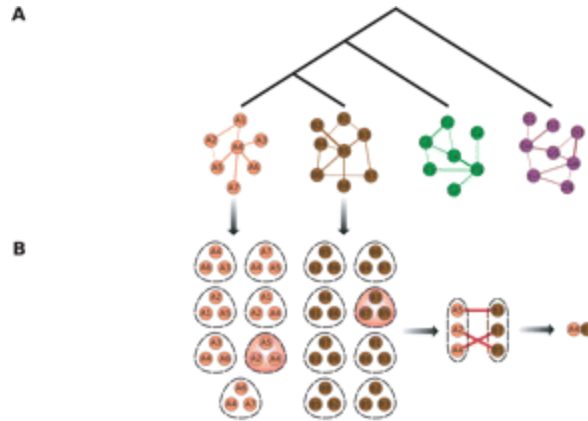
Methods: Edge Scoring

- Refines the matrix as more species are added



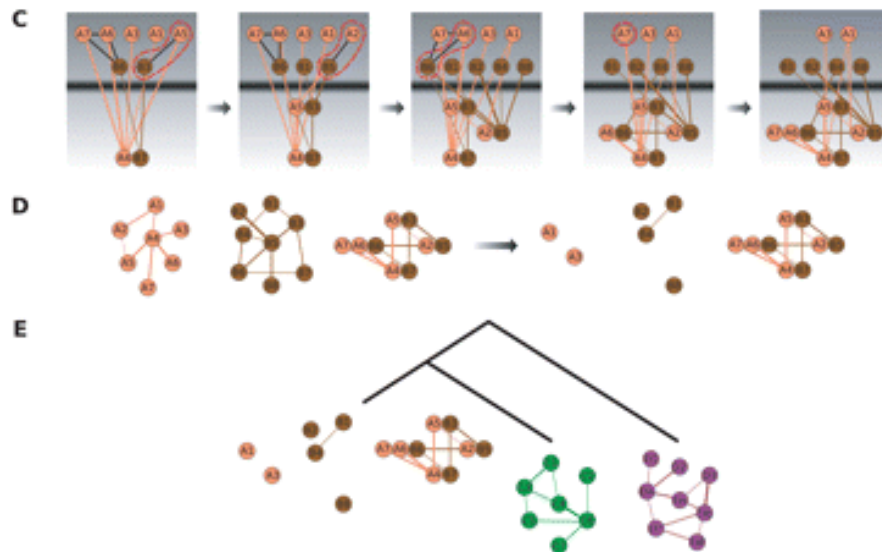
Methods: Alignment

- Pairwise
 - D-cluster = k-mer
 - In ONE species
 - Heuristic!
 - Join pairs of d-clusters over some threshold



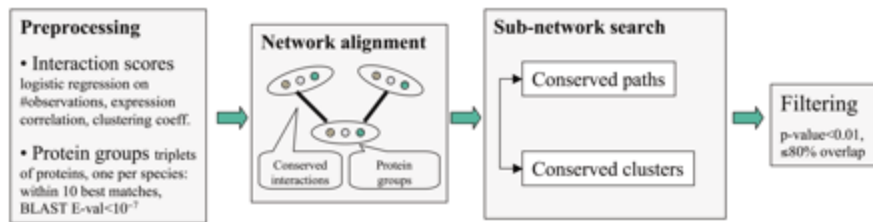
Methods: Alignment

- Expand the aligned d-cluster greedily
 - Stop when no further addition increases the score
- Multiply align by using our combined network as a single network



Results

- Compared to NetworkBLAST and MaWISH
 - NetworkBLAST



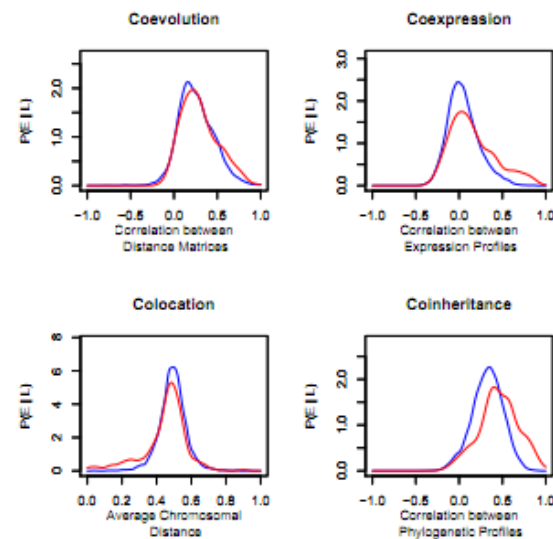
(Sharan et al 2004)

– Mawish

- Builds a fixed alignment graph, starts from nodes and greedily finds maximally scoring subgraphs

Results: Building a network with SRINI

- SRINI uses four factors
- Not from direct experimental data
- Recapitulates KEGGS in prokaryotes only



(b) Evidence vs. Training Set

(Srinivasan et al 2006)

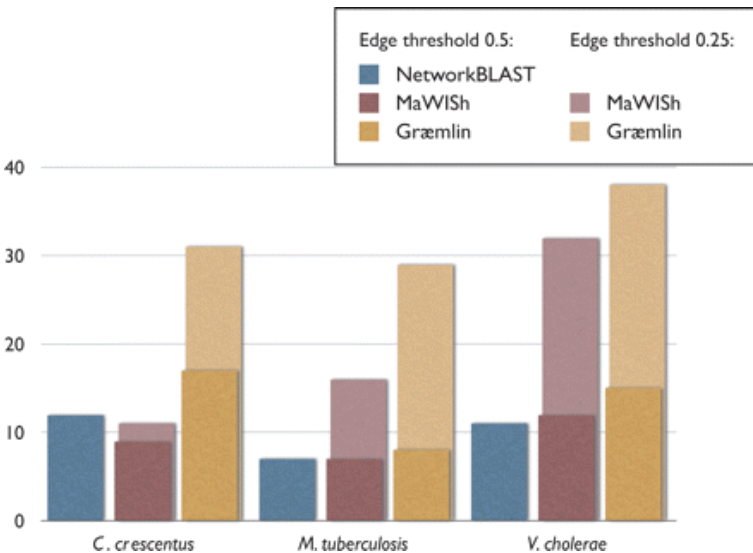
Results

- KEGGs are manually annotated pathways
- Sensitivity: what is a hit?
 - ‘aligned at least three proteins in each species to their correct counterparts in all other species’
 - Sensitivity
- Also a specificity measure: percent of alignments that are GO-enriched

Results: Pairwise

Table 3. Results on pairwise alignment of complete networks thresholded at 0.5

		KEGGs hit	KEGG coverage	Alignments enriched	Running time (sec)
<i>E. coli</i> vs. <i>C. crescentus</i>					
	MaWISH	9 (20%)	32%	72%	3
	NetworkBLAST	6 (14%)	28%	61%	9624
		12 (27%)	49%	72%	
	Græmlin	15 (34%)	47%	68%	21
		17 (39%)	45%	67%	11
<i>E. coli</i> vs. <i>M. tuberculosis</i>					
	MaWISH	7 (13%)	20%	85%	3
	NetworkBLAST	7 (13%)	24%	88%	301
		7 (13%)	32%	88%	
	Græmlin	8 (15%)	36%	89%	11
		8 (15%)	39%	89%	8
<i>E. coli</i> vs. <i>V. cholerae</i>					
	MaWISH	12 (31%)	35%	64%	3
	NetworkBLAST	10 (26%)	35%	58%	8797
		11 (28%)	41%	64%	
	Græmlin	19 (49%)	48%	75%	13
		15 (38%)	55%	74%	12
<i>E. coli</i> vs. <i>S. coelicolor</i>					
	MaWISH	N/A	N/A	N/A	N/A
	NetworkBLAST	6 (14%)	23%	46%	122,168
		10 (23%)	67%	95%	
	Græmlin	8 (19%)	58%	88%	734
		9 (21%)	59%	85%	829



For each pair of species, we performed complete network-to-network alignment using MaWISH and Græmlin. For each tested method, shown, from left, is the total number of KEGG pathways hit by an alignment, the fraction of KEGG pathways hit by an alignment, the average coverage of a KEGG pathway, the percentage of enriched alignments, and the total running time. We calculated the average coverage of KEGGs with respect to only those KEGGs that an aligner hit, and measured running time in CPU-seconds.

Results: Multiple

Table 4. Results on multiple alignment of complete networks

		KEGGs hit	KEGG coverage	Alignments enriched	Running time (sec)
0.25 threshold					
<i>E. coli</i> vs. <i>C. crescentus</i> vs. <i>V. cholerae</i>					
Græmlin	Pathway	27 (57%)	68%	72%	329
	Complex	29 (62%)	71%	79%	251
<i>E. coli</i> vs. <i>C. jejuni</i> vs. <i>H. pylori</i>					
Græmlin	Pathway	16 (57%)	57%	87%	44
	Complex	17 (61%)	63%	89%	43
0.5 threshold					
<i>E. coli</i> vs. <i>C. crescentus</i> vs. <i>V. cholerae</i>					
NetworkBLAST	Pathway	N/A	N/A	N/A	>10 ⁶
	Complex				
Græmlin	Pathway	7 (26%)	67%	72%	63
	Complex	9 (33%)	62%	75%	38
<i>E. coli</i> vs. <i>C. jejuni</i> vs. <i>H. pylori</i>					
NetworkBLAST	Pathway	5 (33%)	41%	94%	32,394
	Complex	4 (27%)	38%	96%	
Græmlin	Pathway	3 (20%)	74%	82%	12
	Complex	3 (20%)	72%	79%	12

We performed three-way multiple network alignment using NetworkBLAST and Græmlin; the columns in this table are analogous to those in Table 3.

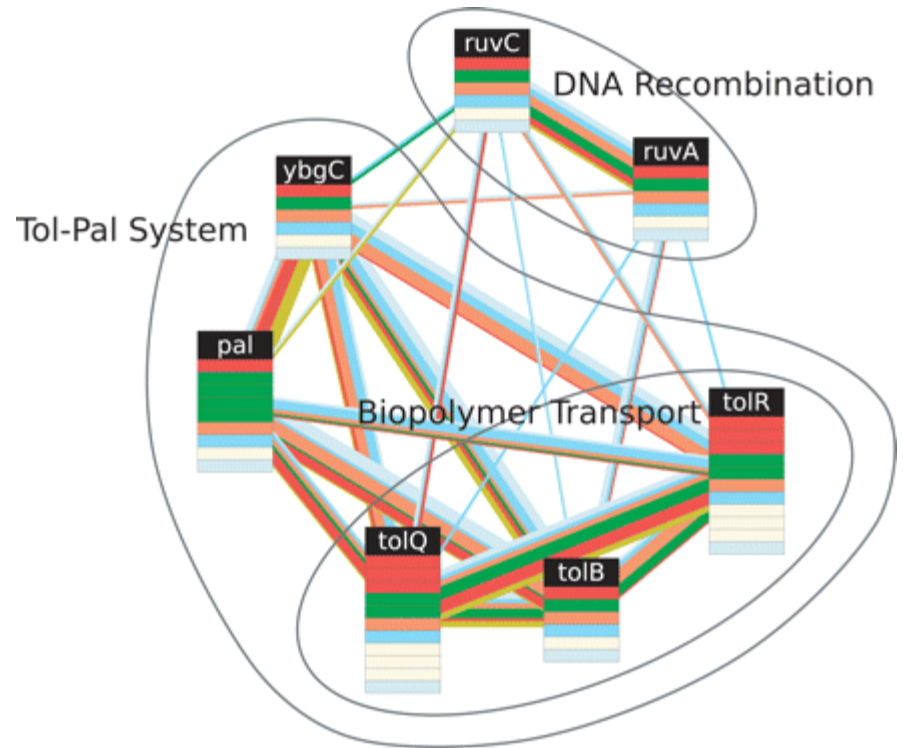
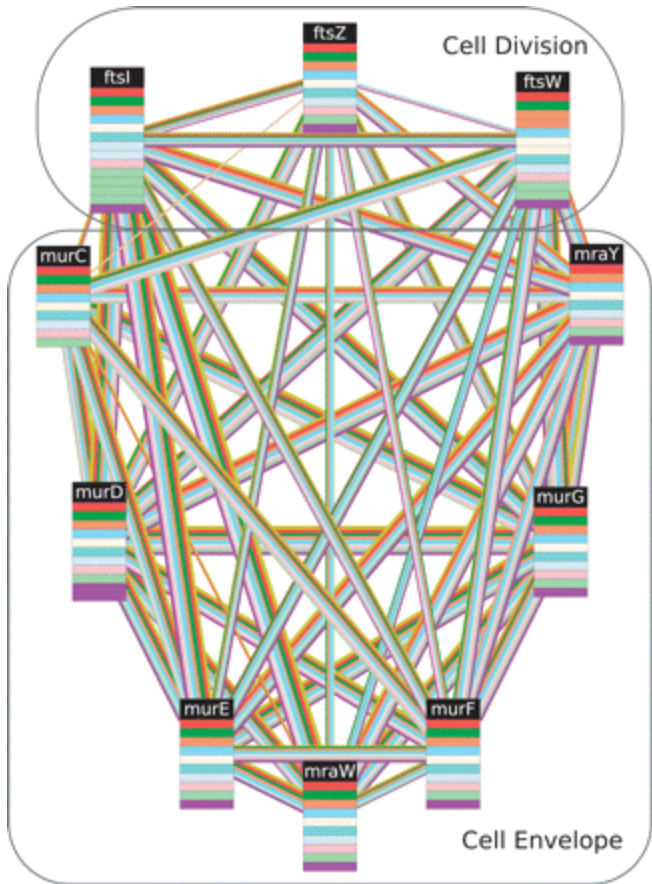
Results: Query

- Complex performs well?!
 - ‘most KEGGs that are alignable are highly connected, making the Complex ESM close to optimal’

Table 5. Results on alignment of a query network to a database thresholded at 0.5

		KEGGs hit	KEGG coverage	Running time (sec)
<i>E. coli</i> vs. <i>C. crescentus</i>				
	MaWISH	15 (34%)	31%	37
	NetworkBLAST			3453
		8 (18%)	32%	
		10 (23%)	49%	
	Græmlin			17
		20 (45%)	45%	
		20 (45%)	47%	3
		20 (45%)	48%	23
<i>C. crescentus</i> vs. <i>E. coli</i>				
	MaWISH	9 (20%)	32%	130
	NetworkBLAST			4788
		10 (23%)	37%	
		10 (23%)	41%	
	Græmlin			6
		15 (34%)	39%	
		15 (34%)	42%	5
		15 (34%)	42%	33
<i>E. coli</i> vs. <i>M. tuberculosis</i>				
	MaWISH	10 (19%)	19%	93
	NetworkBLAST			3947
		12 (22%)	23%	
		12 (22%)	29%	
	Græmlin			3
		17 (31%)	31%	
		17 (31%)	35%	3

Results: Biological discovery



Conclusion

- Built a network alignment algorithm and a repeatable set of assessments
- Demonstrated ability to recover alignments from prokaryotes
- Builds on the BLAST model of alignment in a new context
- Fast; greedy; scales well