

Tesis Doctoral
Improving Co-evolution Based
Methods for Protein-Protein
Interaction Prediction



David Ochoa

Universidad Autónoma de Madrid

Supervisor: Florencio Pazos

Systems Biology and PPI Networks

- Complexity of biological processes
 - Holistic point of view
 - Network approaches
 - Protein-Protein Interactions

Experimental methods for detecting PPIs



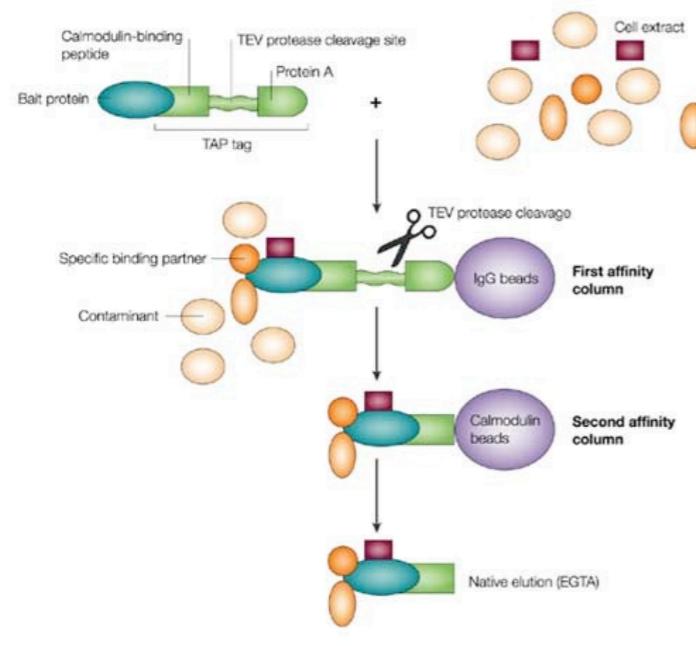
- Yeast two-hybrid

- Affinity purification - Mass Spectrometry

- Synthetic lethality

- Protein chips

- Phage display



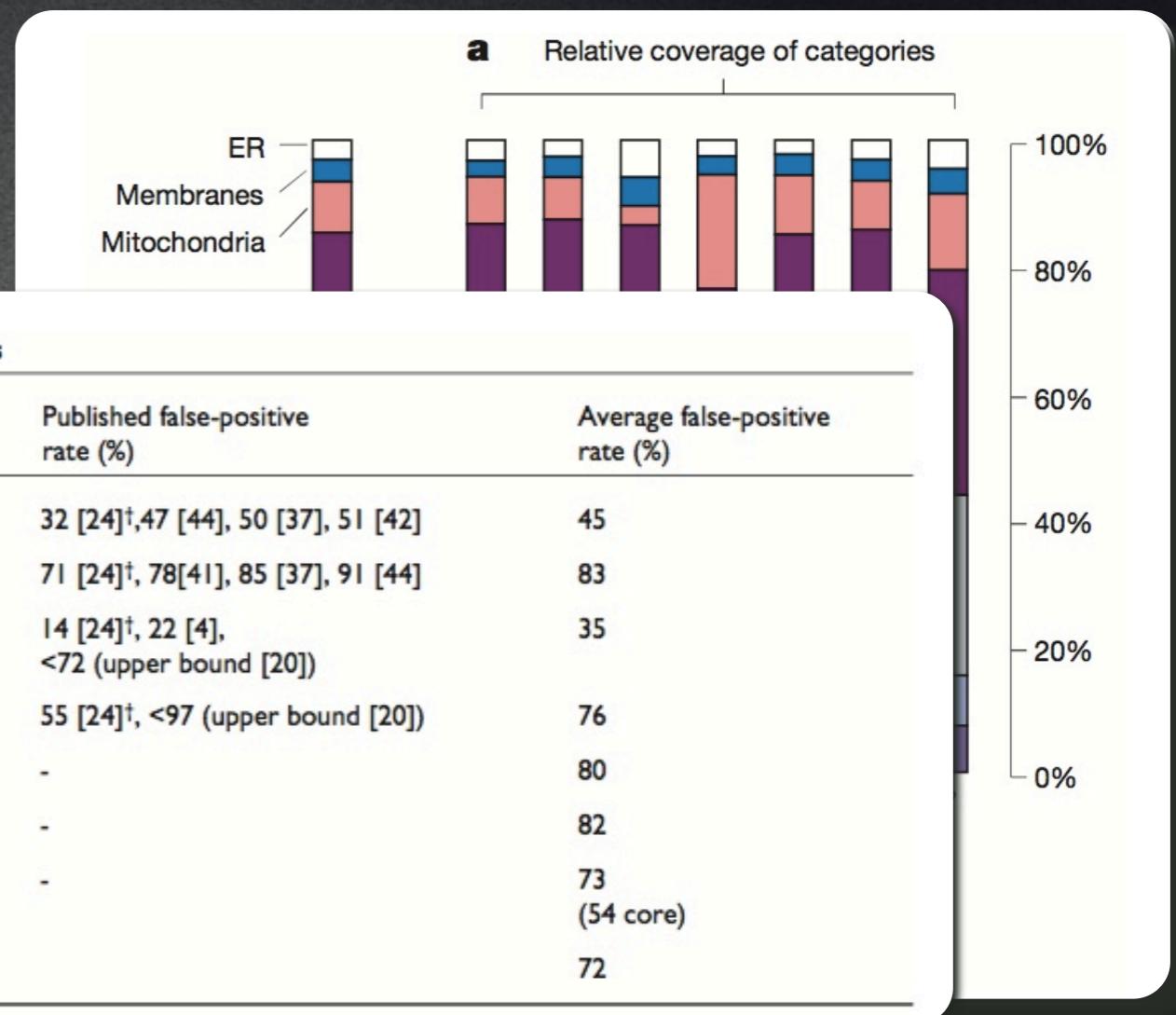
Nature Reviews | Molecular Cell Biology

Experimental problems

- High False-positives Rates

Yeast protein-interaction assay false-positive rates: yeast datasets

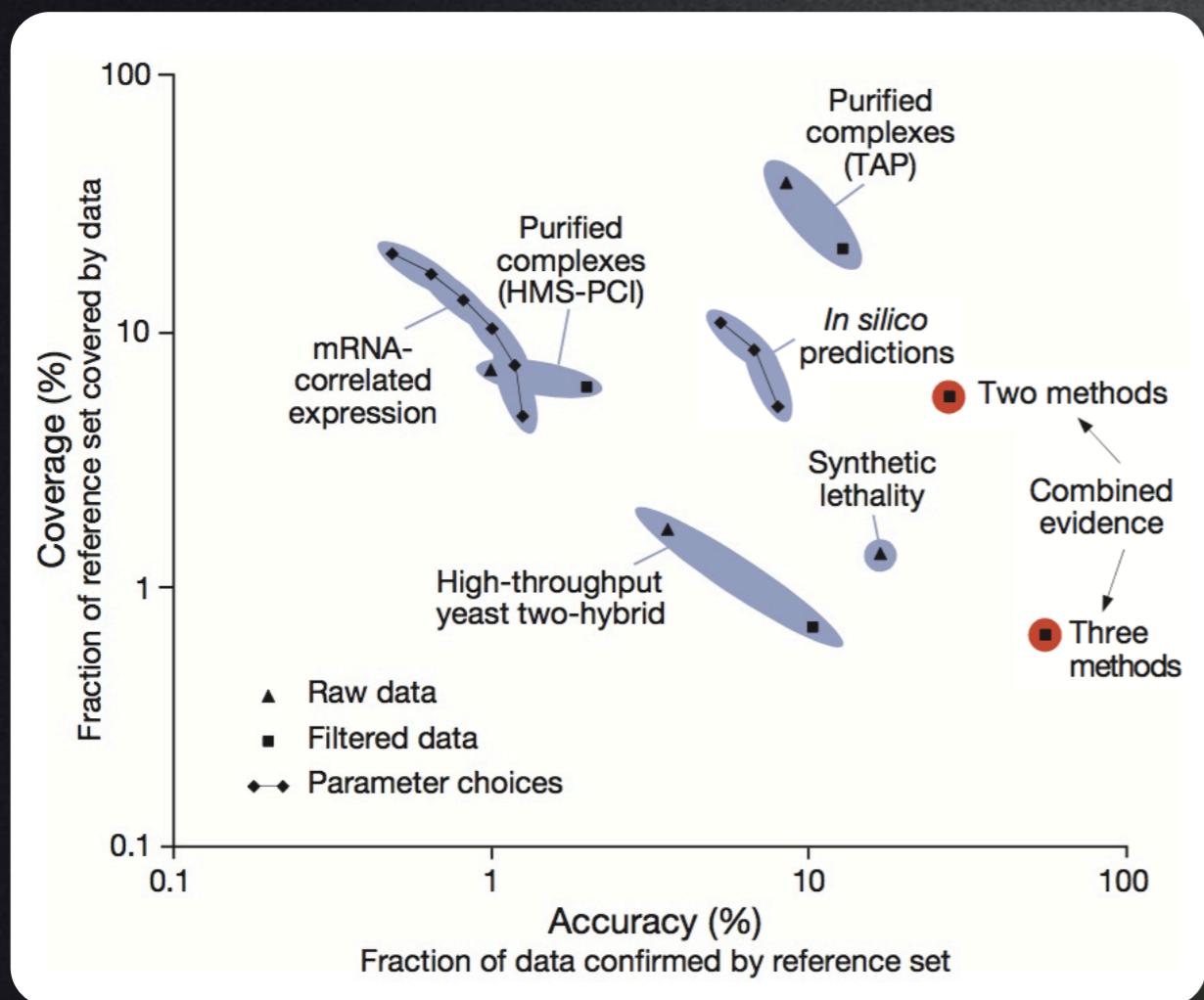
Dataset	Number of interactions	Derived false-positive rate* (%)	Published false-positive rate (%)	Average false-positive rate (%)
Uetz <i>et al.</i> [35]	854	46 [32]	32 [24] [†] , 47 [44], 50 [37], 51 [42]	45
Ito [36]	4,393	89 [32]	71 [24] [†] , 78[41], 85 [37], 91 [44]	83
Gavin <i>et al.</i> [16]	3,180	68 [32]	14 [24] [†] , 22 [4], <72 (upper bound [20])	35
Ho <i>et al.</i> [17]	3,618	83 [32], 81, 82, 80	55 [24] [†] , <97 (upper bound [20])	76
Jansen <i>et al.</i> [22]	15,922	81, 79	-	80
Gavin <i>et al.</i> [27]	18,137	78, 82, 86 [‡]	-	82
Krogan <i>et al.</i> [28]	14,317 (7,123 core)	75, 79, 66 [‡] (59, 65, 37 [‡] core)	-	73 (54 core)
Overall	51,419			72



Hart, G. T., Ramani, A. K., Marcotte, E. M. (2006) How complete are current yeast and human protein-interaction networks? *Genome Biol* 7, 120.

Mering, von, C. et al. Comparative assessment of large-scale data sets of protein-protein interactions. *Nature* 417, 399–403 (2002).

Possible Solutions



- Repeated screenings
- Combining evidences
- Confidence evaluation

Mering, von, C. et al. Comparative assessment of large-scale data sets of protein-protein interactions. Nature 417, 399–403 (2002).

Computational Methods

- Gene Fusion Events
- Conservation of Gene Neighborhood
- Similarity of Phylogenetic Profiles
- Similarity of Phylogenetic Trees



Enright et al. Protein interaction maps for complete genomes based on gene fusion events. Nature (1999) vol. 402 (6757) pp. 86-90

Overbeek et al. Use of contiguity on the chromosome to predict functional coupling. In Silico Biol (1999) vol. 1 (2) pp. 93-108

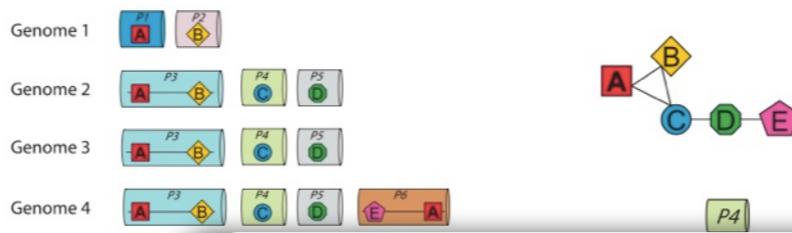
Dandekar et al. Conservation of gene order: a fingerprint of proteins that physically interact. Trends Biochem Sci (1998) vol. 23 (9) pp. 324-328

Pellegrini et al. Assigning protein functions by comparative genome analysis: protein phylogenetic profiles. Proc Natl Acad Sci U S A (1999) vol. 96 (8) pp. 4285-4288

Date y Marcotte. Discovery of uncharacterized cellular systems by genome-wide analysis of functional linkages. Nat Biotechnol (2003) vol. 21 (9) pp. 1055-1062

	1	2
organism 1	1	1
organism 2	1	1
organism 3	0	1
organism 4	1	1
organism 5	1	1
organism 6	1	0
organism 7	1	1
organism 8	1	1
organism 9	1	1
organism 10	1	1

Local Profiles: Domains

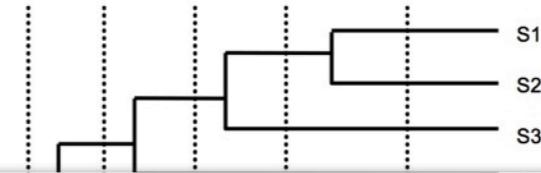


Selection of Organisms

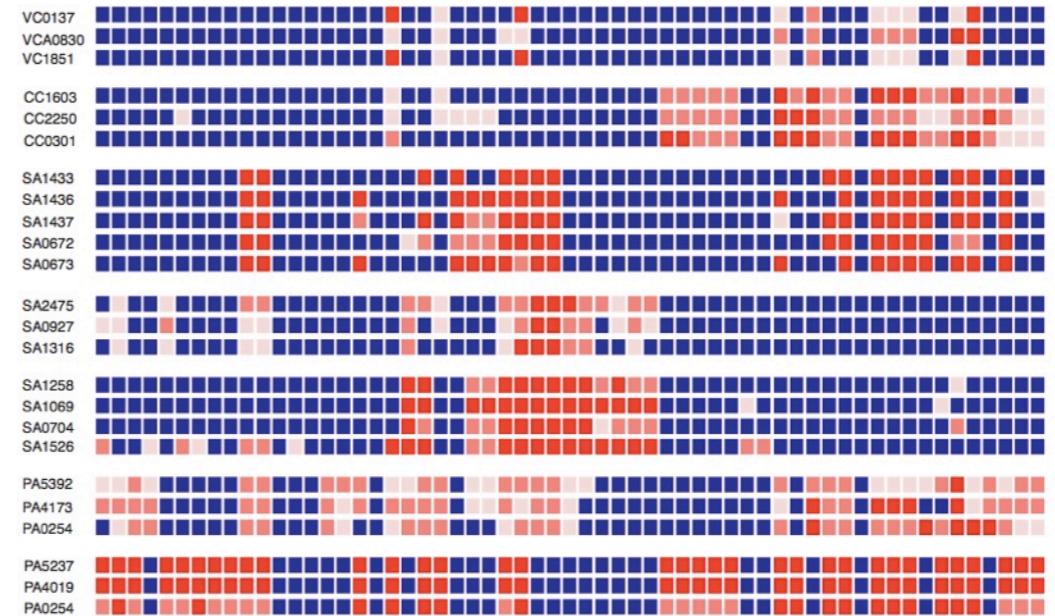
	A
Genome 1	1
Genome 2	1
Genome 3	1
Genome 4	1
Genome 5	0
Genome 6	0

Pagel, P., Wong,
phylogenetic p

Sun, J., Li, Y. &
protein intera
Res Commun

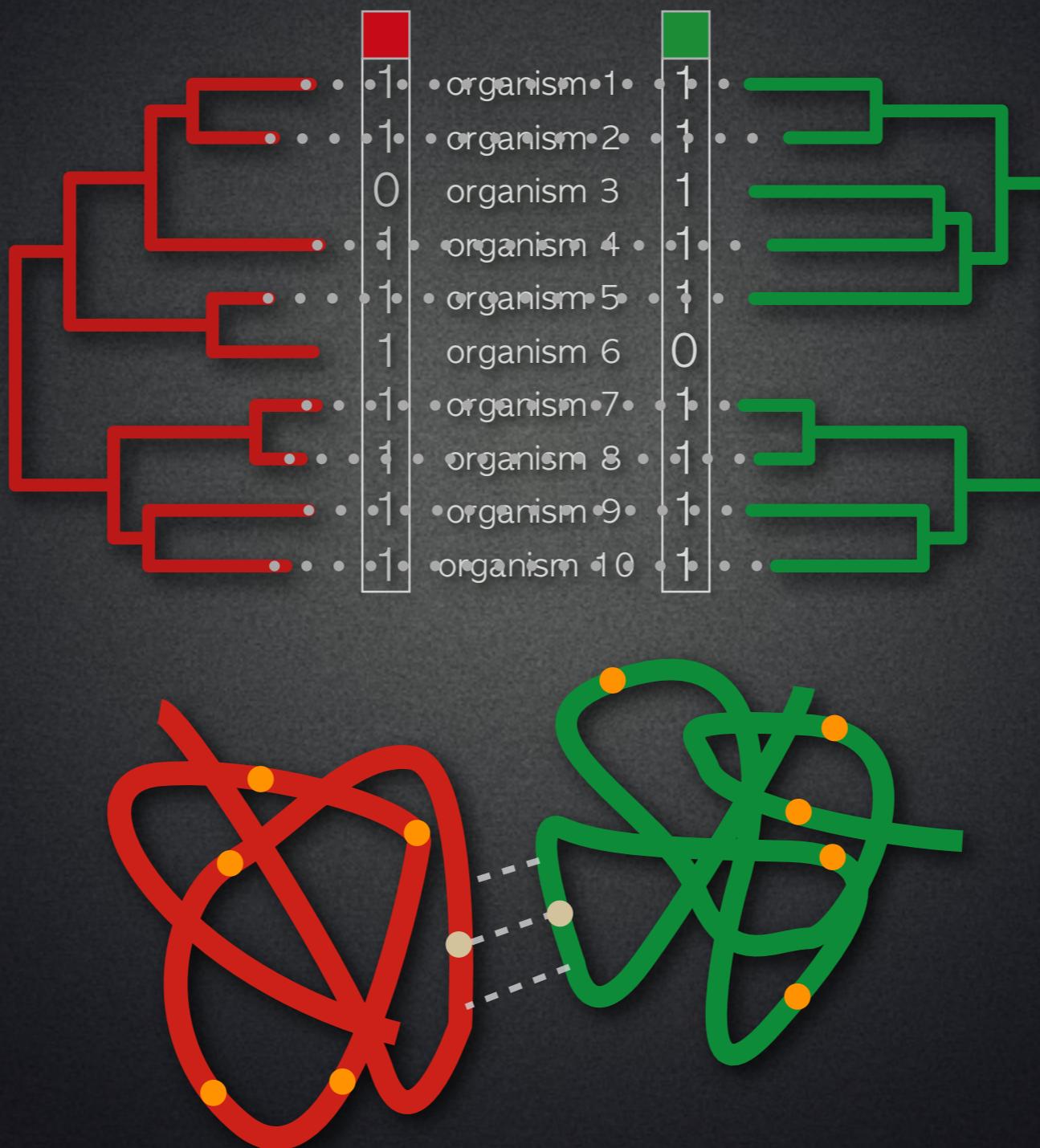


Quantitative Profiles



Date, S. V. & Marcotte, E. M. Discovery of uncharacterized cellular systems by genome-wide analysis of functional linkages. Nat Biotechnol 21, 1055–1062 (2003)

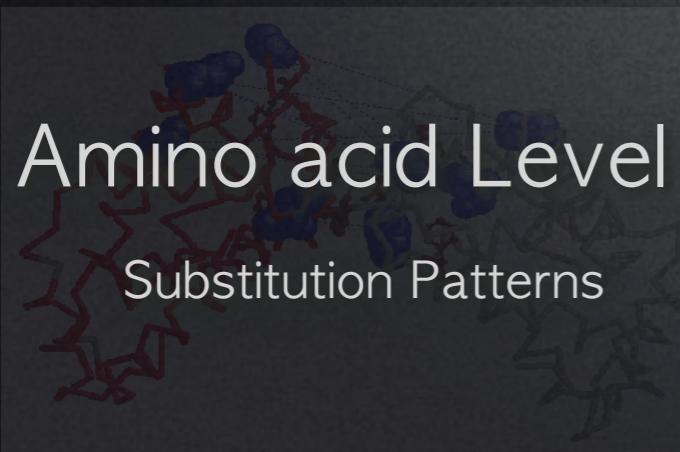
Similarity of Phylogenetic Trees



Fryxell, K. J. The coevolution of gene family trees. Trends Genet 12, 364–369 (1996).

Co-evolution

“reciprocal evolutionary change in interacting species” Thompson 1994

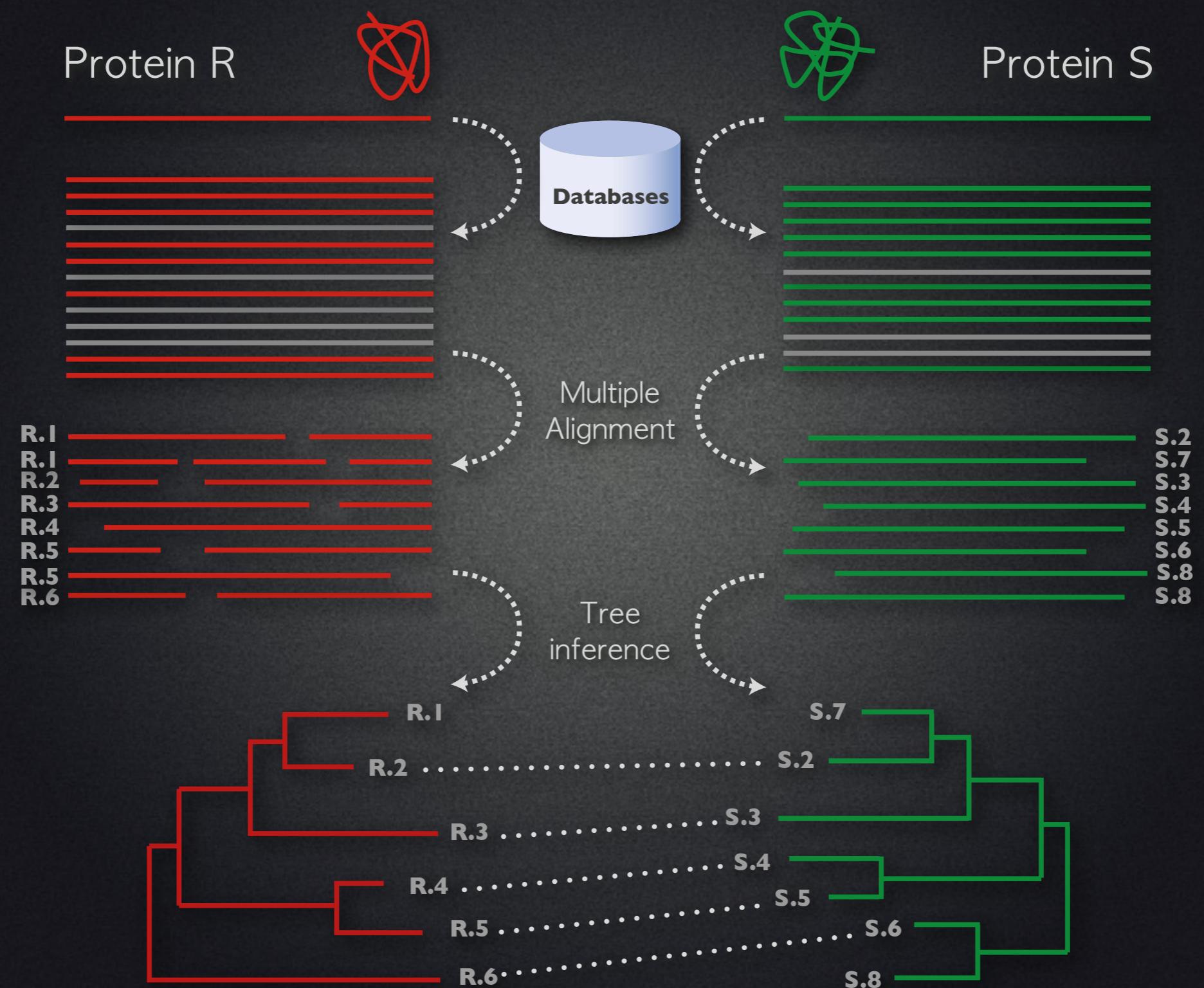


Ecological Factors
Trophic networks
Mutualisms
Symbioses
...

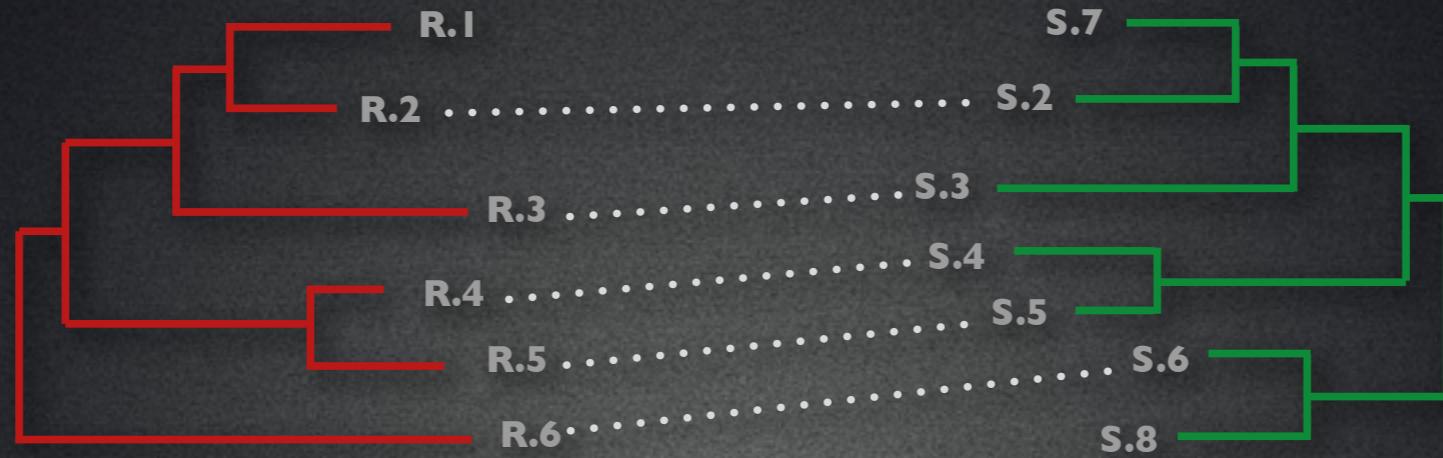
Protein Structure
Protein Function
Protein Networks
...

Local Structural Effects
Amino acid Contact Networks
Functional Sites
Protein-Protein Interfaces
...

MirrorTree

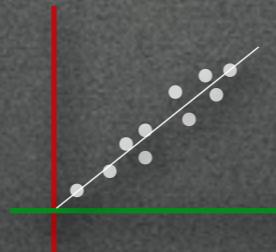


MirrorTree



Distances matrix

	R1	R2	R3	R4	R5	R6
R1						
R2						
R3						
R4						
R5						
R6						



Pearson Correlation

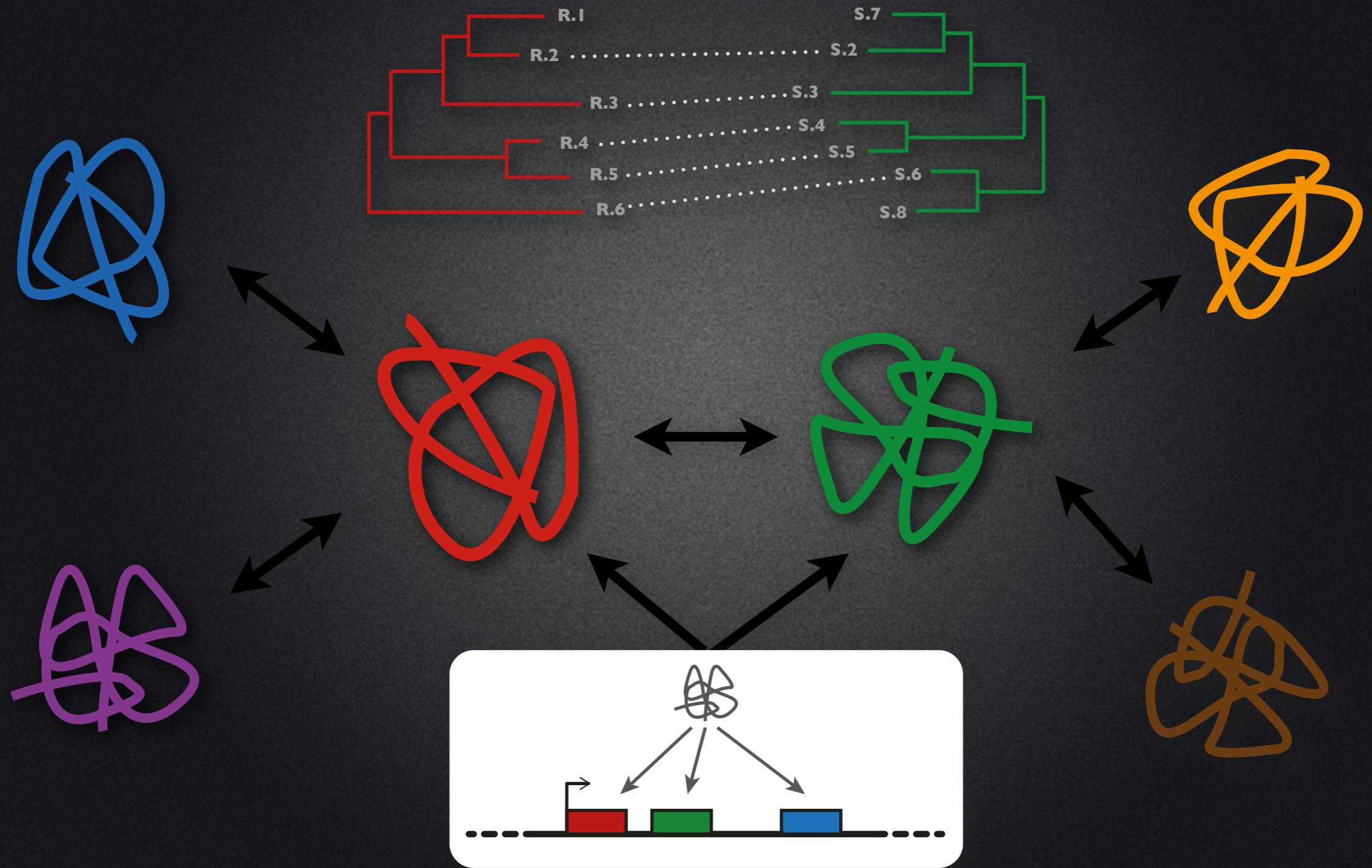
$$r = \frac{\sum_{i=1}^n (R_i - \bar{R})(S_i - \bar{S})}{\sqrt{\sum_{i=1}^n (R_i - \bar{R})^2} \sqrt{\sum_{i=1}^n (S_i - \bar{S})^2}}$$

Distances matrix

	S2	S3	S4	S5	S6	S7	S8
S2							
S3							
S4							
S5							
S6							
S7							
S8							

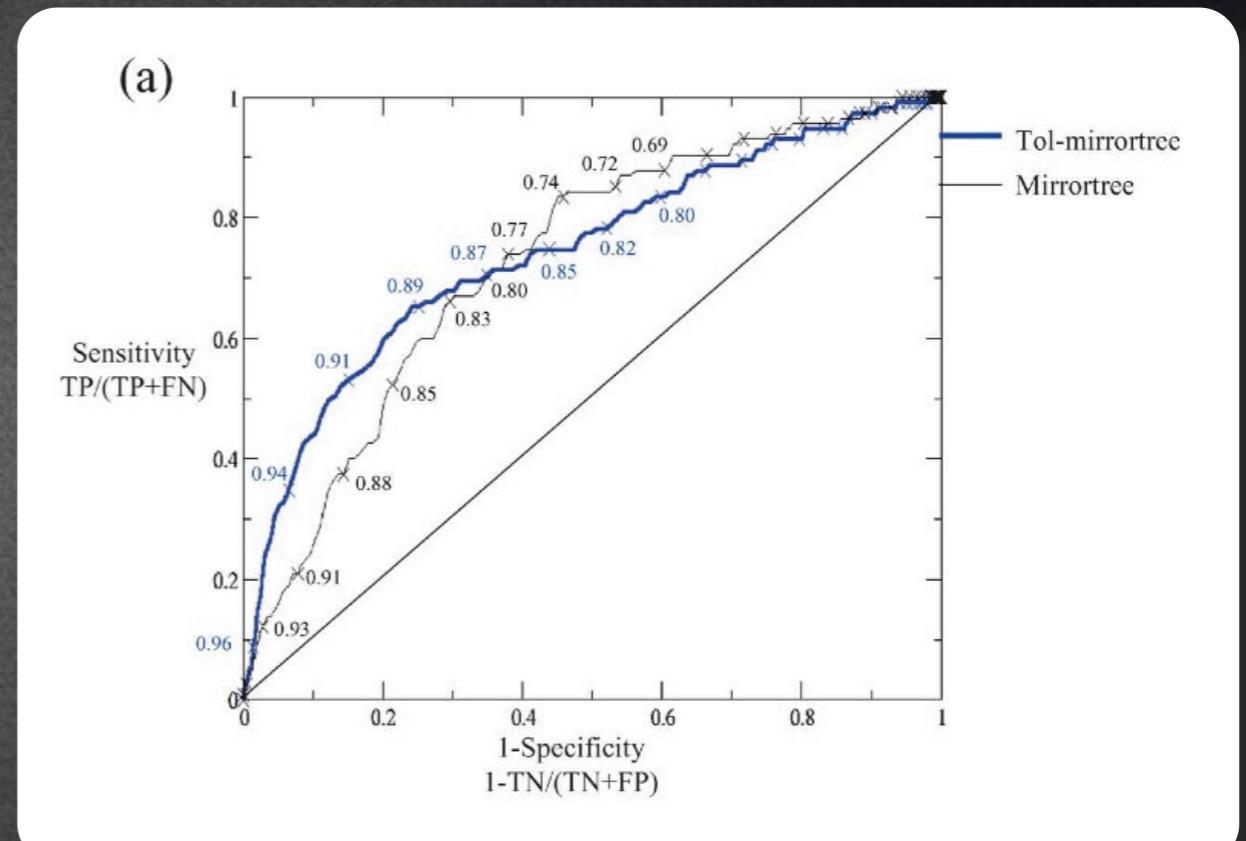
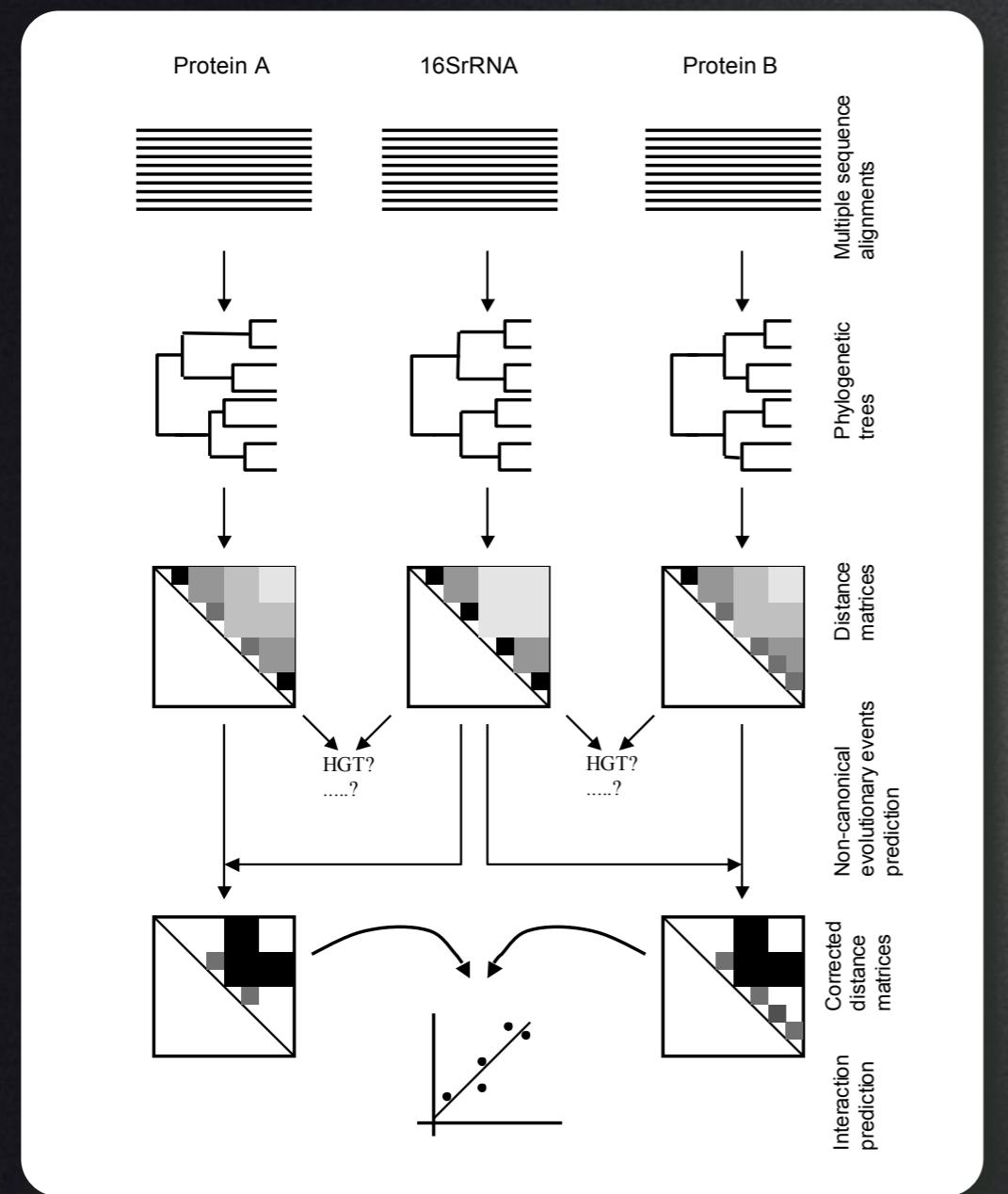
Pazos, F. & Valencia, A. Similarity of phylogenetic trees as indicator of protein-protein interaction. *Protein Eng* 14, 609–614 (2001)

Possible causes of observed similarity



Pazos, F. & Valencia, A. Protein co-evolution, co-adaptation and interactions. EMBO J 27, 2648–2655 (2008).

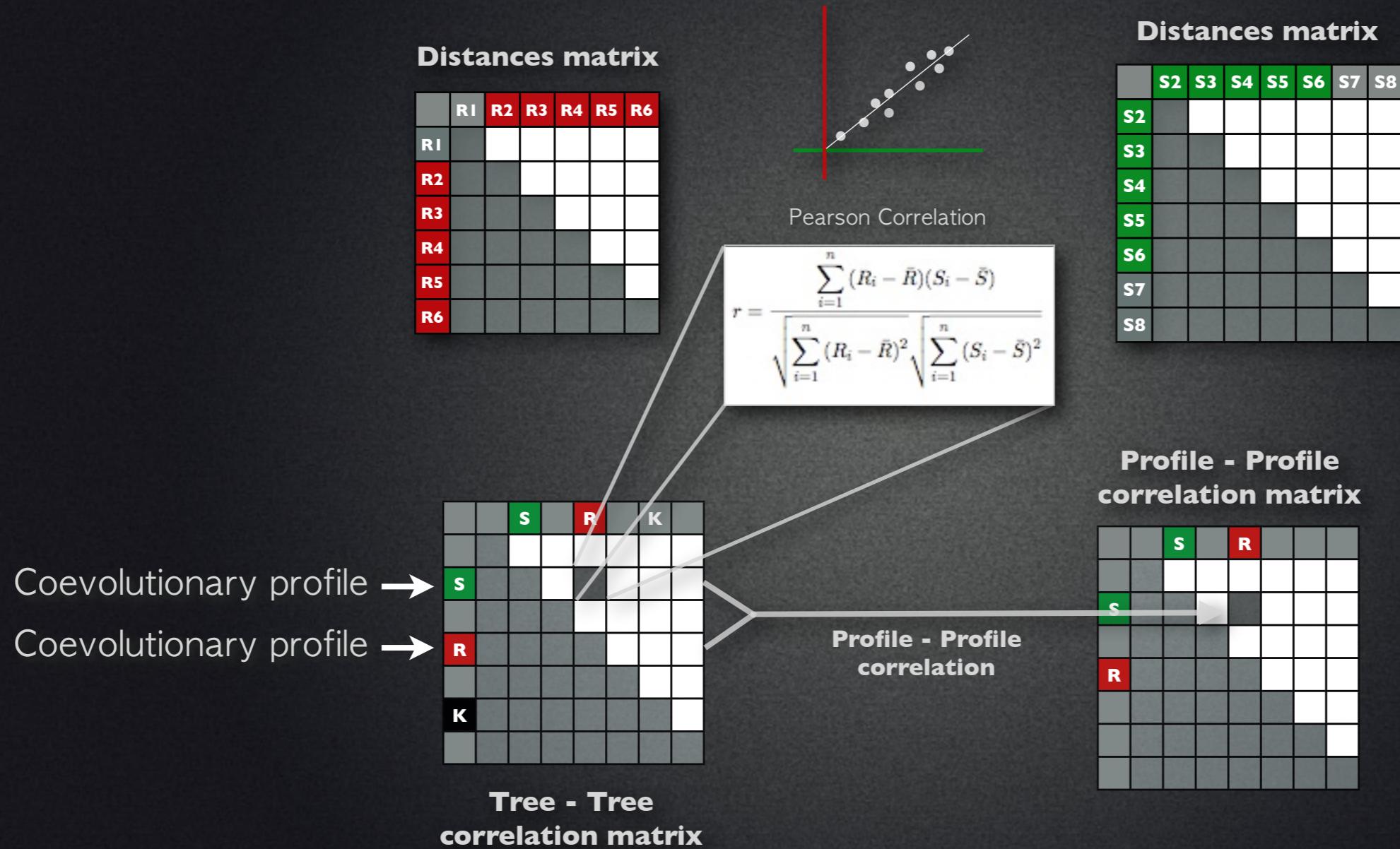
TOL-MirrorTree



Sato, T., Yamanishi, Y., Kanehisa, M. & Toh, H. The inference of protein-protein interactions by co-evolutionary analysis is improved by excluding the information about the phylogenetic relationships. *Bioinformatics* 21, 3482–3489 (2005).

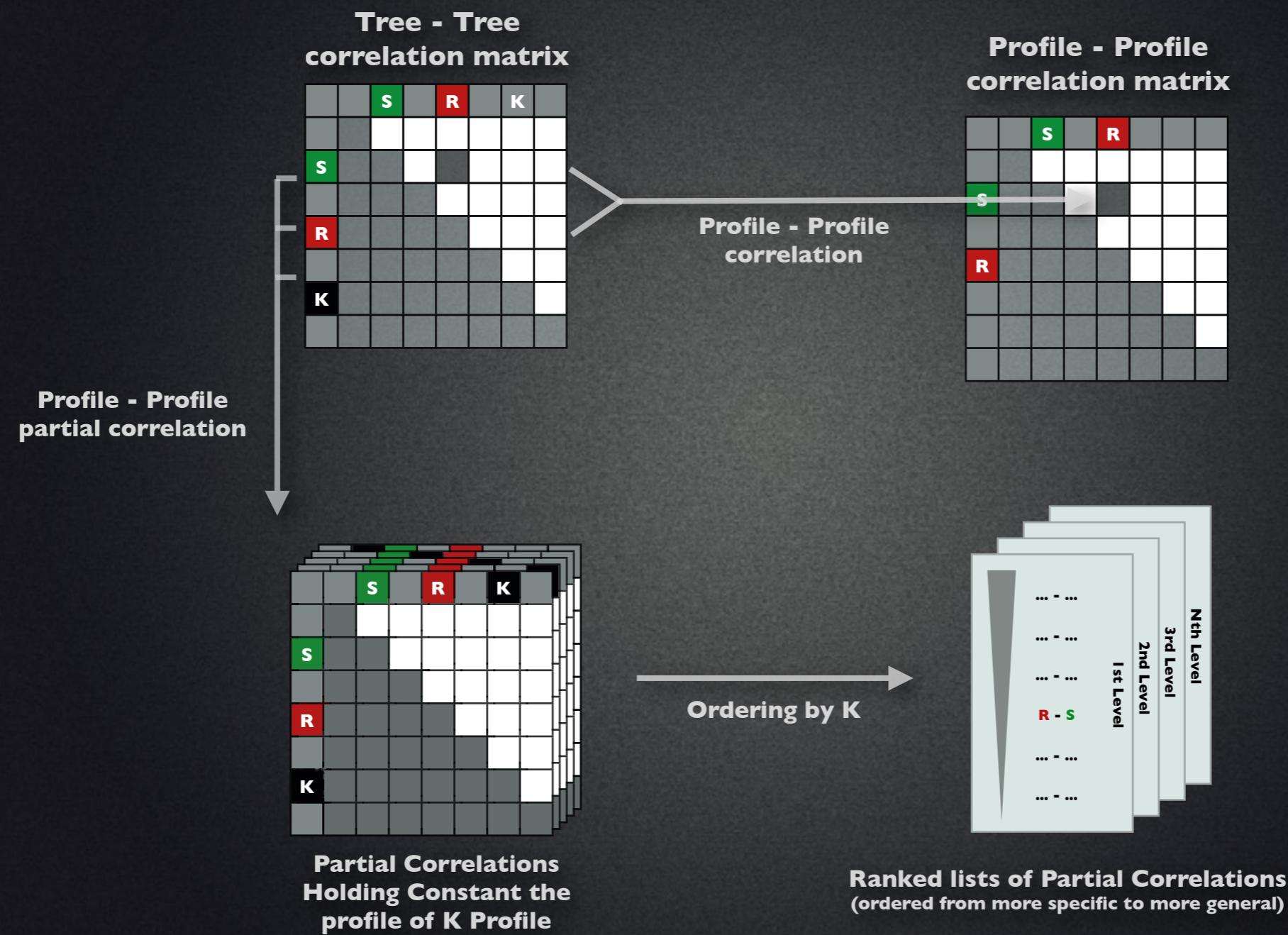
Pazos, F., Ranea, J. A. G., Juan, D. & Sternberg, M. J. E. Assessing protein co-evolution in the context of the tree of life assists in the prediction of the interactome. *J Mol Biol* 352, 1002–1015 (2005).

Profile Correlation



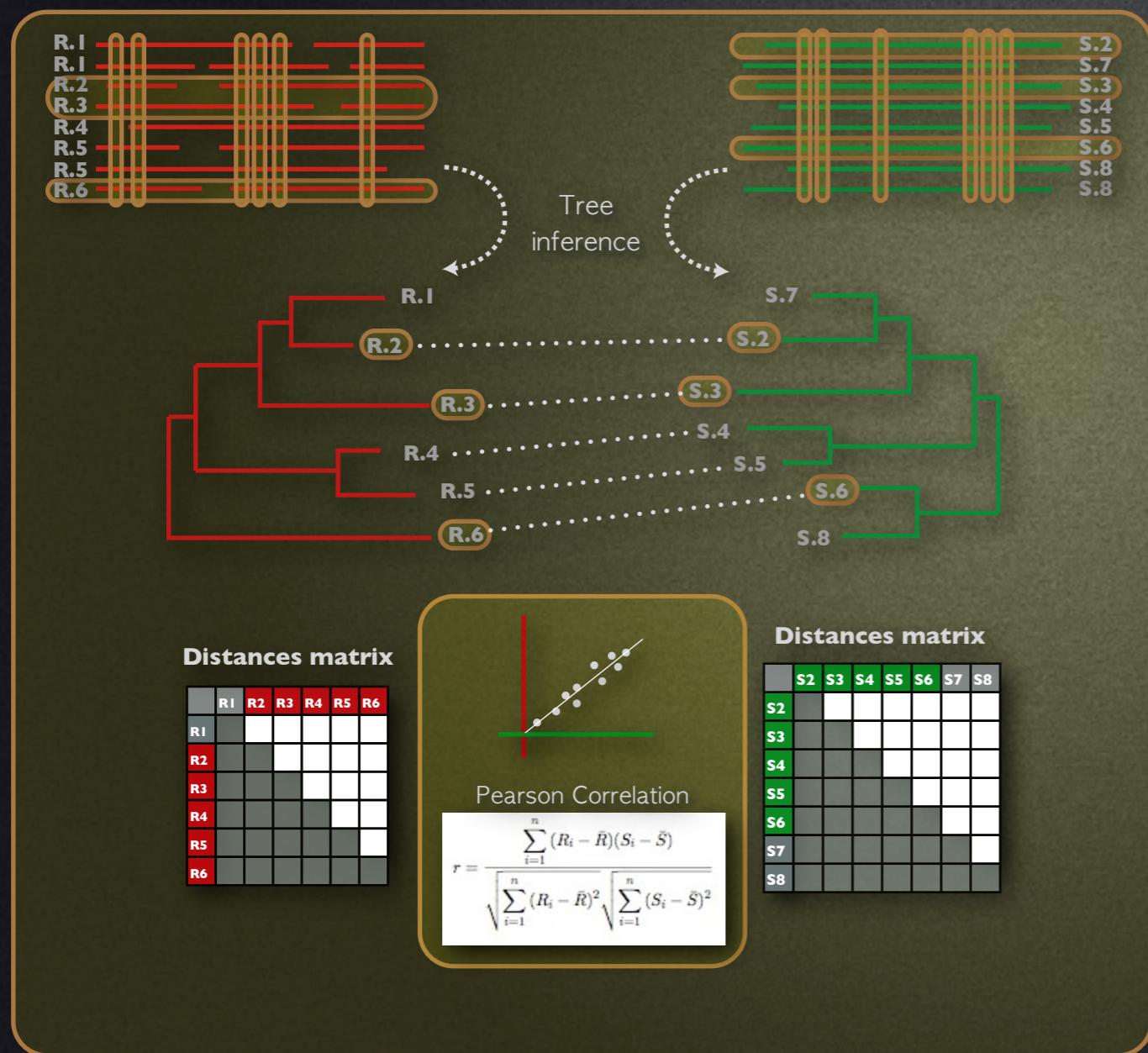
Juan, D., Pazos, F. & Valencia, A. High-confidence prediction of global interactomes based on genome-wide coevolutionary networks. *Proc Natl Acad Sci U S A* **105**, 934–939 (2008)

ContextMirror

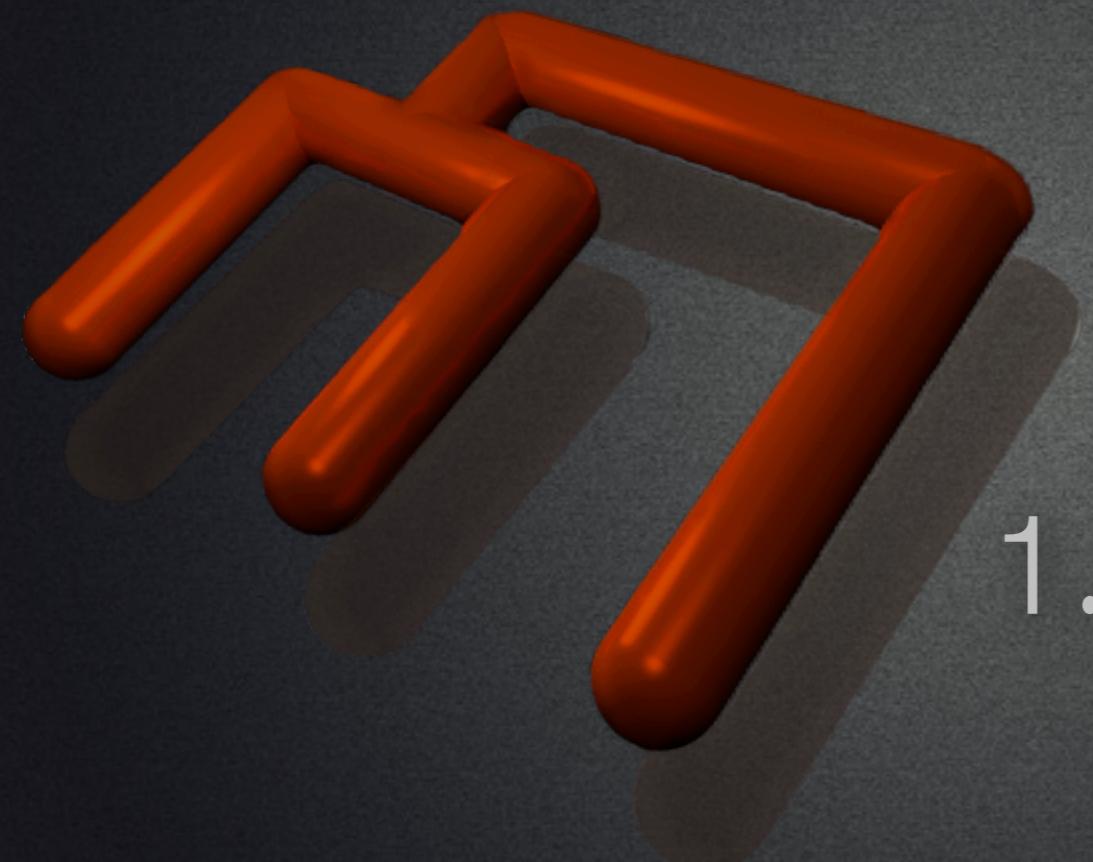


Juan, D., Pazos, F. & Valencia, A. High-confidence prediction of global interactomes based on genome-wide coevolutionary networks. *Proc Natl Acad Sci U S A* **105**, 934–939 (2008)

Objectives



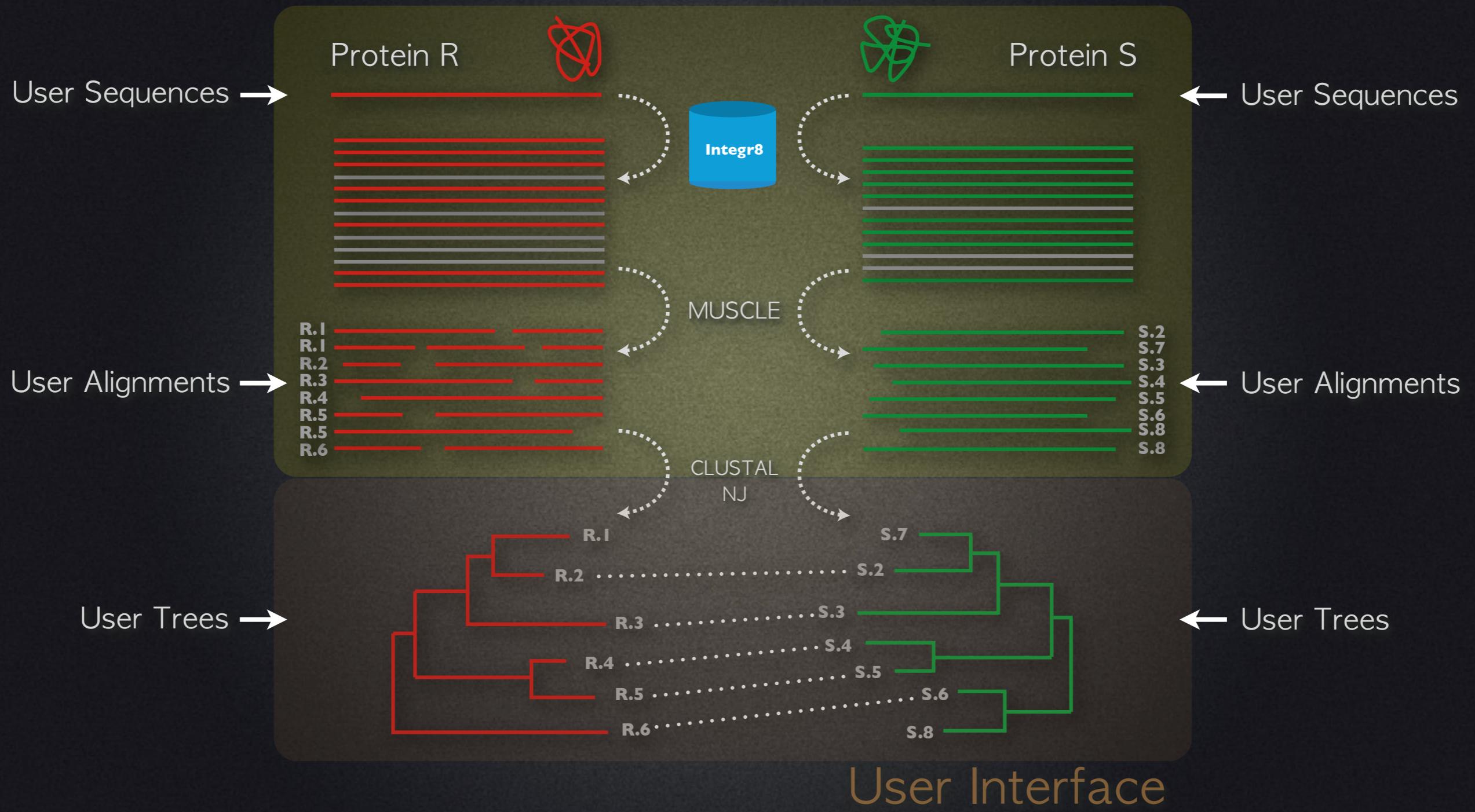
1. MirrorTree Server
2. Incorporating predicted accessibility
3. Selection of organisms
4. Co-evolution significance



1. MirrorTree Server

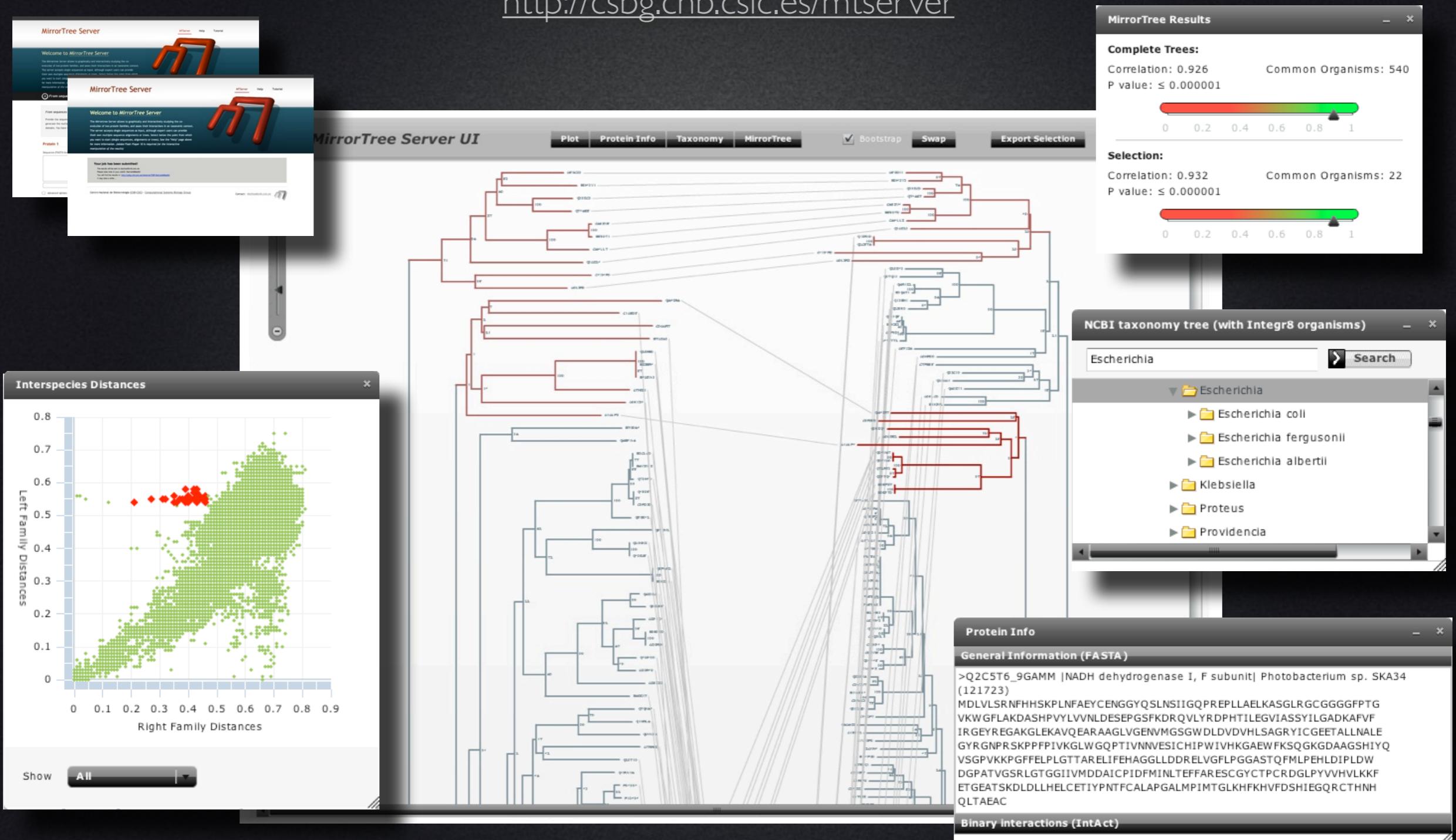
Ochoa, D. & Pazos, F. Studying the co-evolution of protein families with the Mirrortree web server.
Bioinformatics **26**, 1370–1371 (2010)

Automatic Pipeline

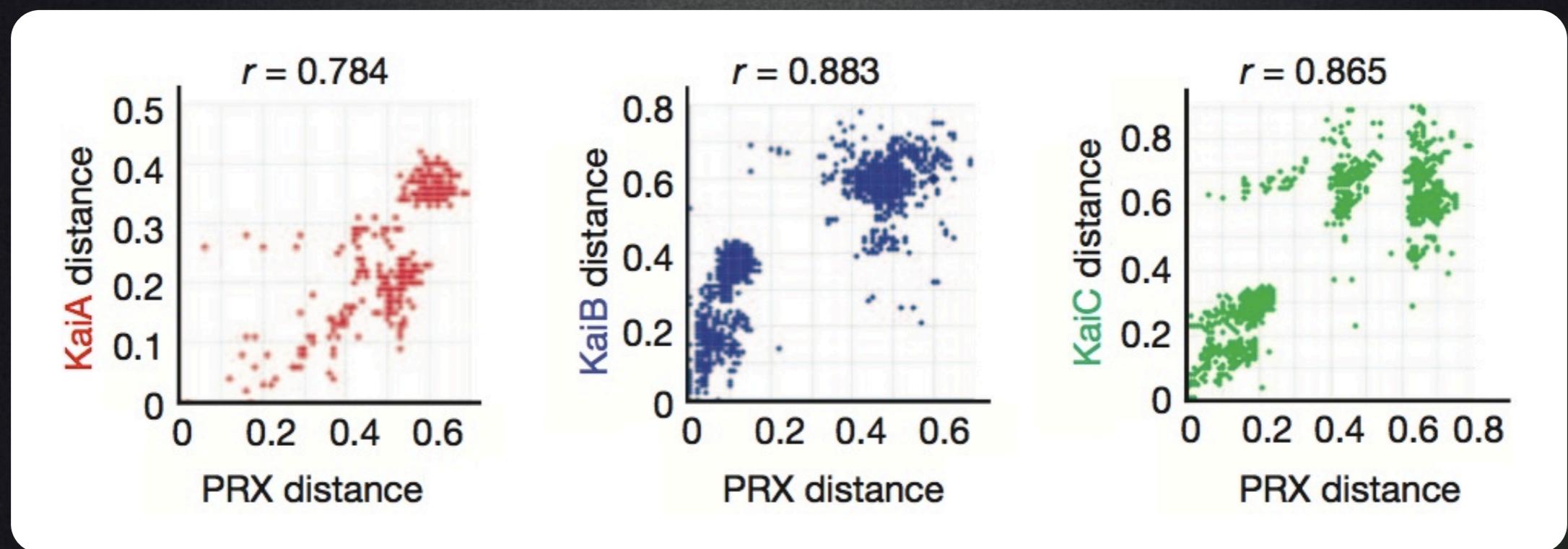


MirrorTree Server

<http://csbg.cnb.csic.es/mtserver>



Peroxiredoxins



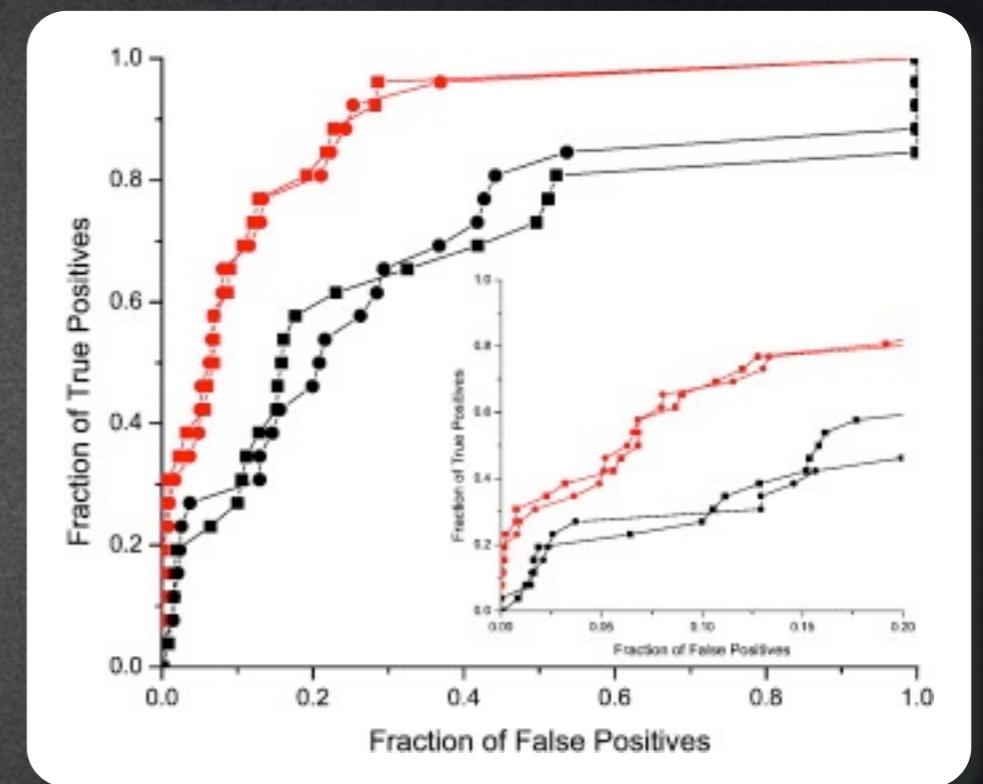
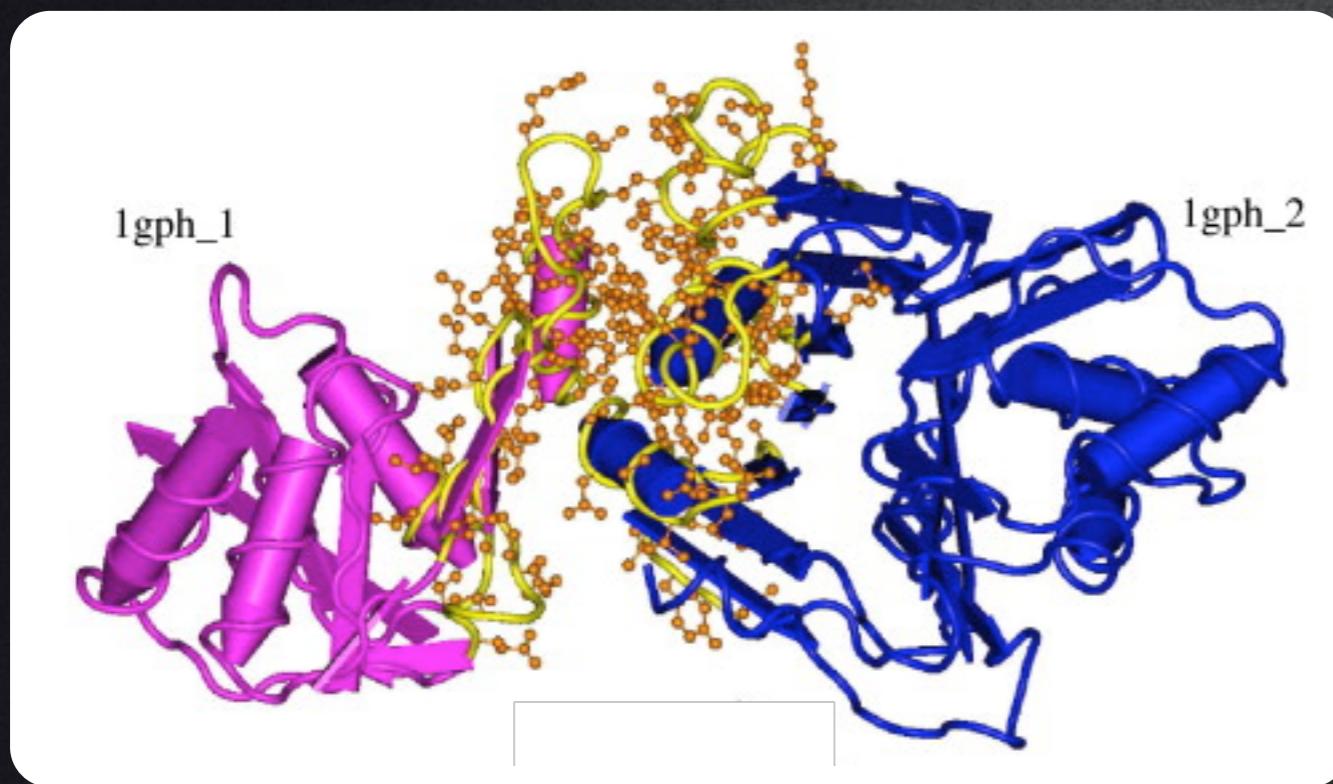
Edgar, R. S. et al. Peroxiredoxins are conserved markers of circadian rhythms. Nature 485, 459–464 (2012).



2. Incorporating Predicted Accessibility

Ochoa, D., García-Gutiérrez, P., Juan, D., Valencia, A. & Pazos, F. Incorporating information on predicted solvent accessibility to the co-evolution-based study of protein interactions. *Molecular bioSystems* 9, 70–76 (2013).

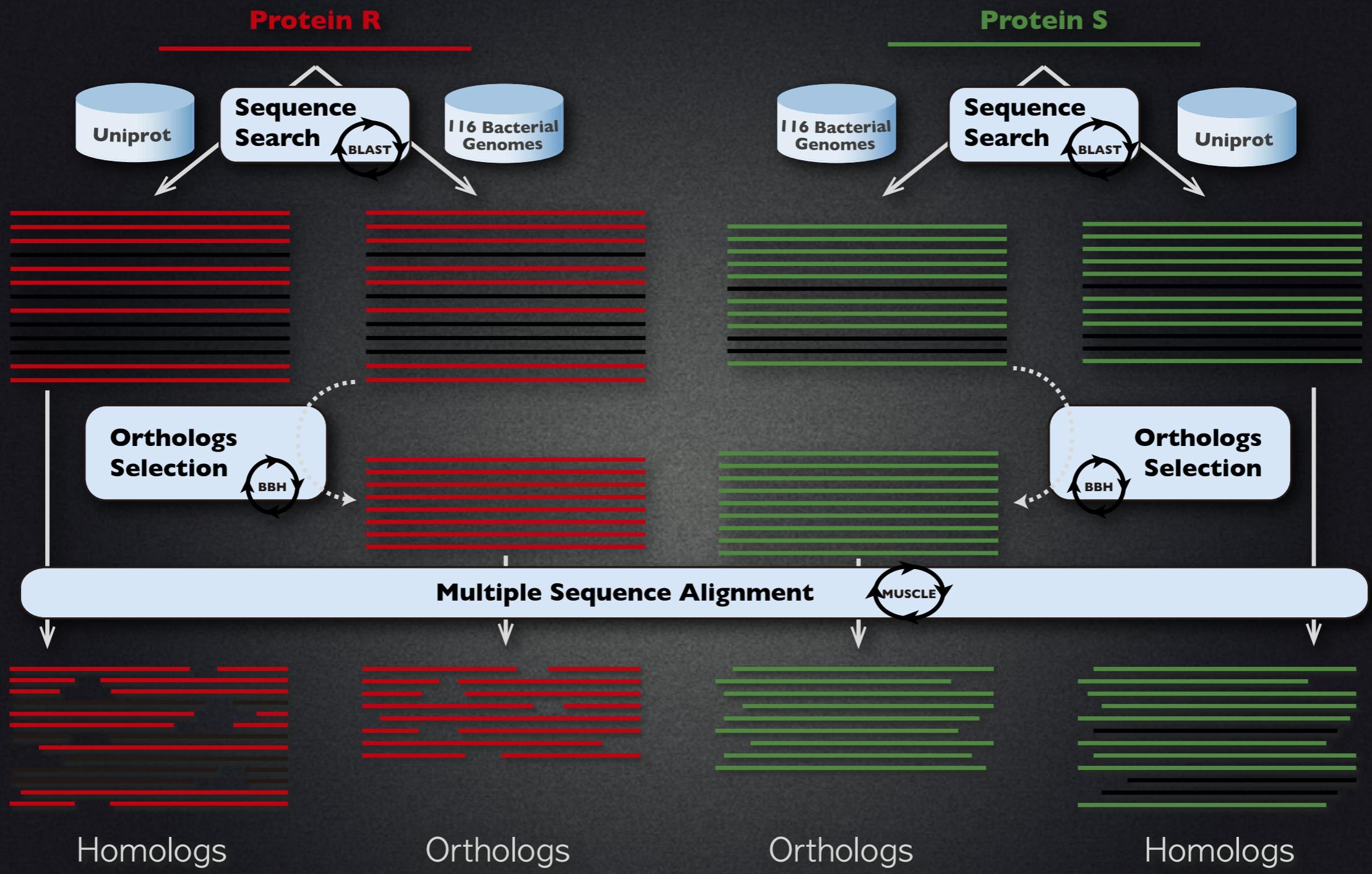
Residue filtering

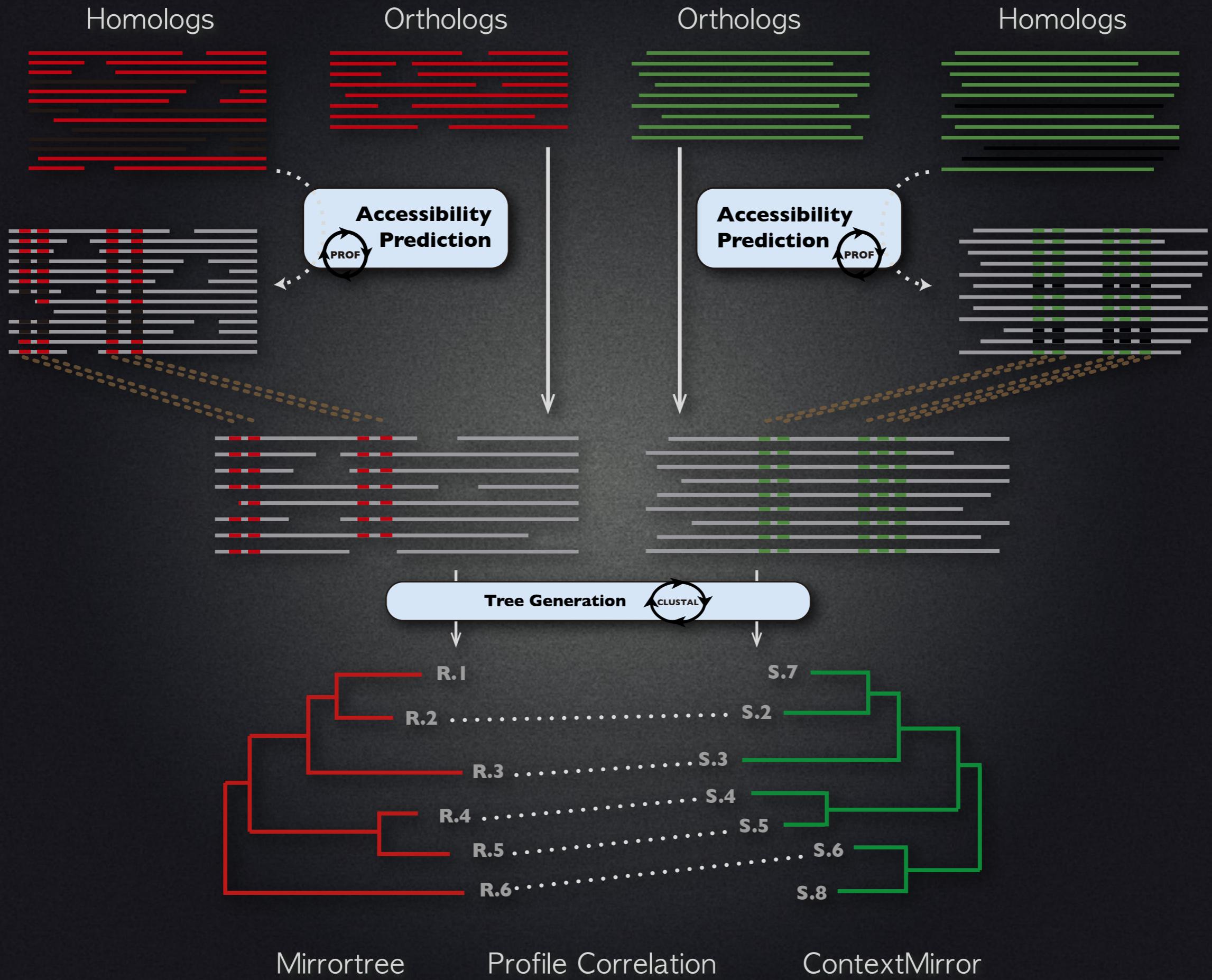


Kann, M. G., Shoemaker, B. A., Panchenko, A. R. & Przytycka, T. M. Correlated evolution of interacting proteins: looking behind the mirrortree. *J Mol Biol* 385, 91–98 (2009).

Objective

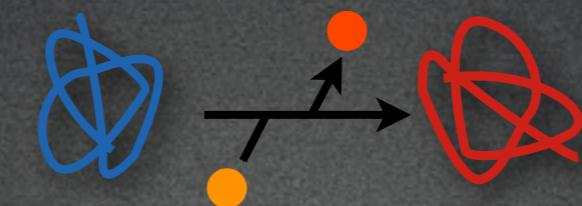
Asses the effect of including predicted solvent accessibility on the performance of the PPI inference based on similarity of phylogenetic trees





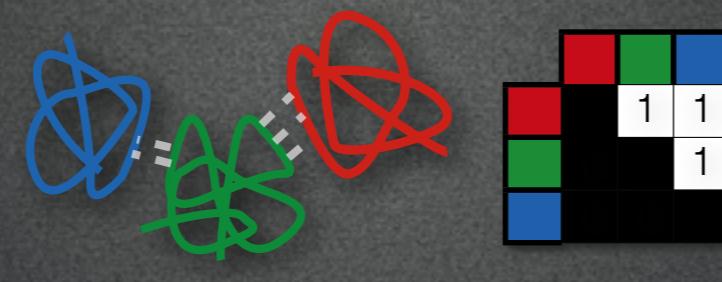
Gold Standards

Pathways
EcoCyc



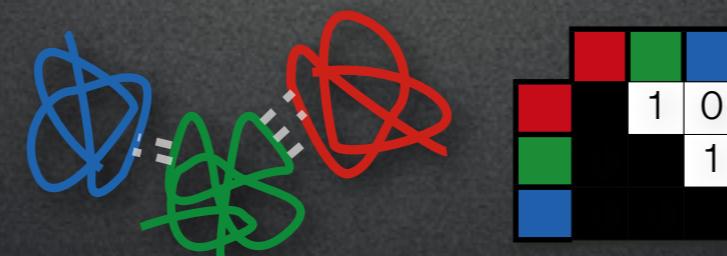
719 proteins
4,491 evidences

Complexes
EcoCyc



591 proteins
1,354 evidences

Binary Physical
MPIDB



1,268 proteins
1,626 evidences

Keseler, I. M. et al. EcoCyc: a comprehensive database of *Escherichia coli* biology. Nucleic Acids Res 39, D583–90 (2011).

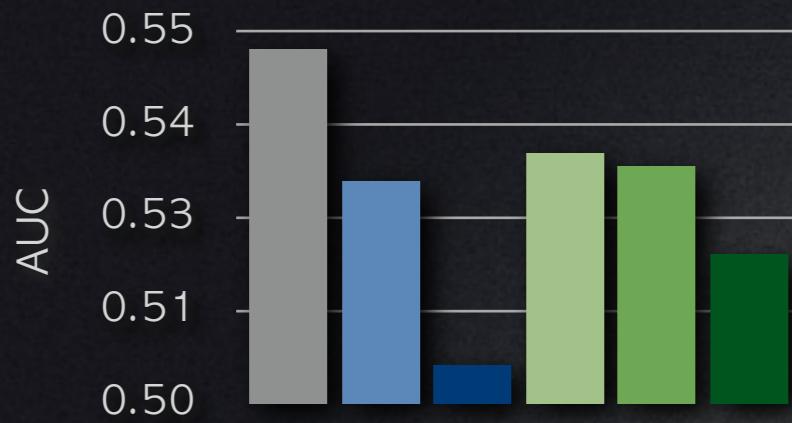
Goll, J. et al. MPIDB: the microbial protein interaction database. Bioinformatics 24, 1743–1744 (2008).

Binary Physical

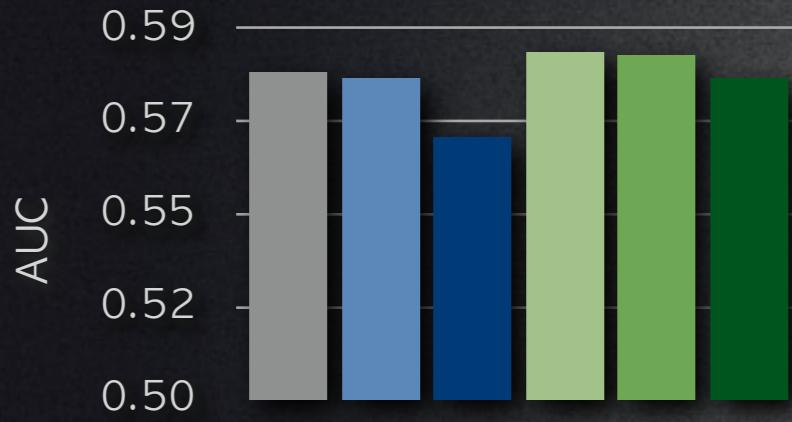
Complexes

Pathways

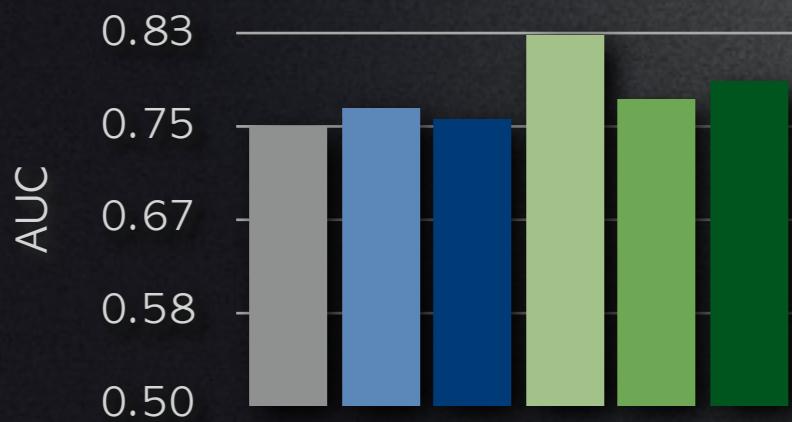
mirrOrtree



Profile Correlation



ContextMirror
(level 1)



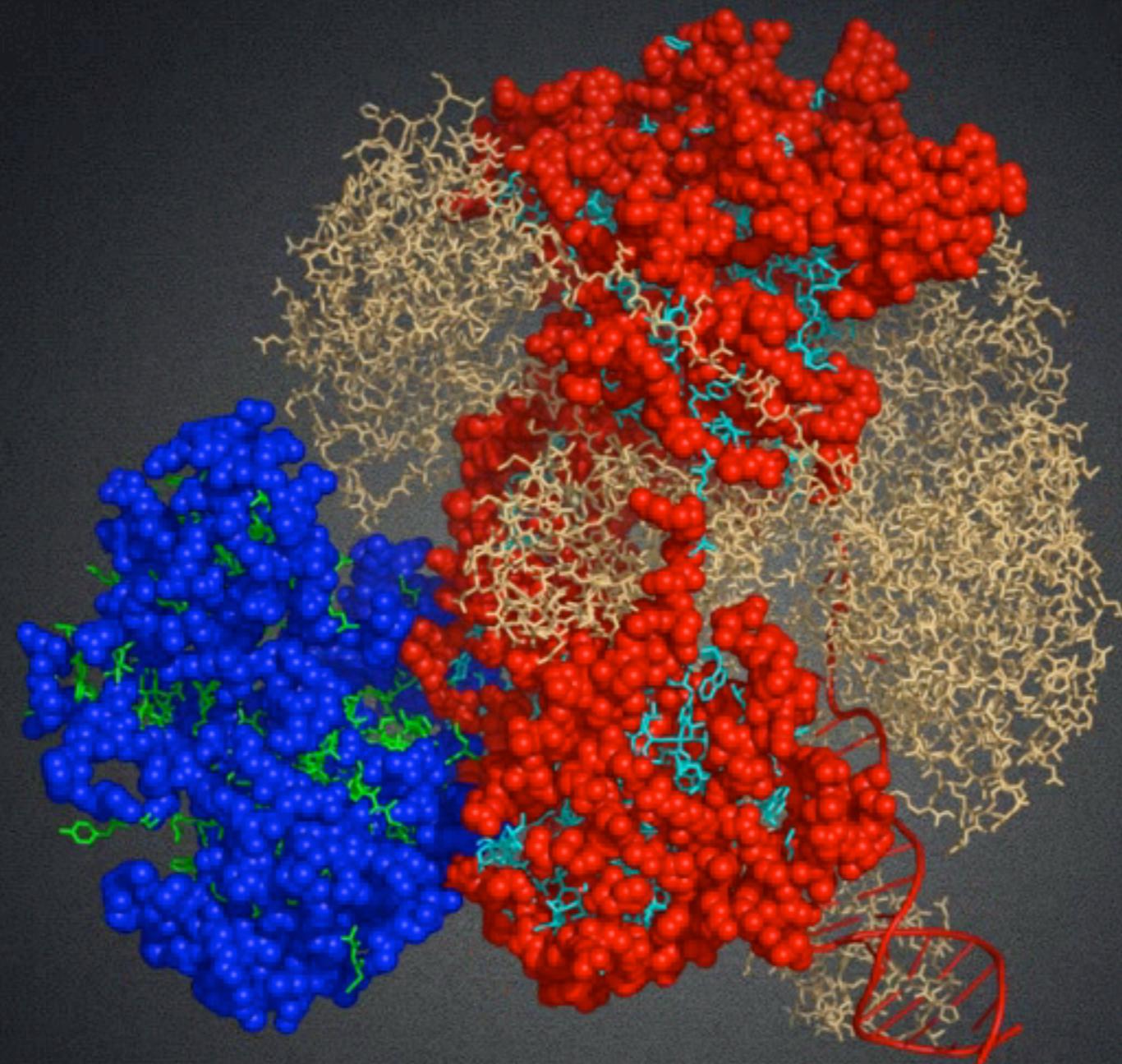
Legend: all (grey), eRIA0 (blue), eRIA3 (dark blue), pACC2 (light green), pACC12 (medium green), pACC50 (dark green)

Ochoa, D., García-Gutiérrez, P., Juan, D., Valencia, A. & Pazos, F. Incorporating information on predicted solvent accessibility to the co-evolution-based study of protein interactions. Molecular bioSystems 9, 70–76 (2013).

 RecB
 RecC
 RecD

 pACC2





ATP-dependent helicase/nuclease complex

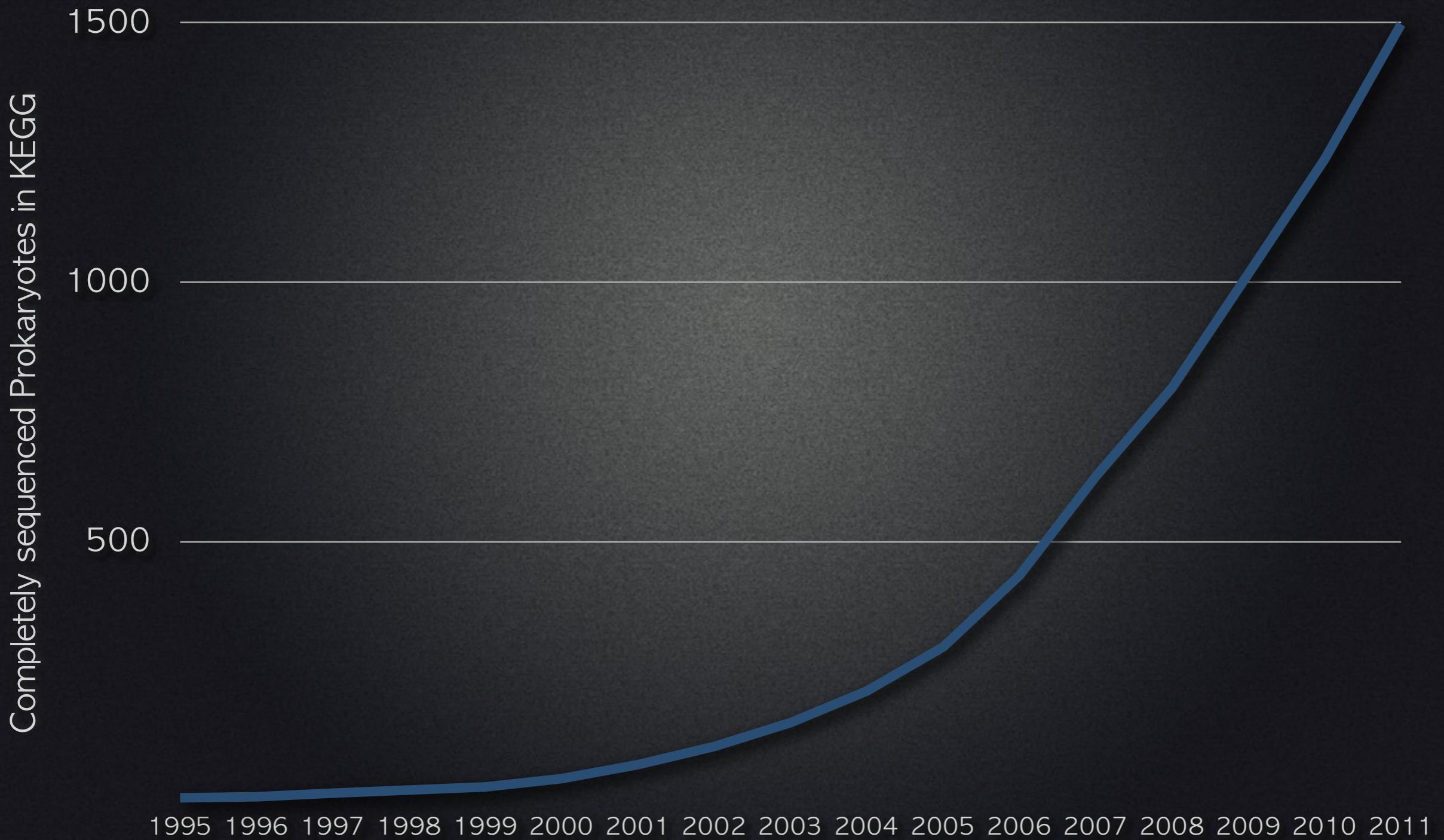
	ALL	eRIA0	eRIA3	pACC2	pACC12	pACC50
✓ RecC-RecD	0.701	0.768	Not significant	0.739	0.750	0.806
✗ RecB-RecD	0.427	Not significant				



3. Selection of organisms

Herman, Dorota, David Ochoa, David Juan, Daniel Lopez, Alfonso Valencia, and Florencio Pazos. 2011. ‘Selection of organisms for the co-evolution-based study of protein interactions.’, BMC Bioinformatics, 12, 363

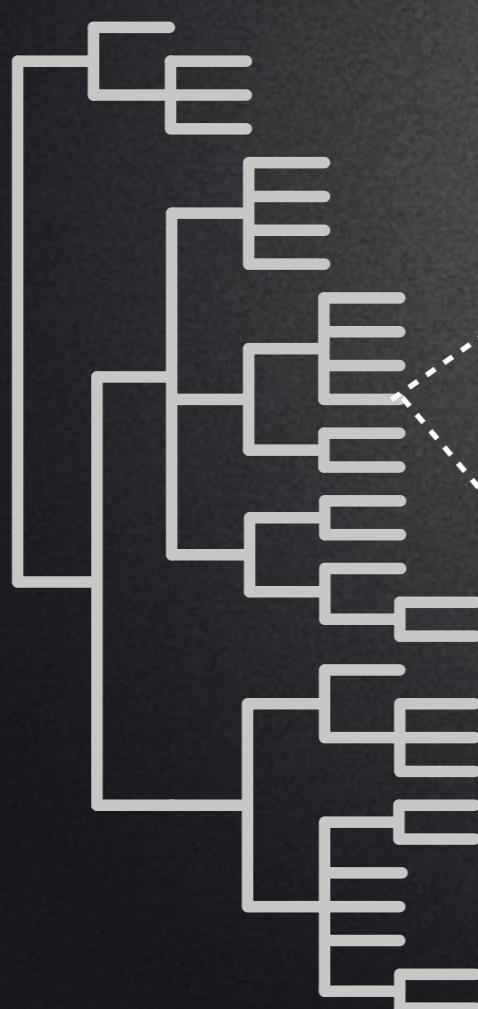
Sequence explosion



The redundancy problem

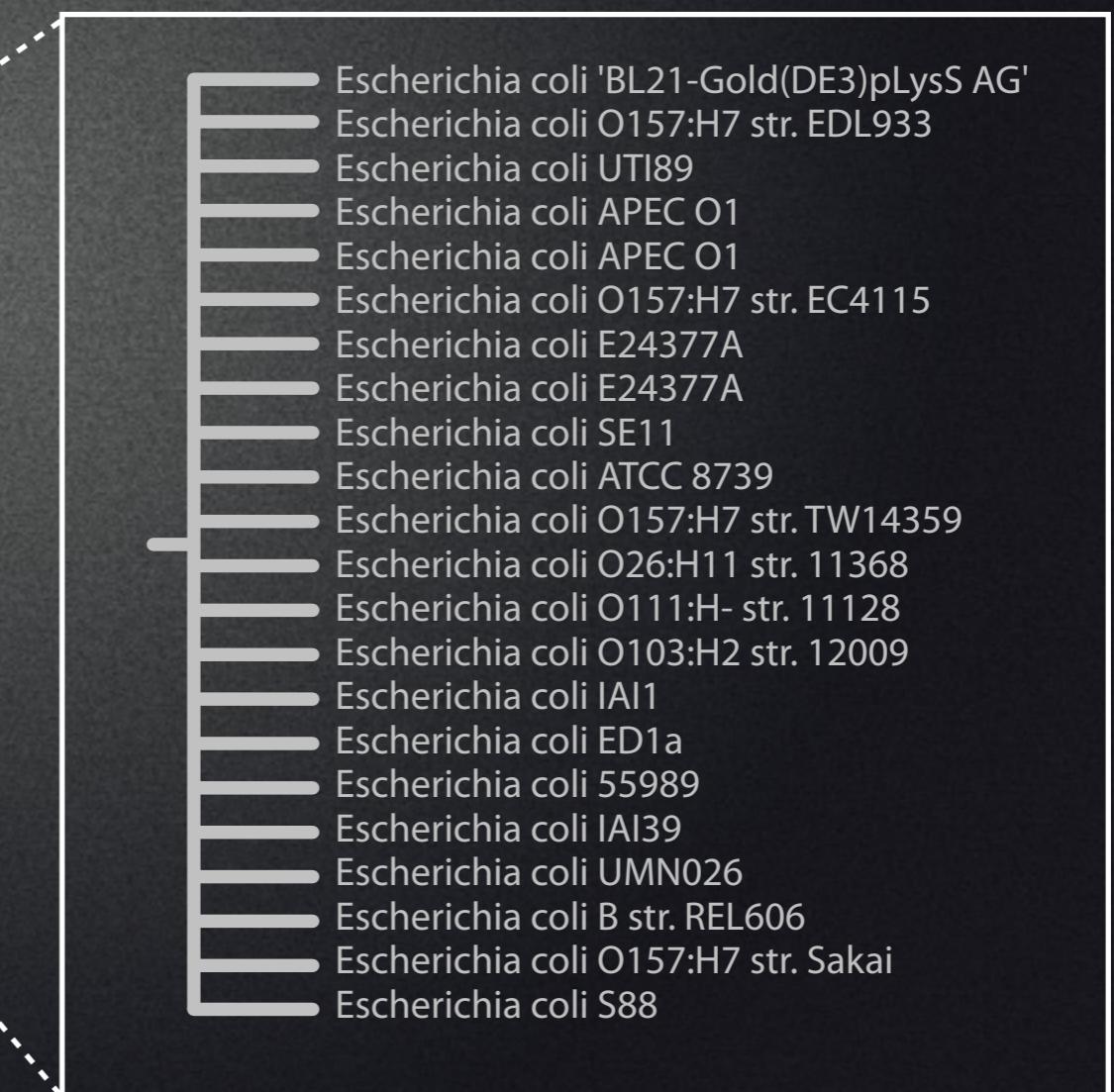
2001

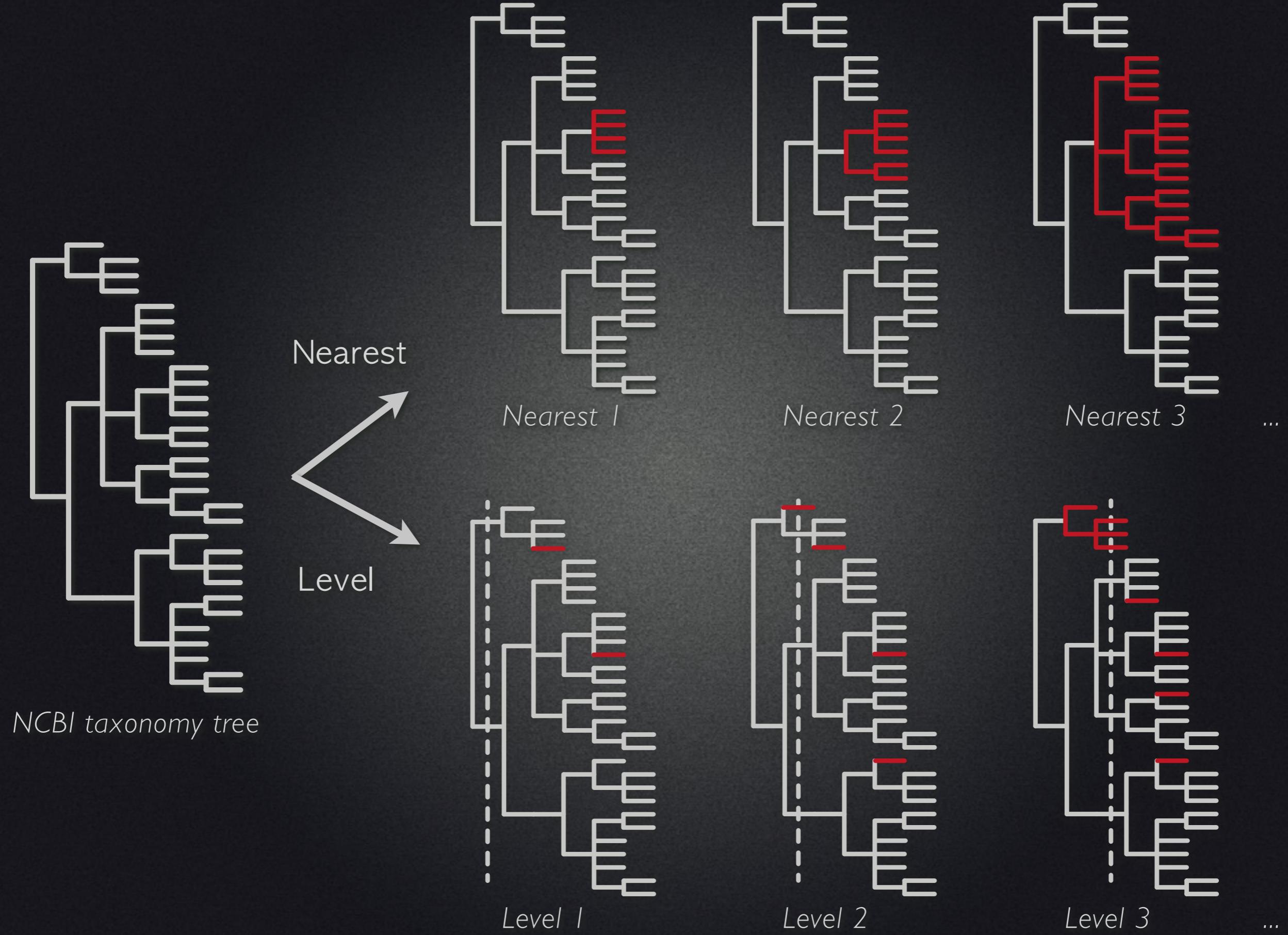
70 sequenced prokaryotes

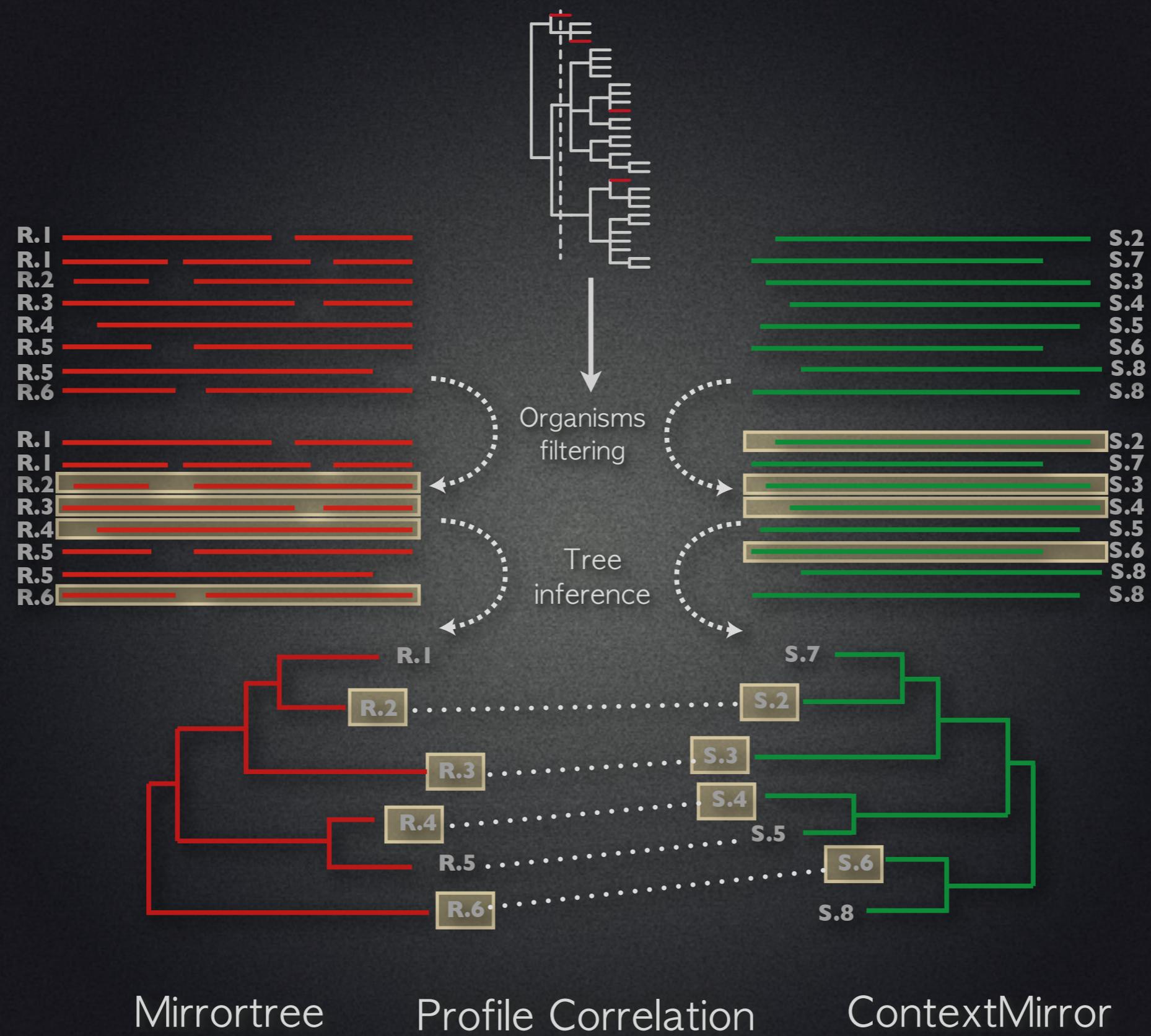


2011

1,498 sequenced prokaryotes







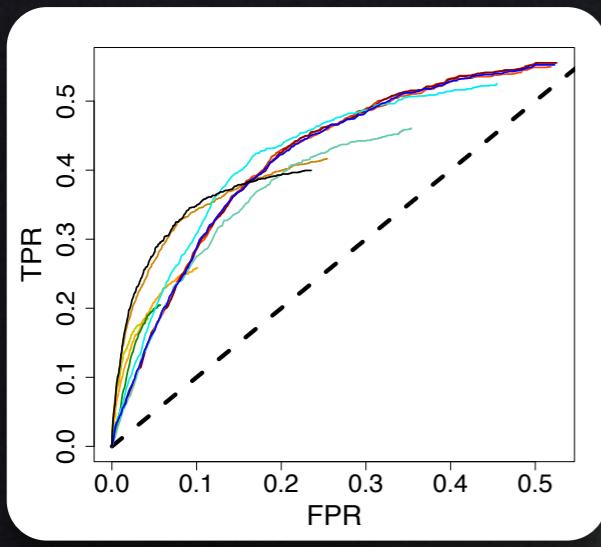
Mirrortree

Profile Correlation

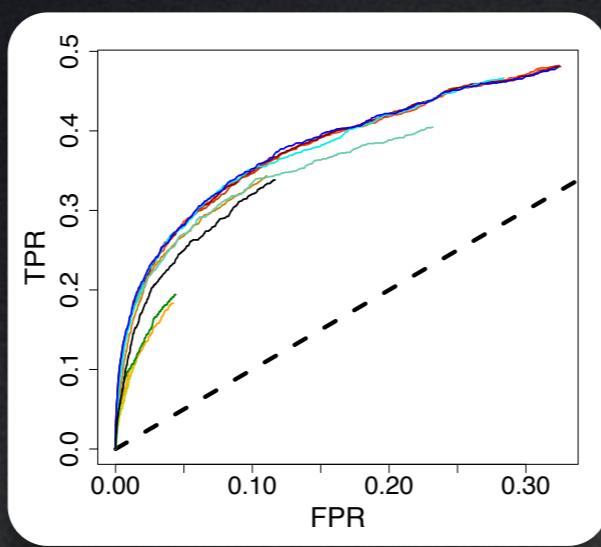
ContextMirror

Complexes

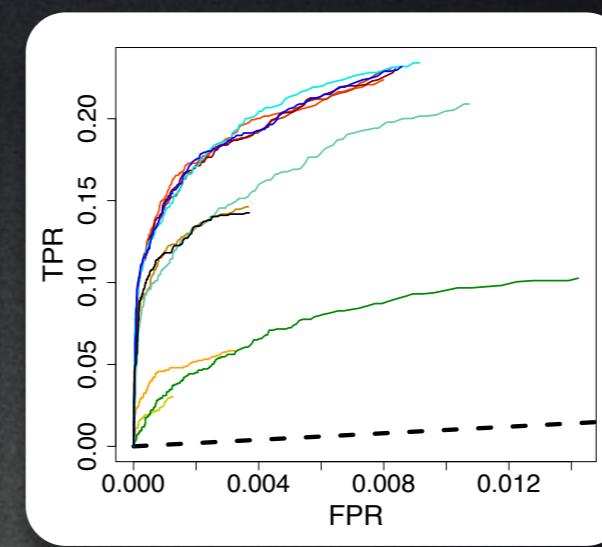
Mirrortree



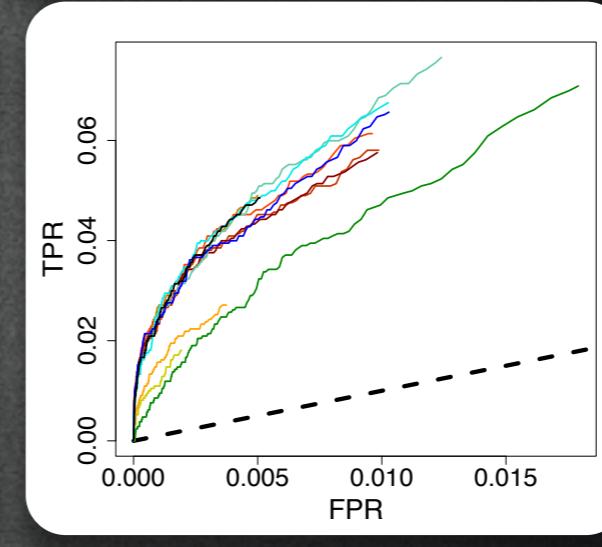
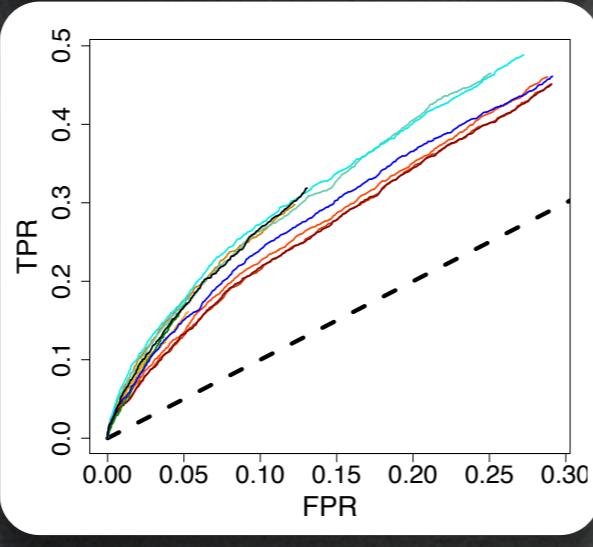
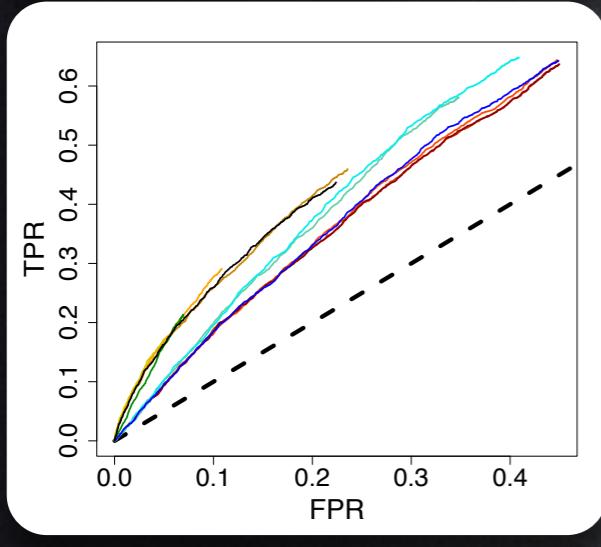
Profile Correlation



ContextMirror

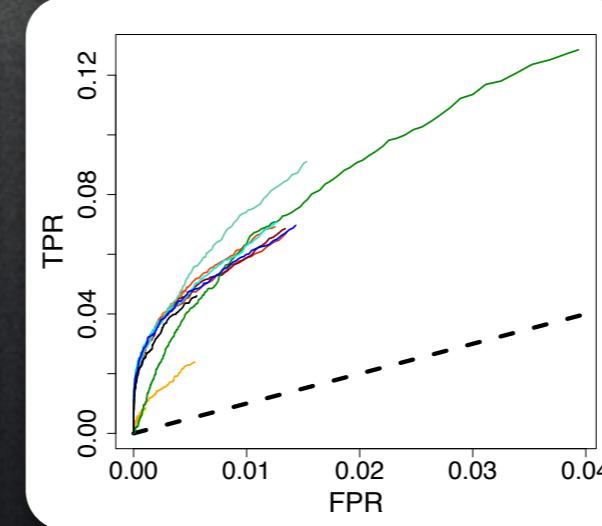
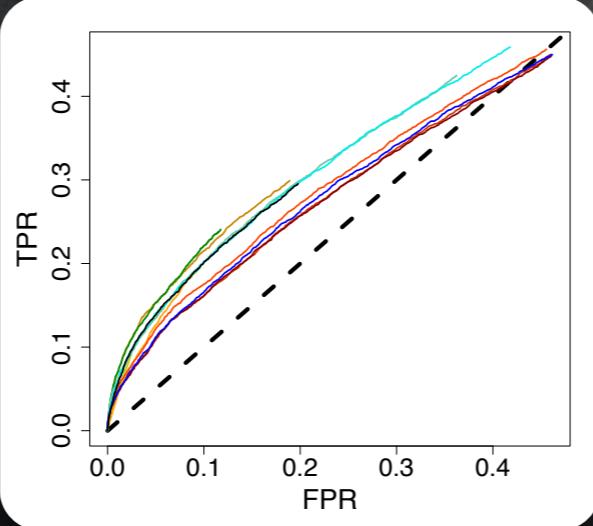
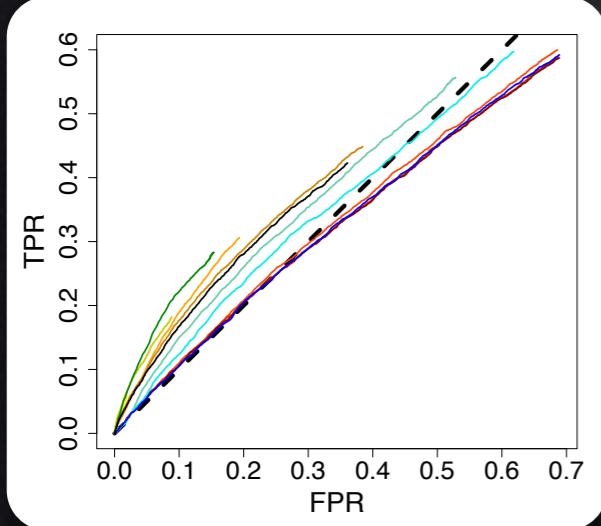


Binary Physical

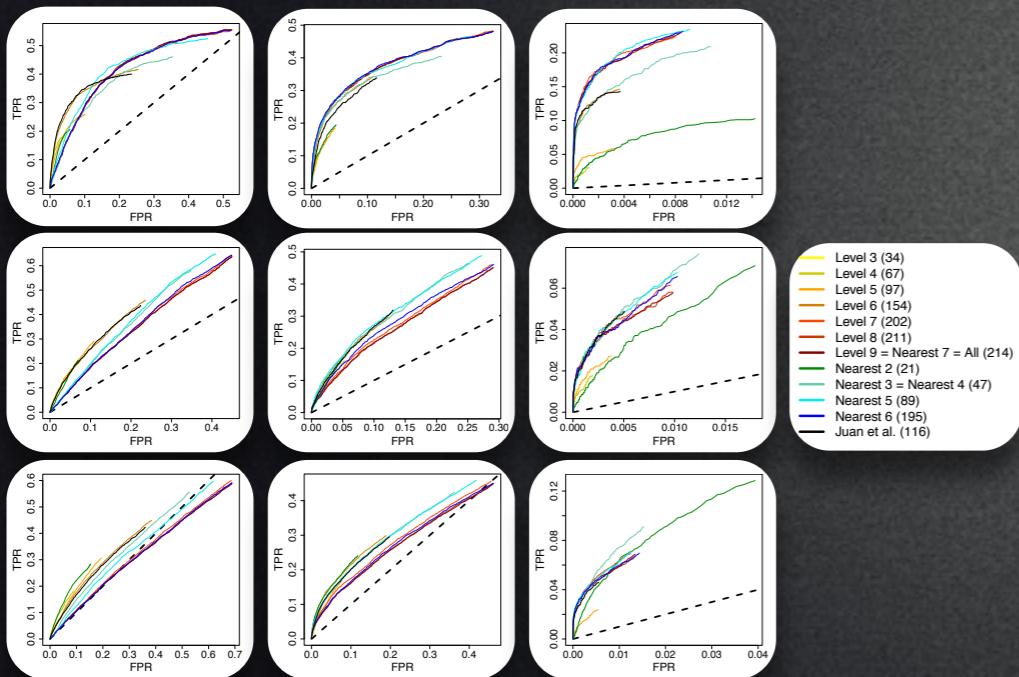


- Level 3 (34)
- Level 4 (67)
- Level 5 (97)
- Level 6 (154)
- Level 7 (202)
- Level 8 (211)
- Level 9 = Nearest 7 = All (214)
- Nearest 2 (21)
- Nearest 3 = Nearest 4 (47)
- Nearest 5 (89)
- Nearest 6 (195)
- Juan et al. (116)

Pathways

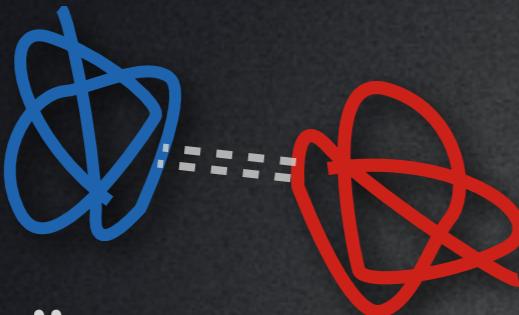
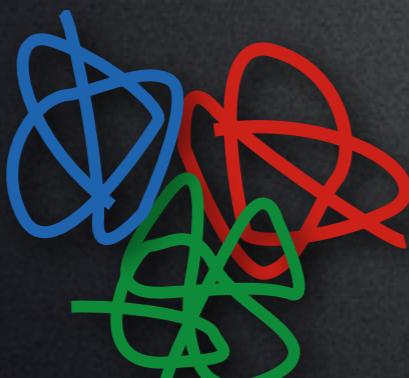


Conclusions



- Not always more organisms is better
- Redundancy negatively affect mirrortree predictions
- Sampling close/distant organisms influences the prediction of different types of interactions

AUCs

		Level9 (=all)	Nearest2
 “recent”	MINE_ECOLI	0.12	0.83
	PABA_ECOLI	0.28	0.96
	DHAS_ECOLI	0.17	0.81
	GSHB_ECOLI	0.30	0.93
 “old”	DPO3A_ECOLI	0.70	0.11
	DPO3B_ECOLI	0.64	0.22
	RPOB_ECOLI	0.82	0.48
	RPOA_ECOLI	0.81	0.48
	ZNUB_ECOLI	1.00	0.36
	ZNUC_ECOLI	0.99	0.41
	ZNUA_ECOLI	0.98	0.79



4. Co-evolution significance

Confidence evaluation

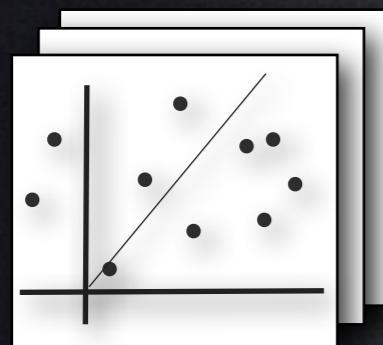
Distances matrix A

	R1	R2	R3	R4	R5	R6
R1						
R2						
R3						
R4						
R5						
R6						

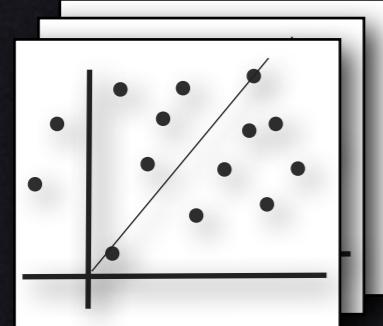
Distances matrix B

	S2	S3	S4	S5	S6	S7	S8
S2							
S3							
S4							
S5							
S6							
S7							
S8							

$n=10$

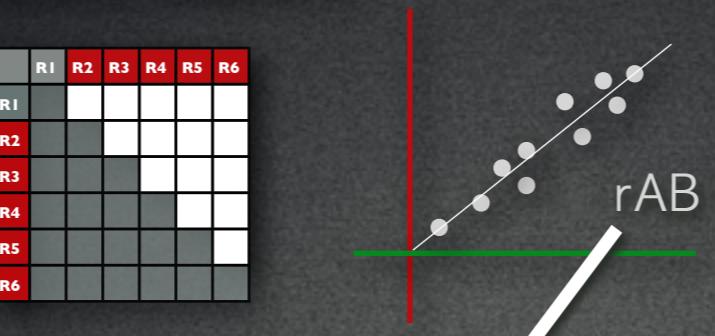


$n=15$

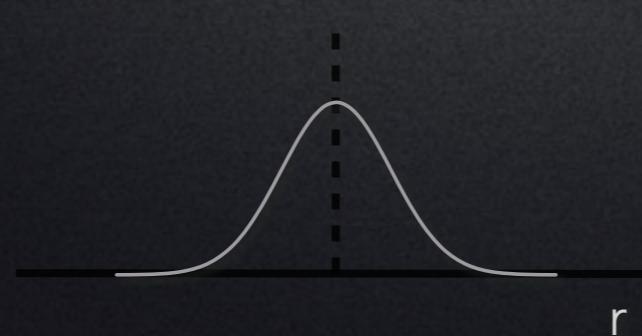
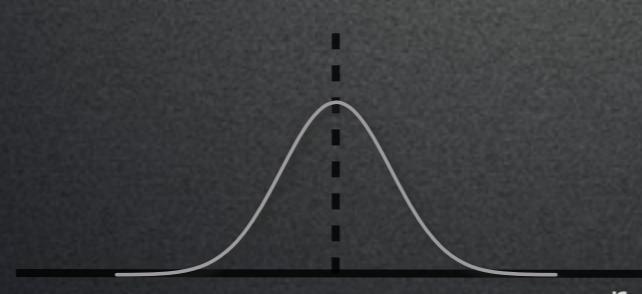
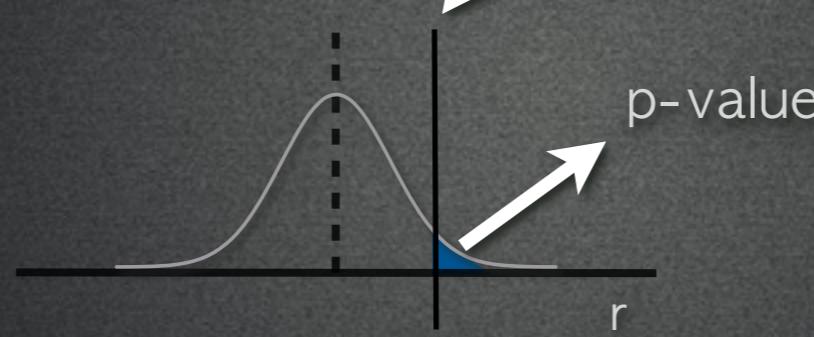


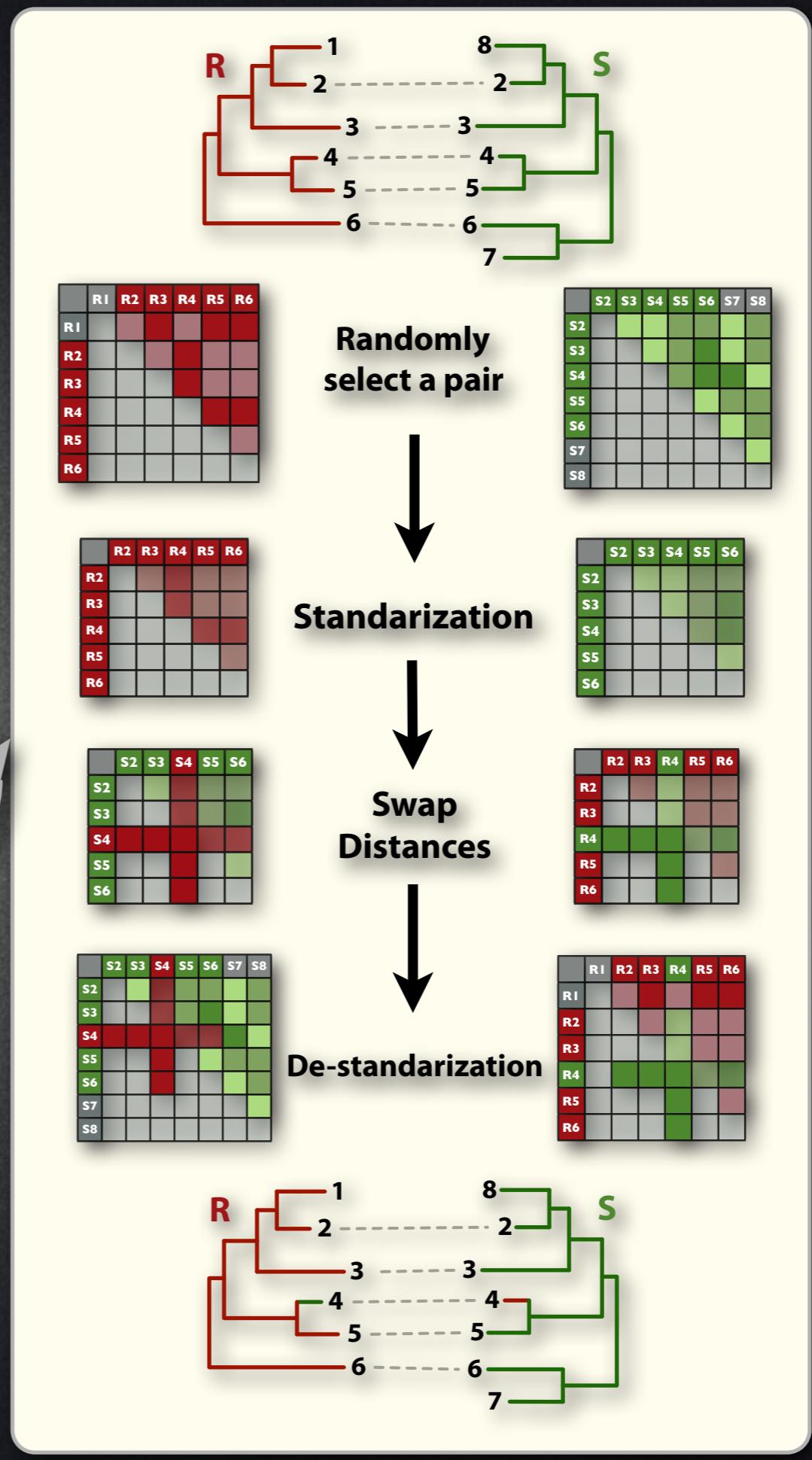
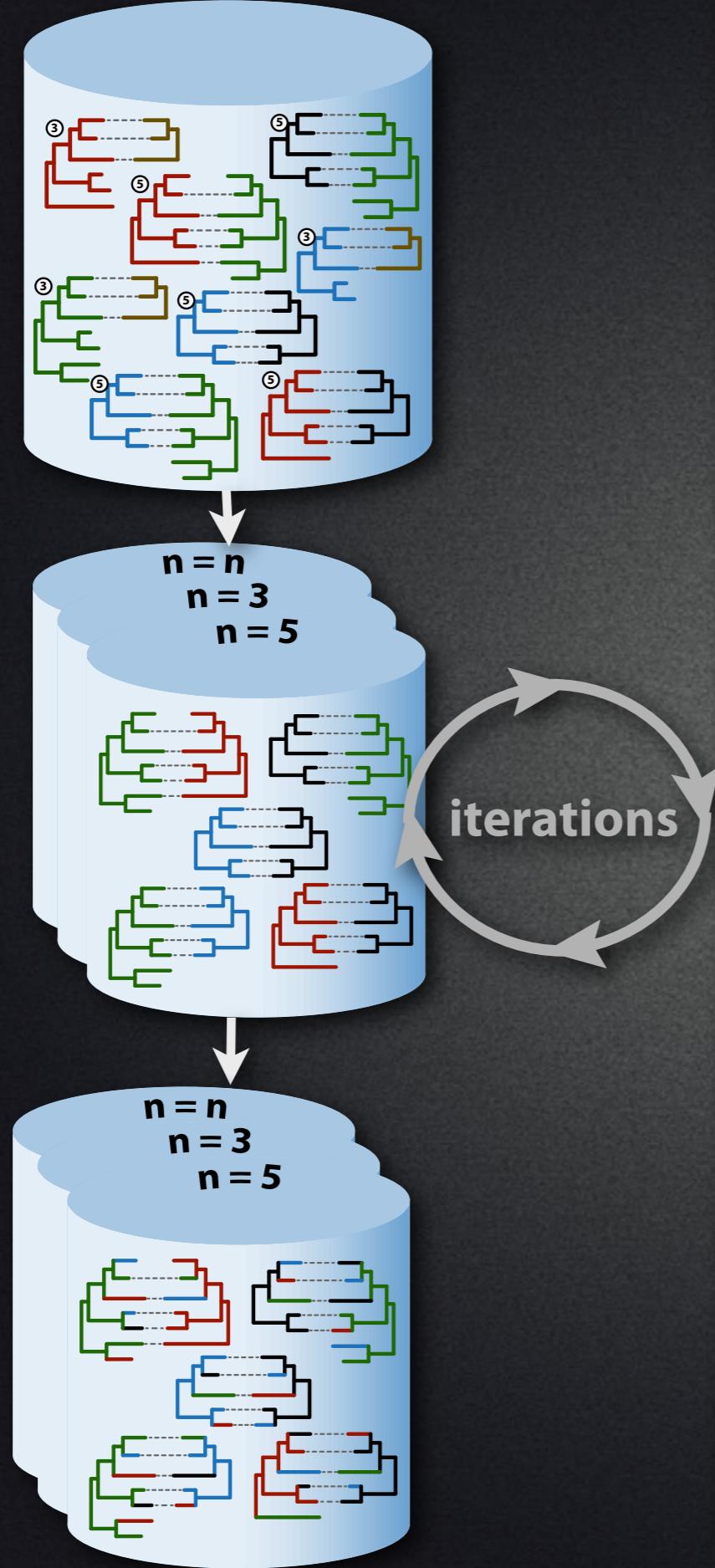
$n=n$

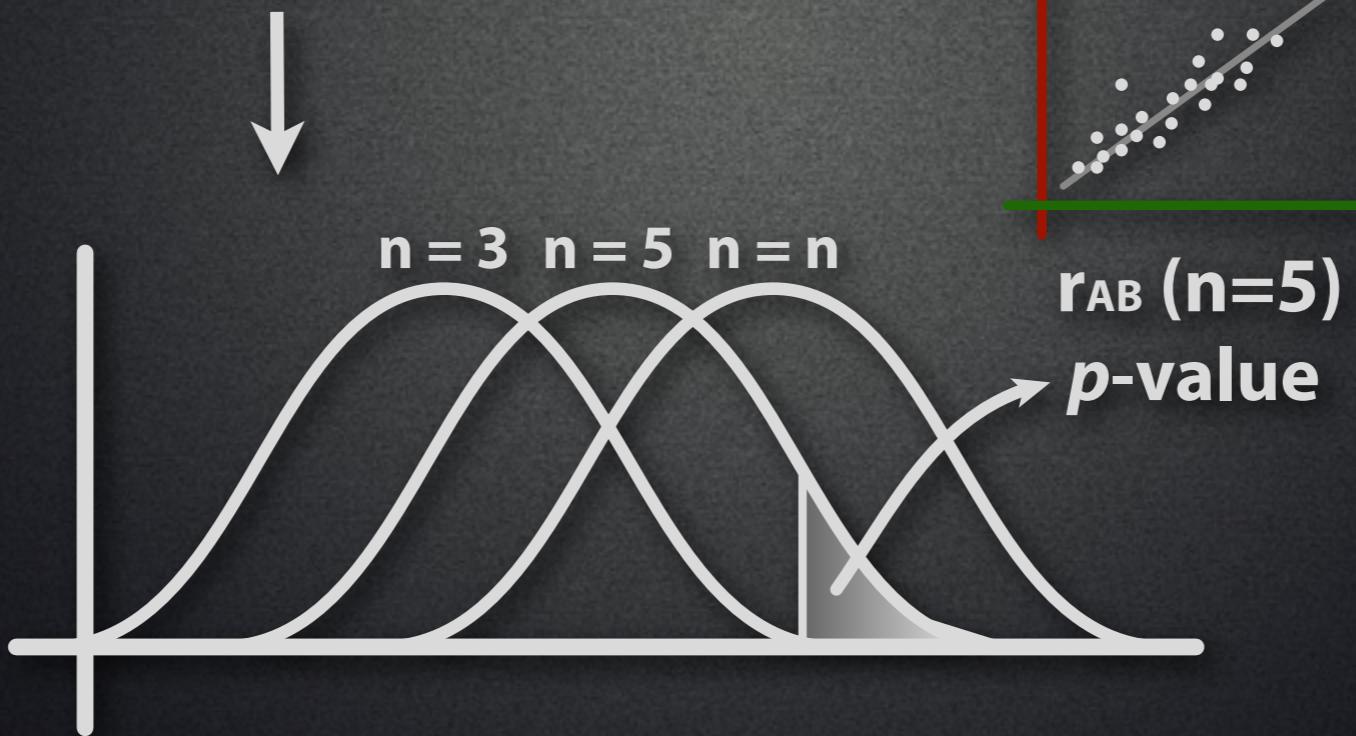
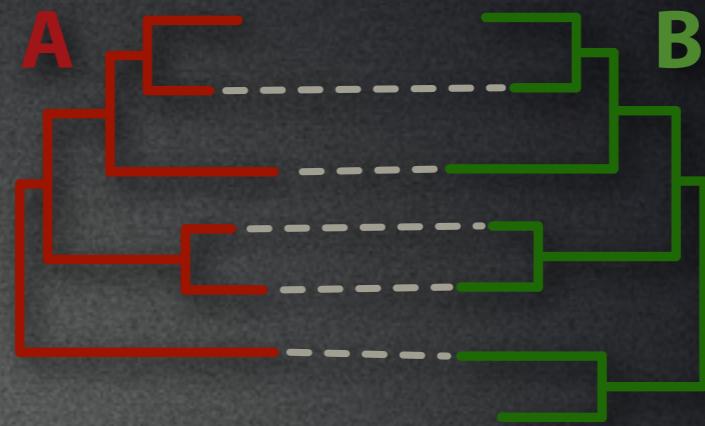
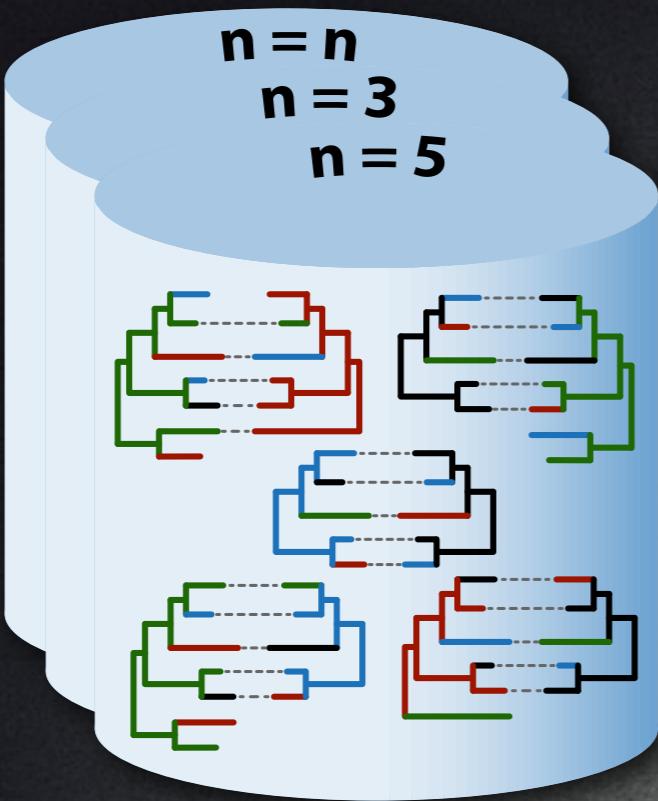
...



Non-specific null model

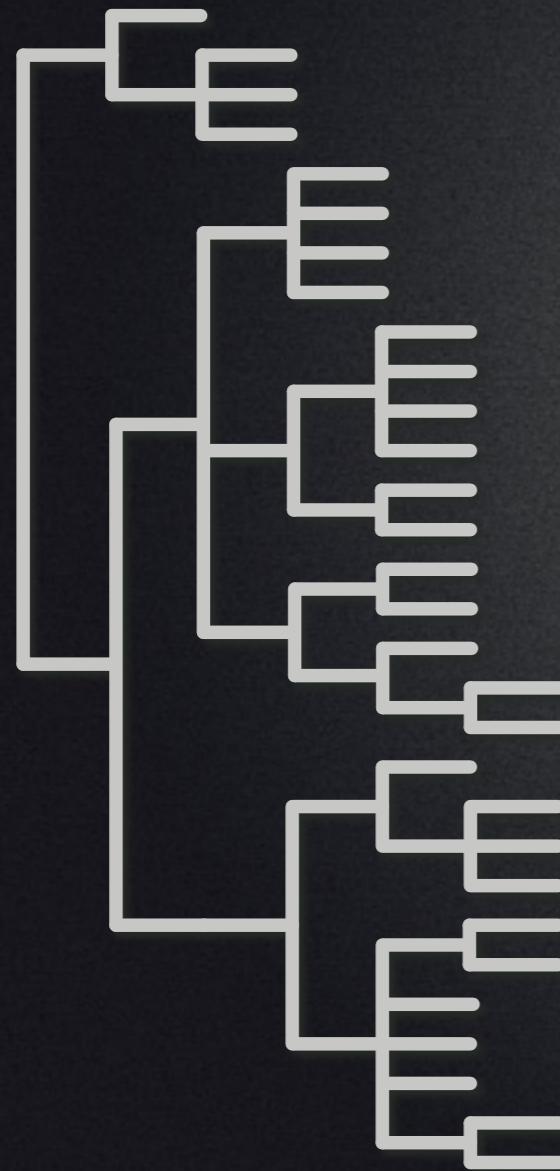




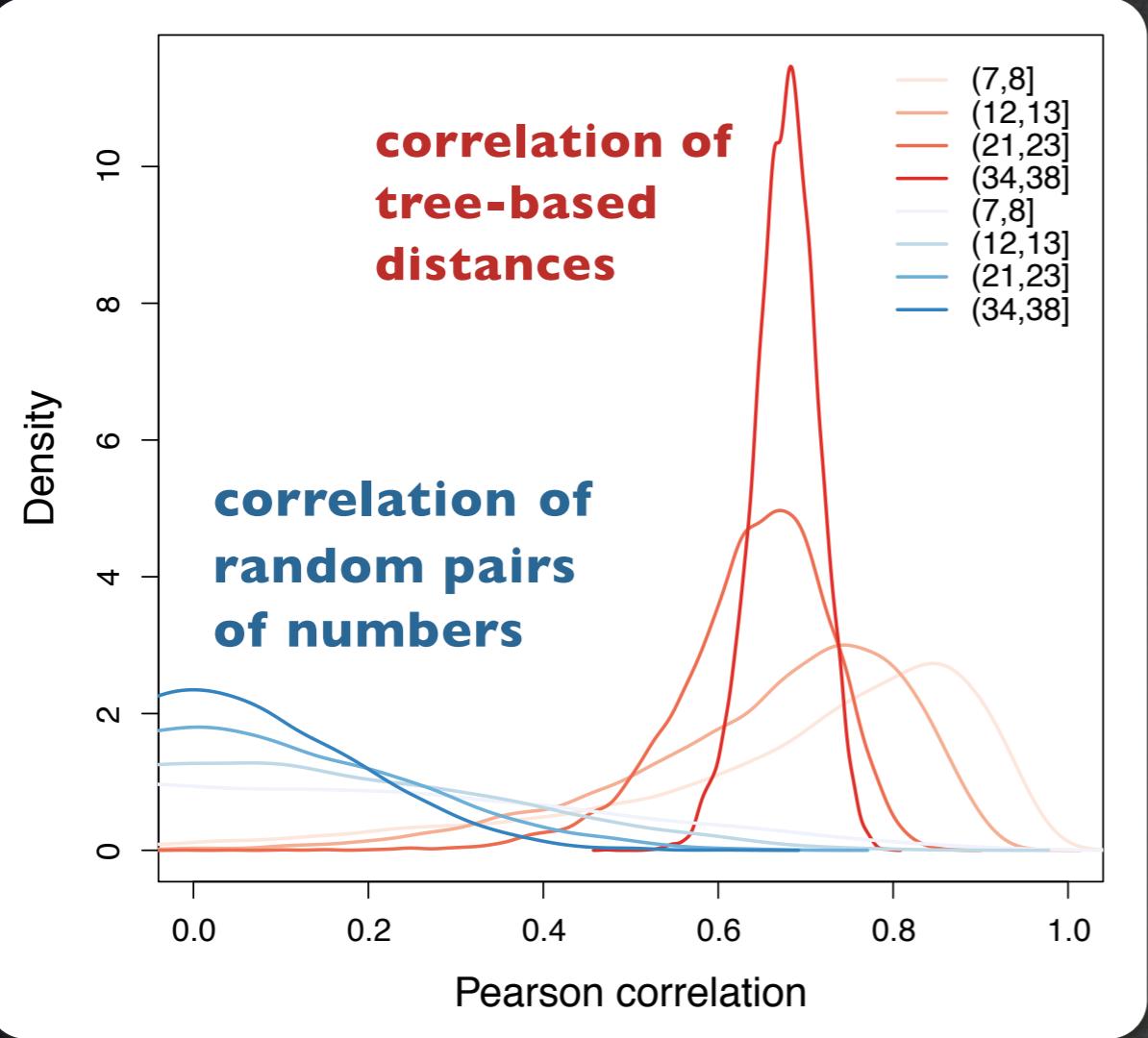


p-mirrortree

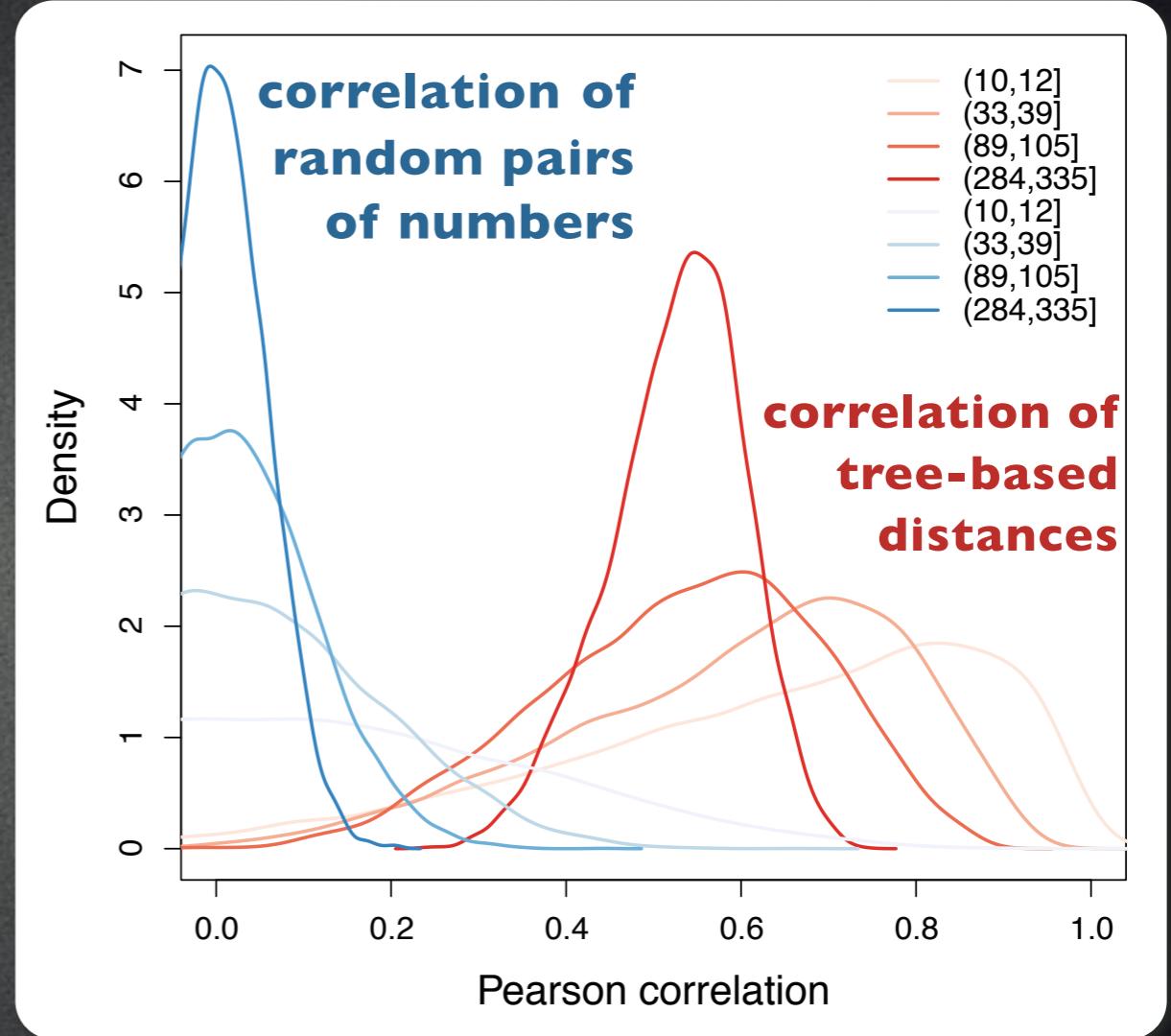
Reference set of organisms



- Prokaryotes in KEGG 2000-2010
- Redundant vs Non-redundant

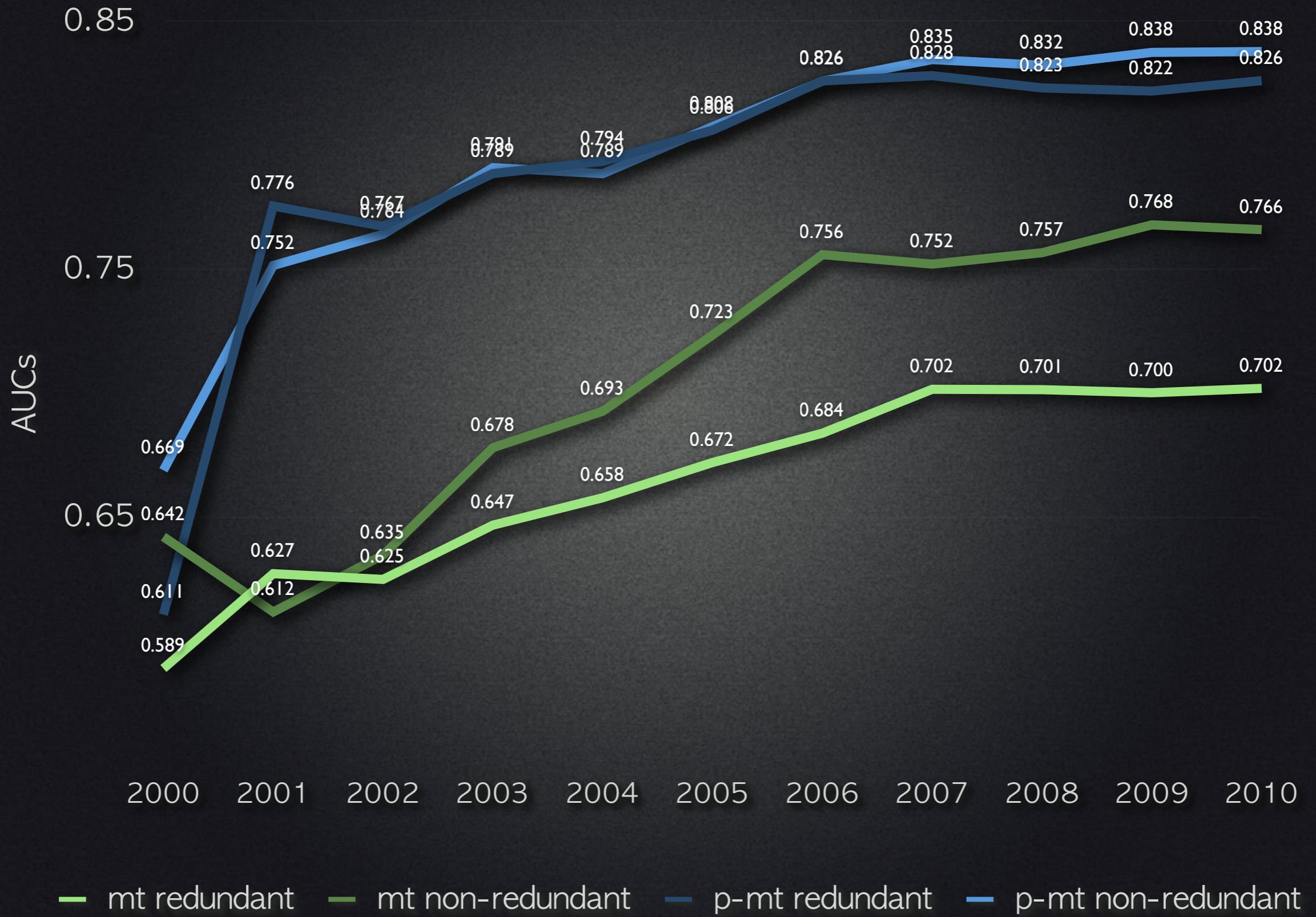


2001



2010

Complexes

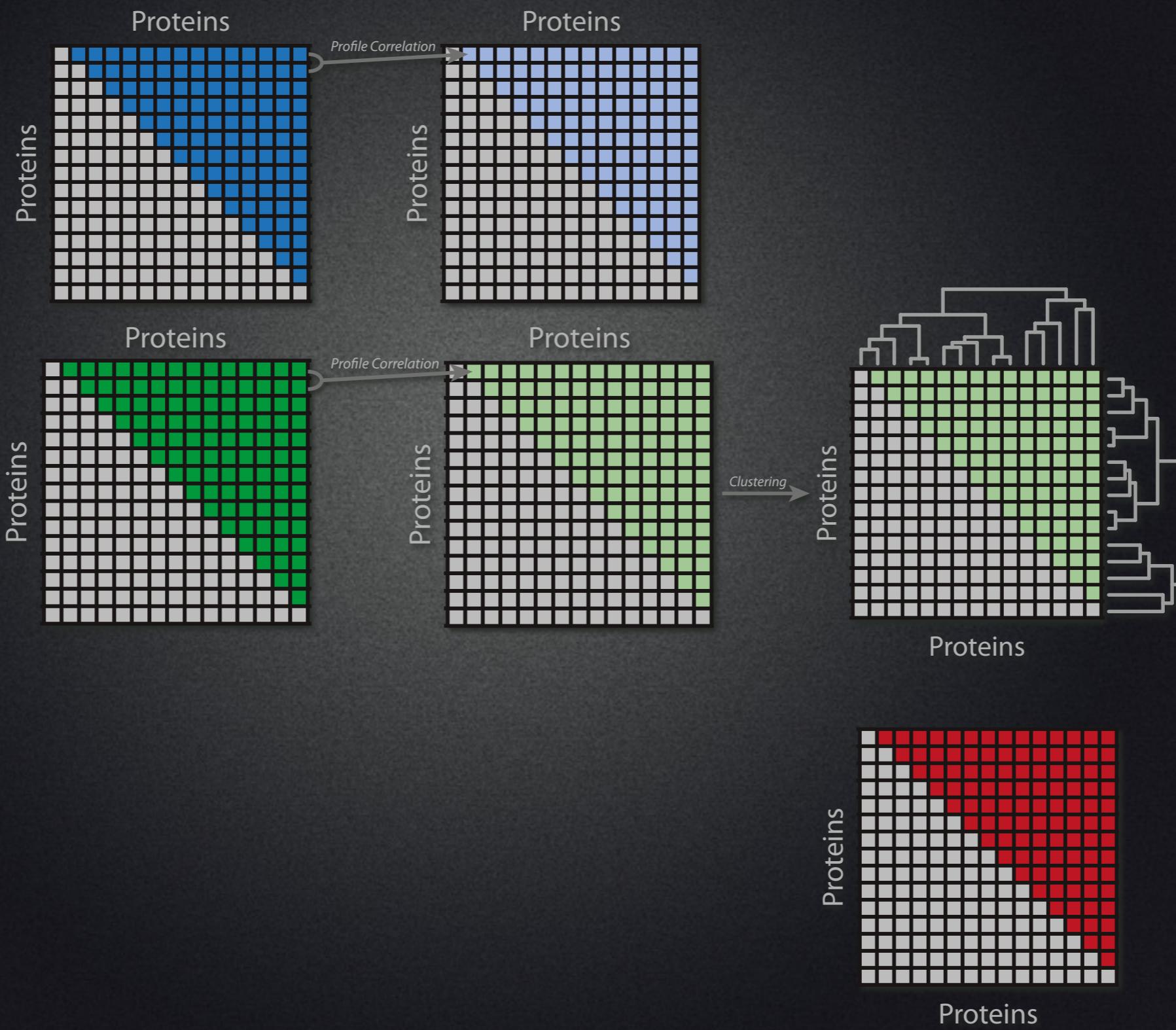


Context methods using p-values

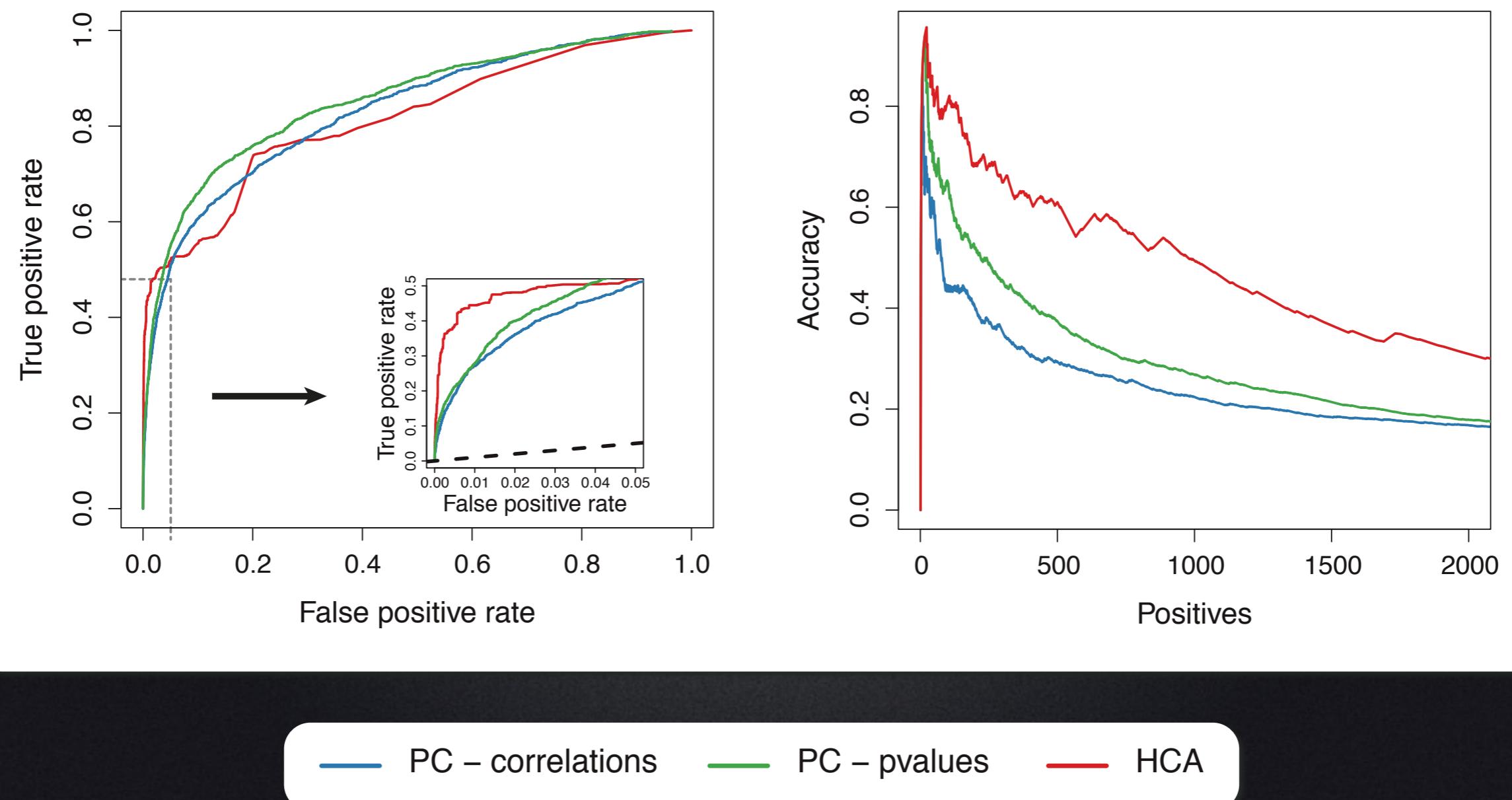
Profile
Correlation
mirrortree

Profile
Correlation
p-mirrortree

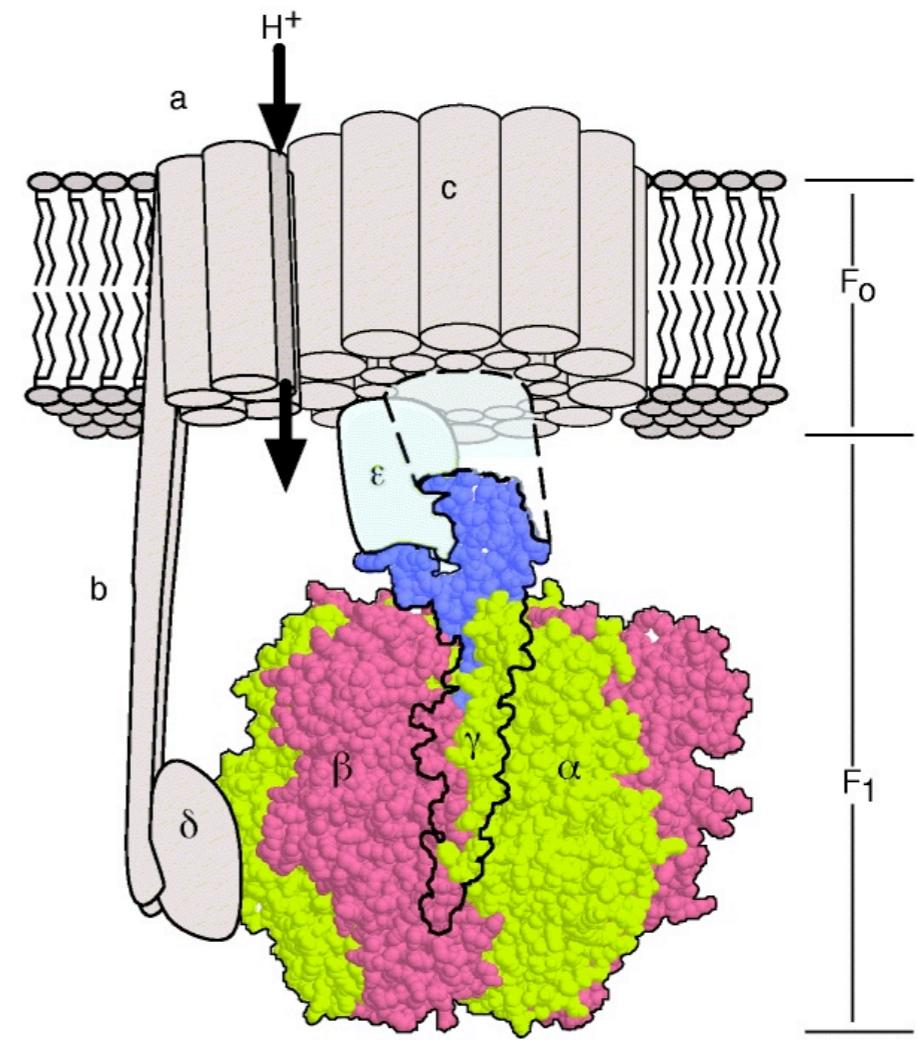
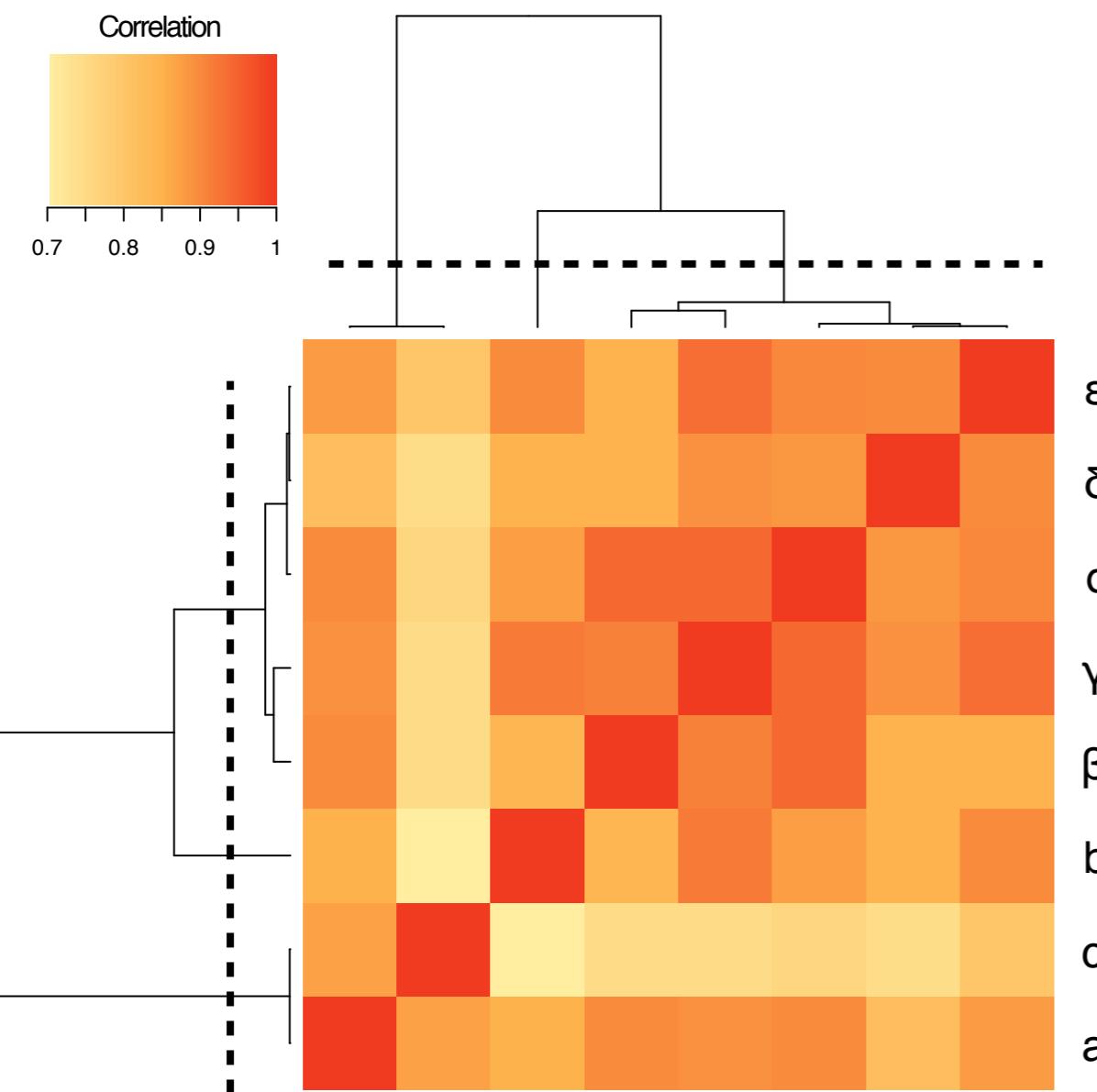
Hierarchical
Clustering
Analysis
(HCA)



Context methods using p-values



ATPase coevolutionary map



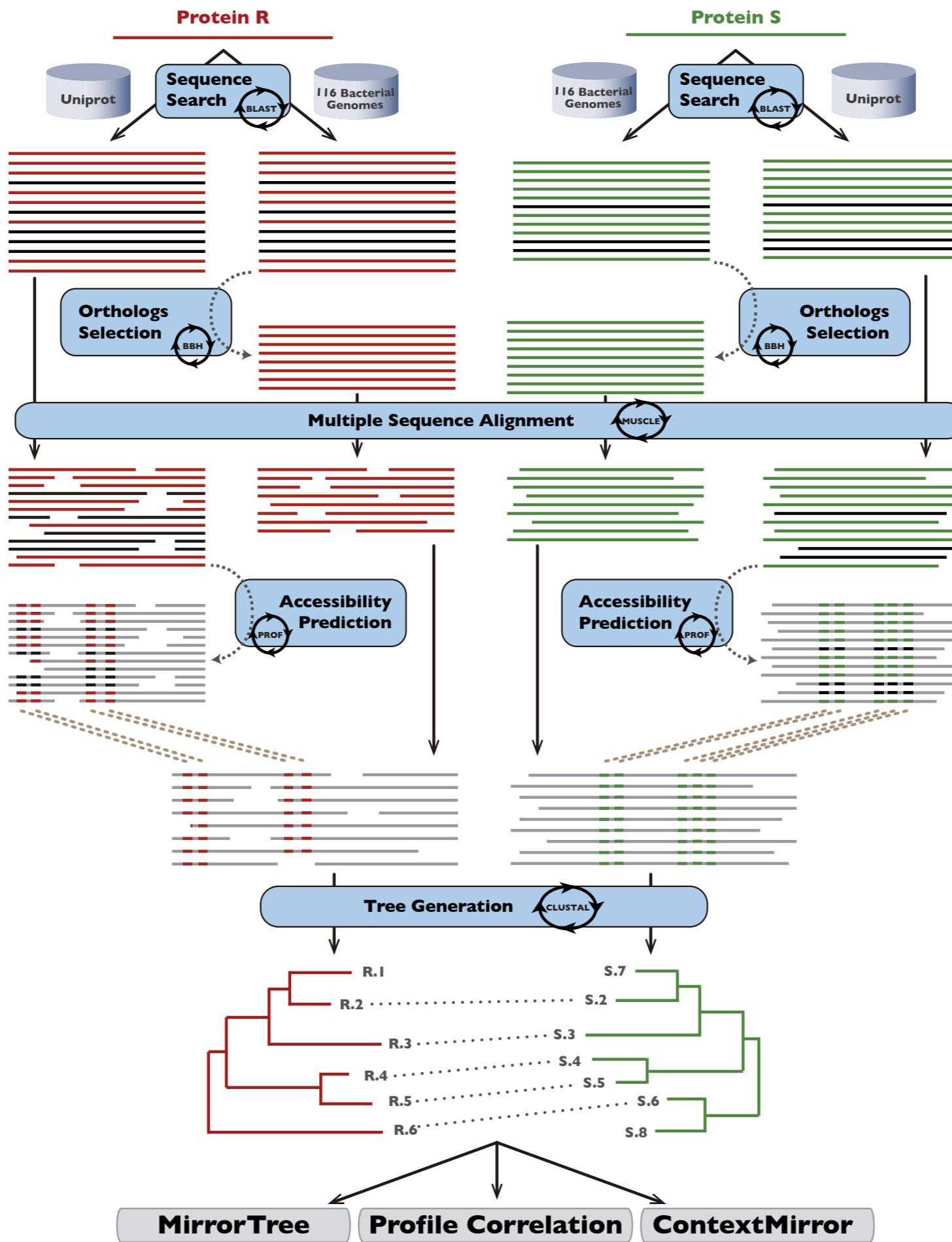
Conclusions

- MirrorTree Server allows to explore the co-evolution of protein families in a taxonomic context, independently of the user's expertise
- In agreement with previous observations, mirrortree family of methods proved being particularly accurate predicting different types of interactions at large scale
- Incorporation of predicted solvent accessibility helps the co-evolution-based prediction of binary physical interactions when context-based methodologies are applied
- The set of organisms used to generate the phylogenetic trees conditions the final performance and the type of interaction predicted by the methods
- The novel methodology, p-mirrortree, outperforms the existing mirrortree-based approaches and extends their range of applicability. The improvement is independent of the set of organisms and might help the input of context-based approaches.

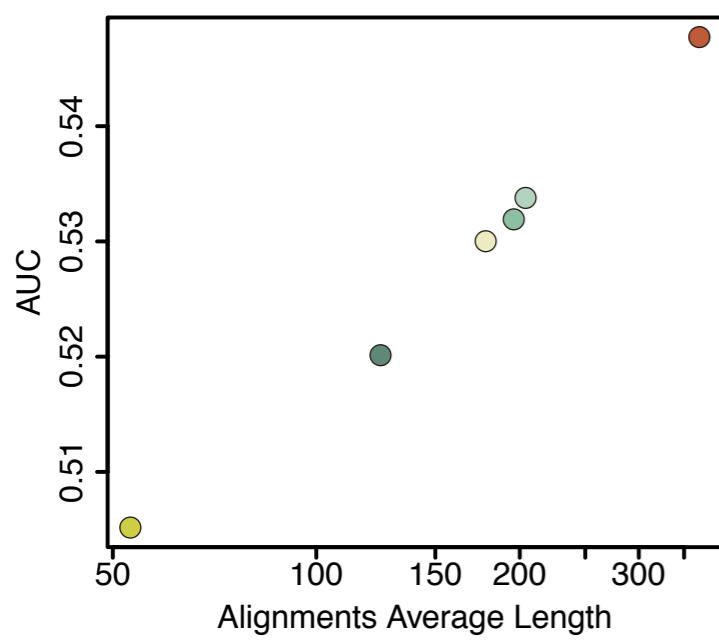
Tesis Doctoral
Improving Co-evolution Based
Methods for Protein-Protein
Interaction Prediction



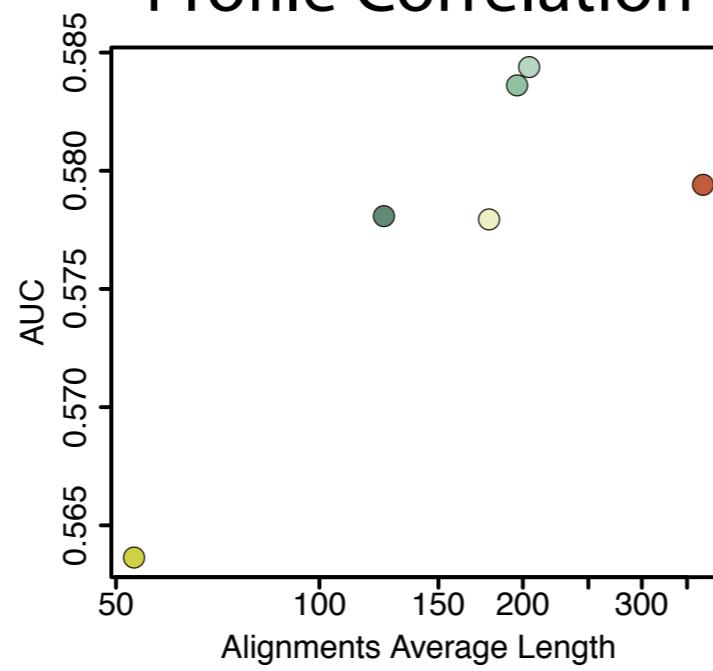
David Ochoa
Universidad Autónoma de Madrid



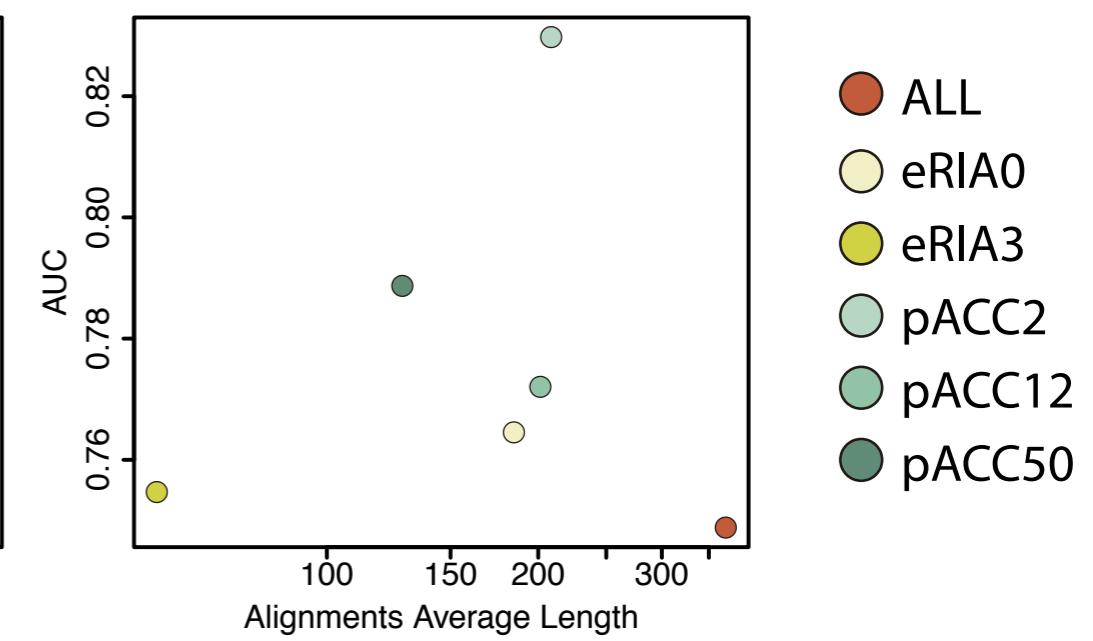
MirrorTree



Profile Correlation



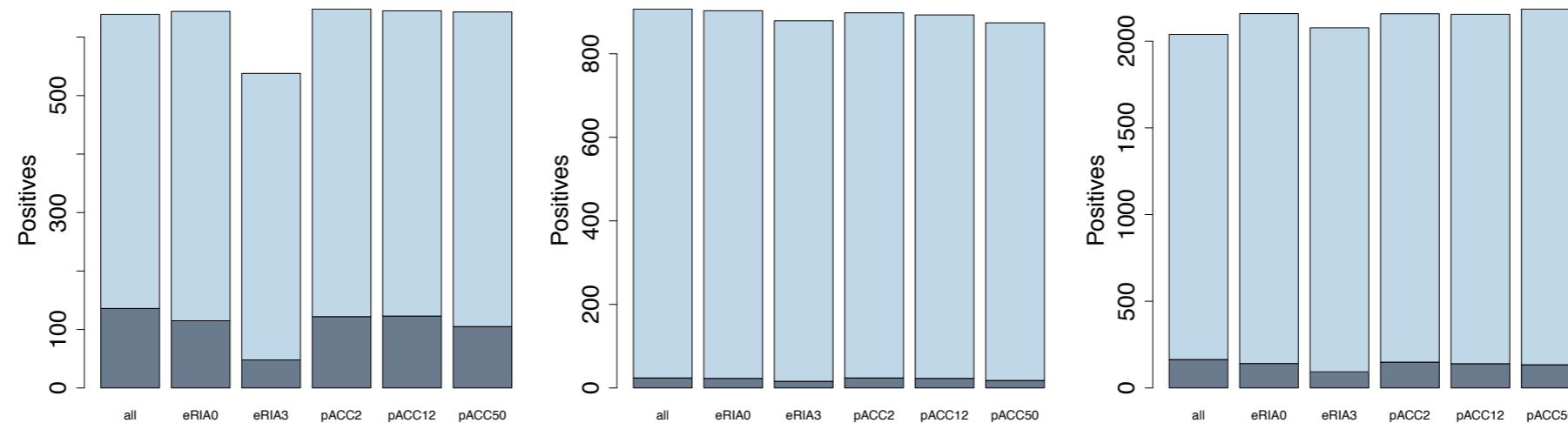
ContextMirror



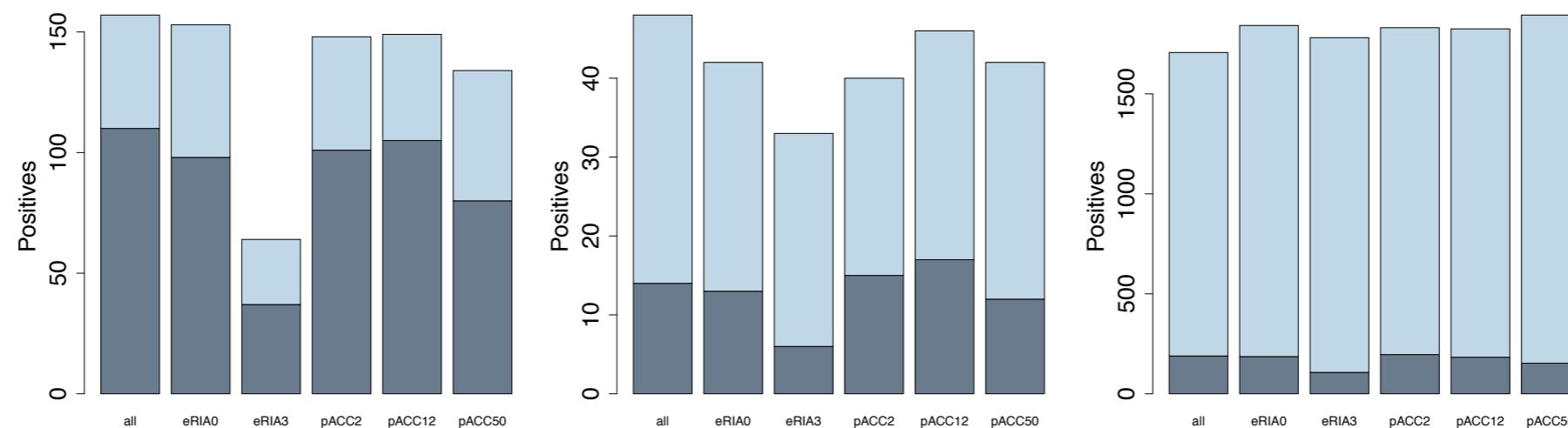
- ALL
- eRIA0
- eRIA3
- pACC2
- pACC12
- pACC50

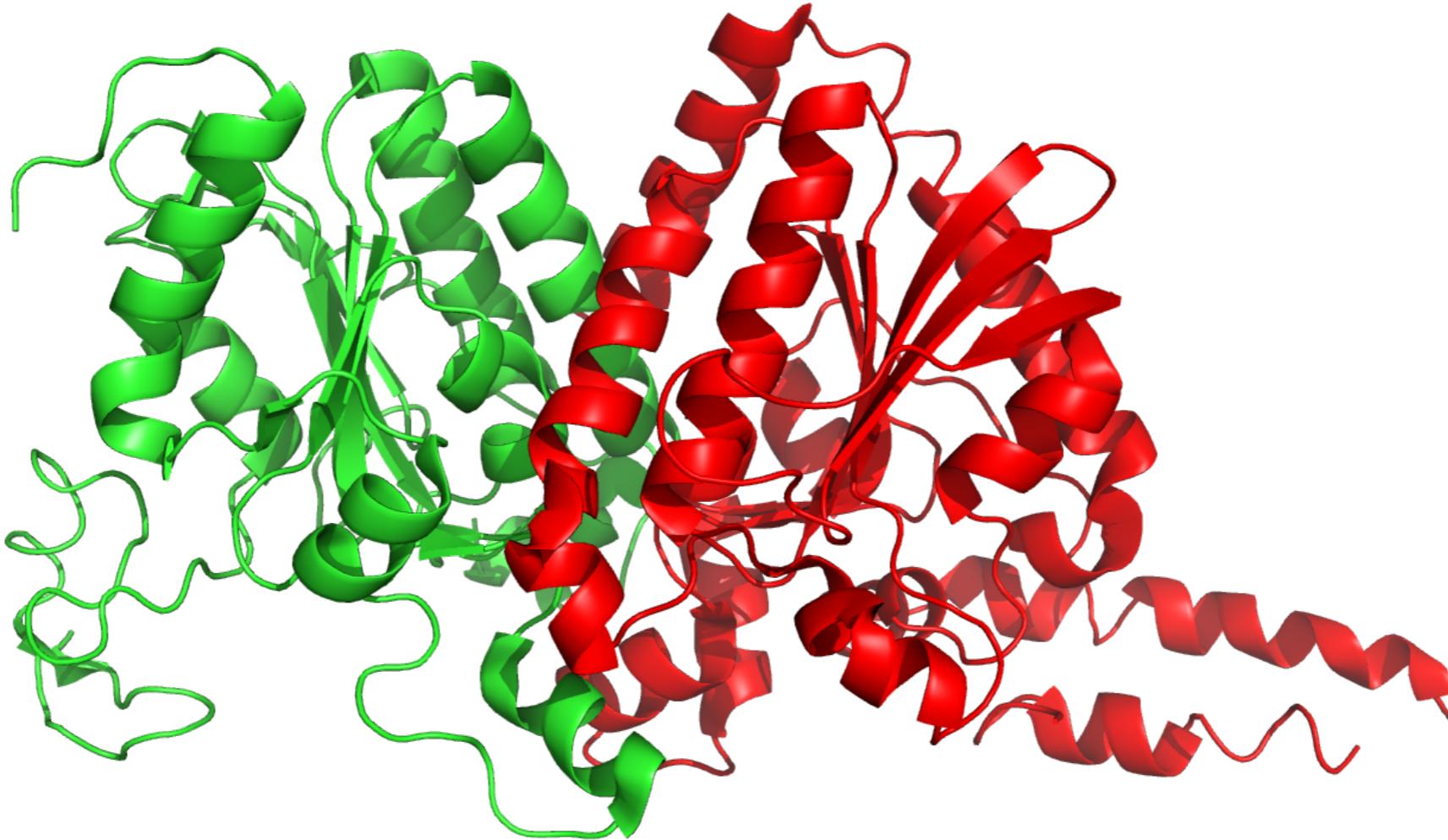
Profile Correlation

MirrorTree

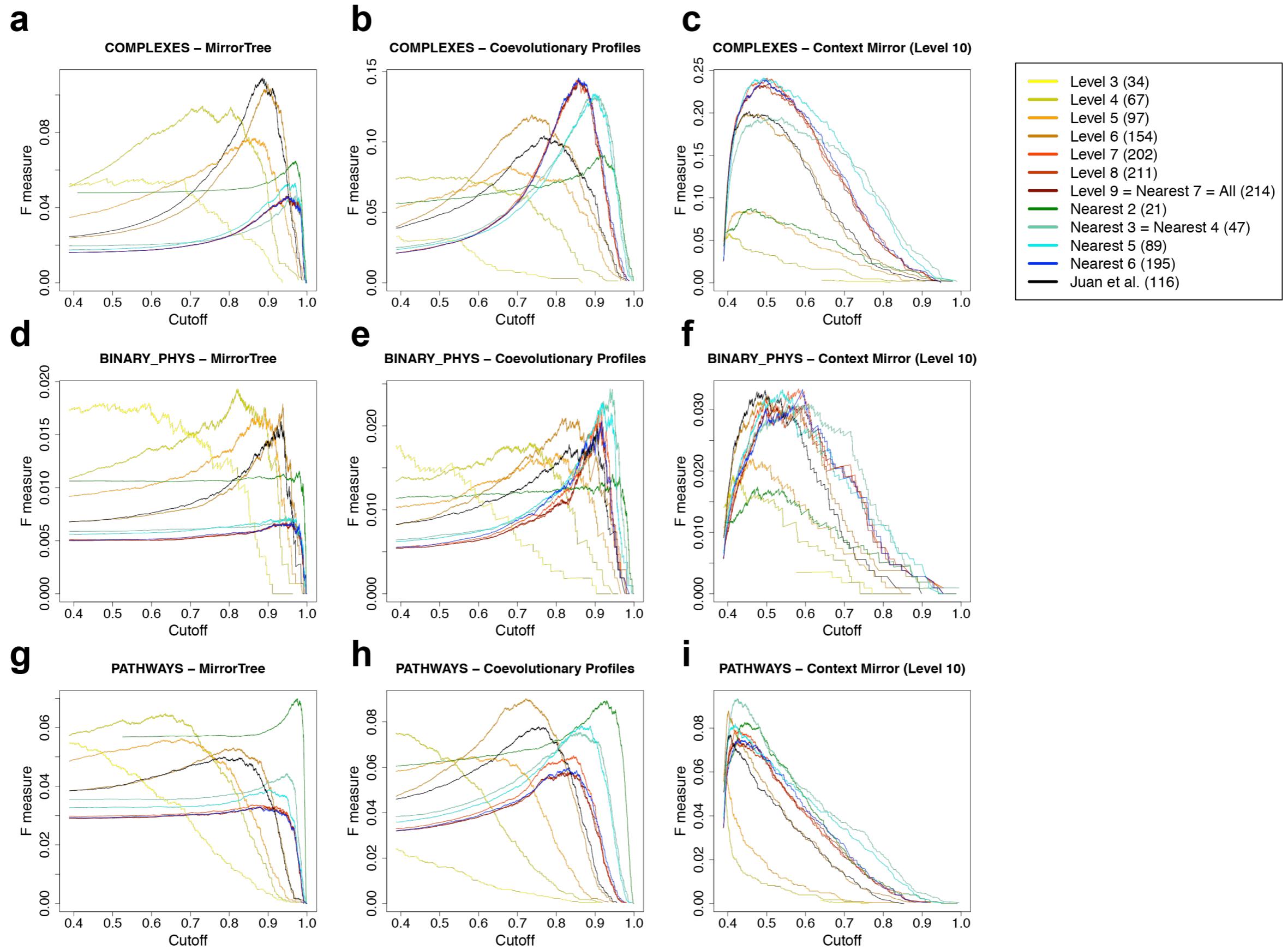


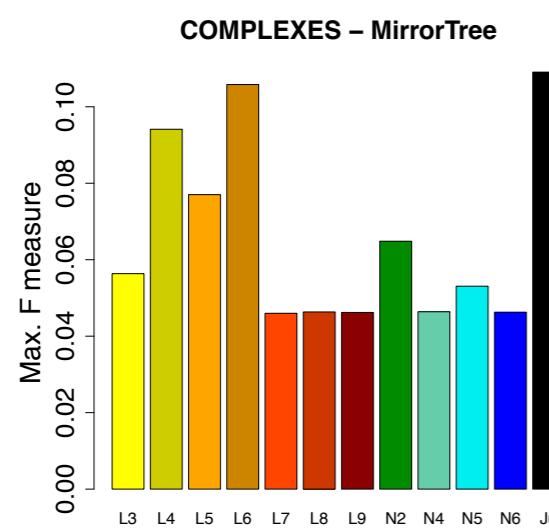
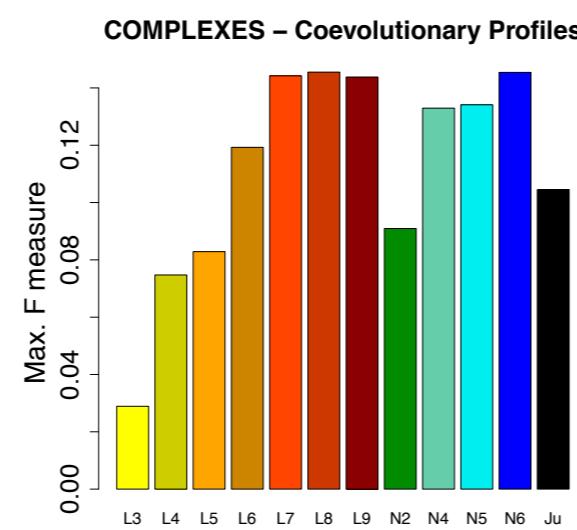
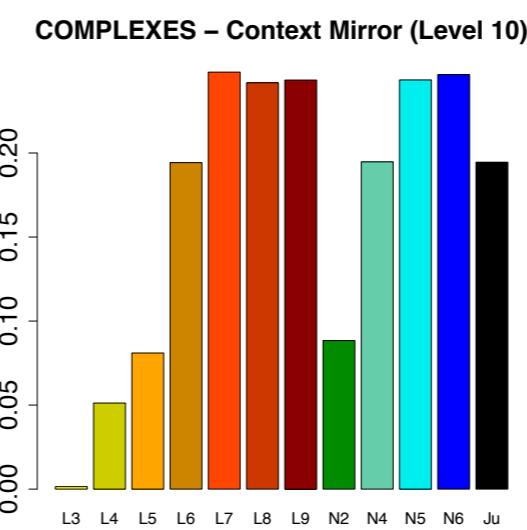
ContextMirror



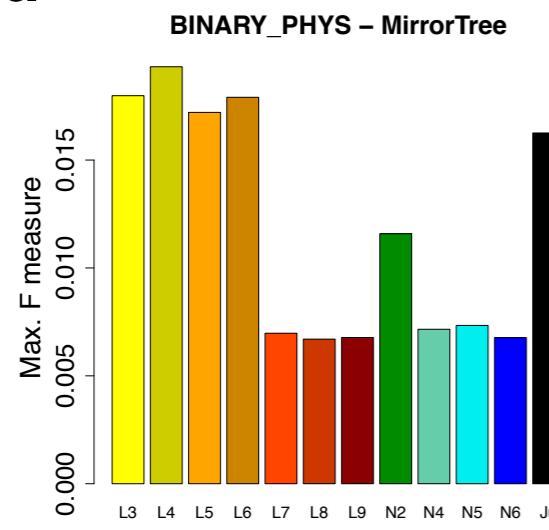
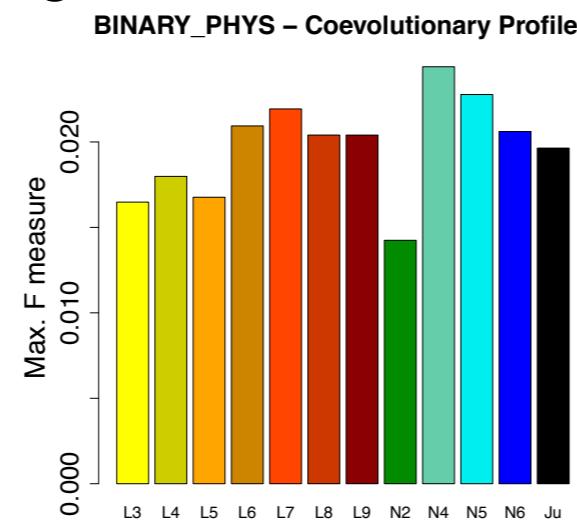
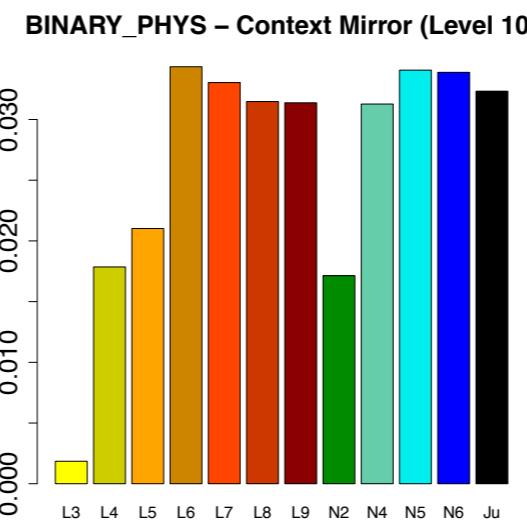
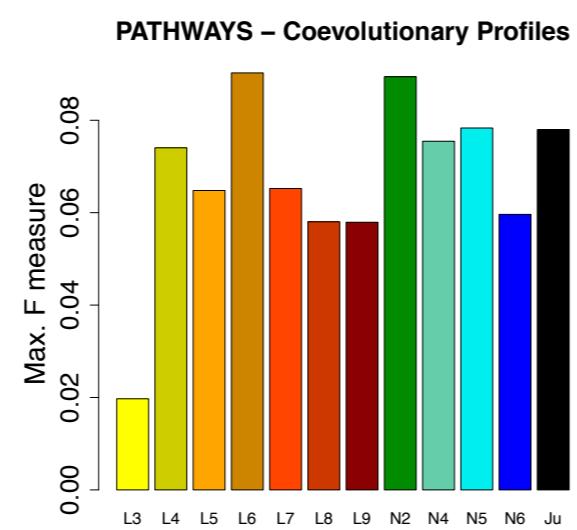
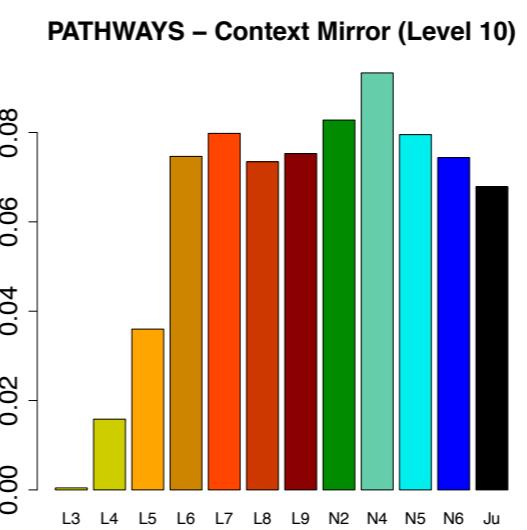


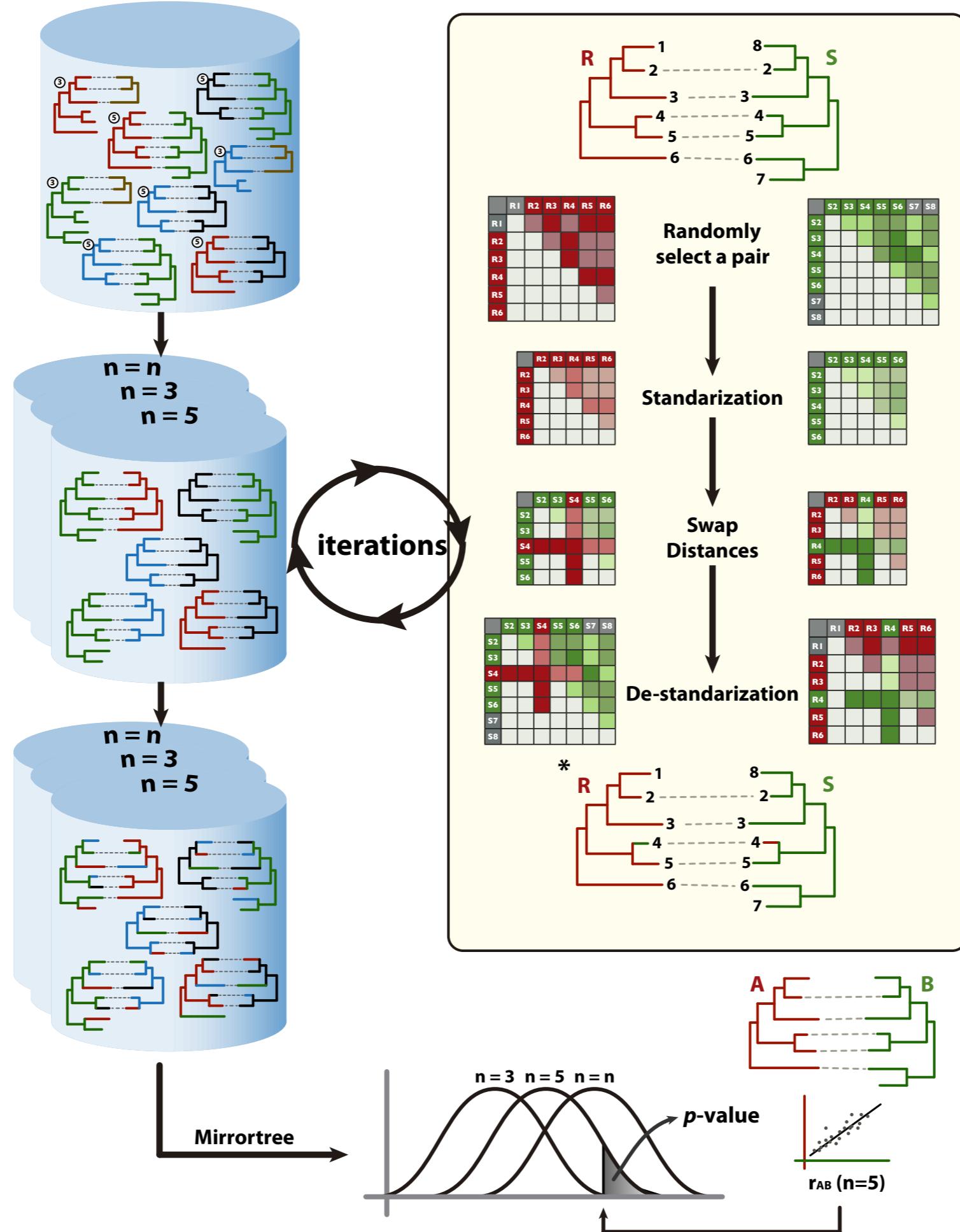
	ALL	eRIA0	eRIA3	pACC2	pACC12	pACC50
MirrorTree	0.9083	0.8982	0.8331	0.8998	0.9069	0.8924
Profile Correlation	0.9516	0.9531	0.9435	0.9605	0.9592	0.9328
ContextMirror	0.6068	0.6650	0.5410	0.6818	0.6828	0.5407



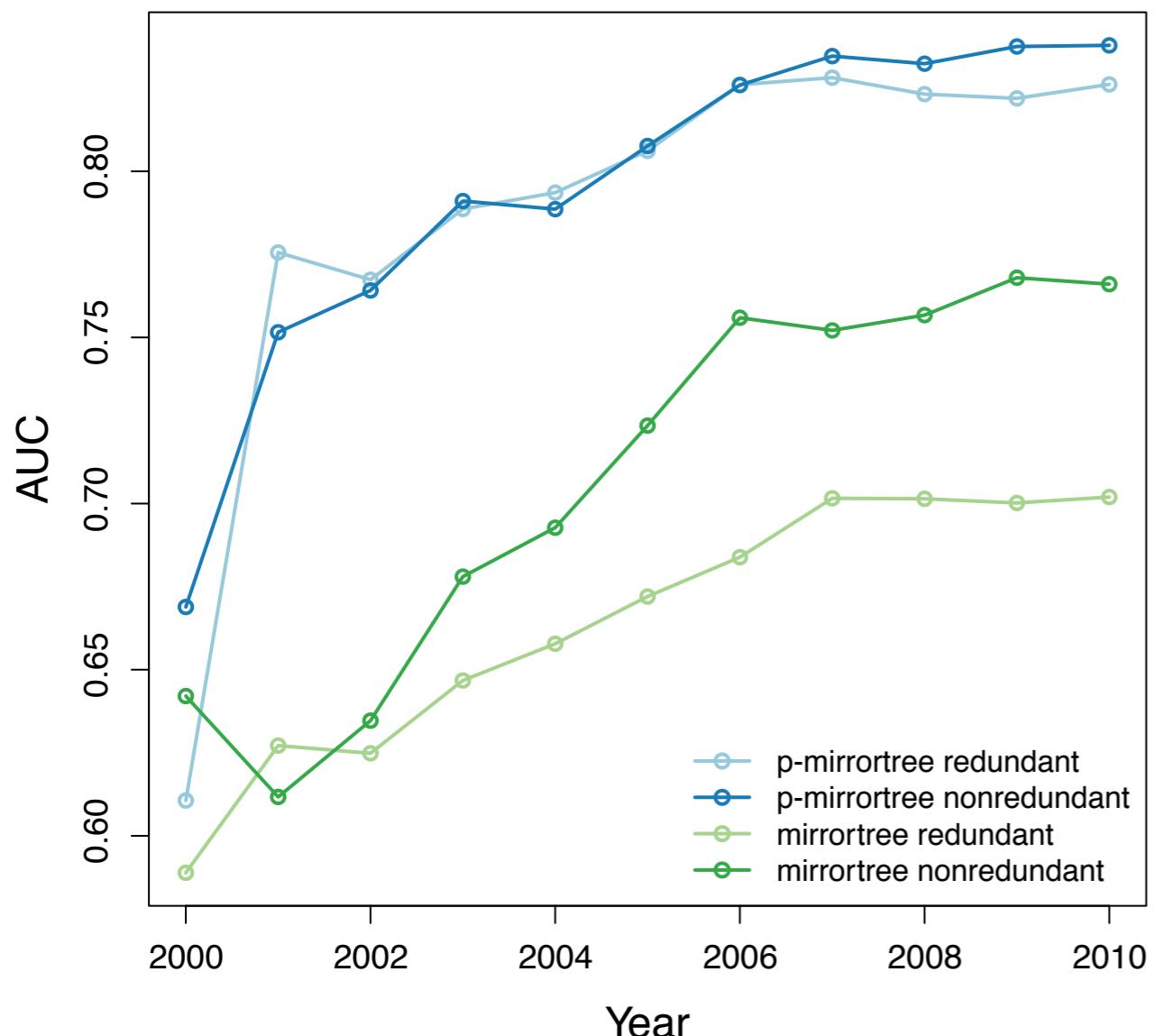
a**b****c**

- Level 3 (34)
- Level 4 (67)
- Level 5 (97)
- Level 6 (154)
- Level 7 (202)
- Level 8 (211)
- Level 9 = Nearest 7 = All (214)
- Nearest 2 (21)
- Nearest 3 = Nearest 4 (47)
- Nearest 5 (89)
- Nearest 6 (195)
- Juan et al. (116)

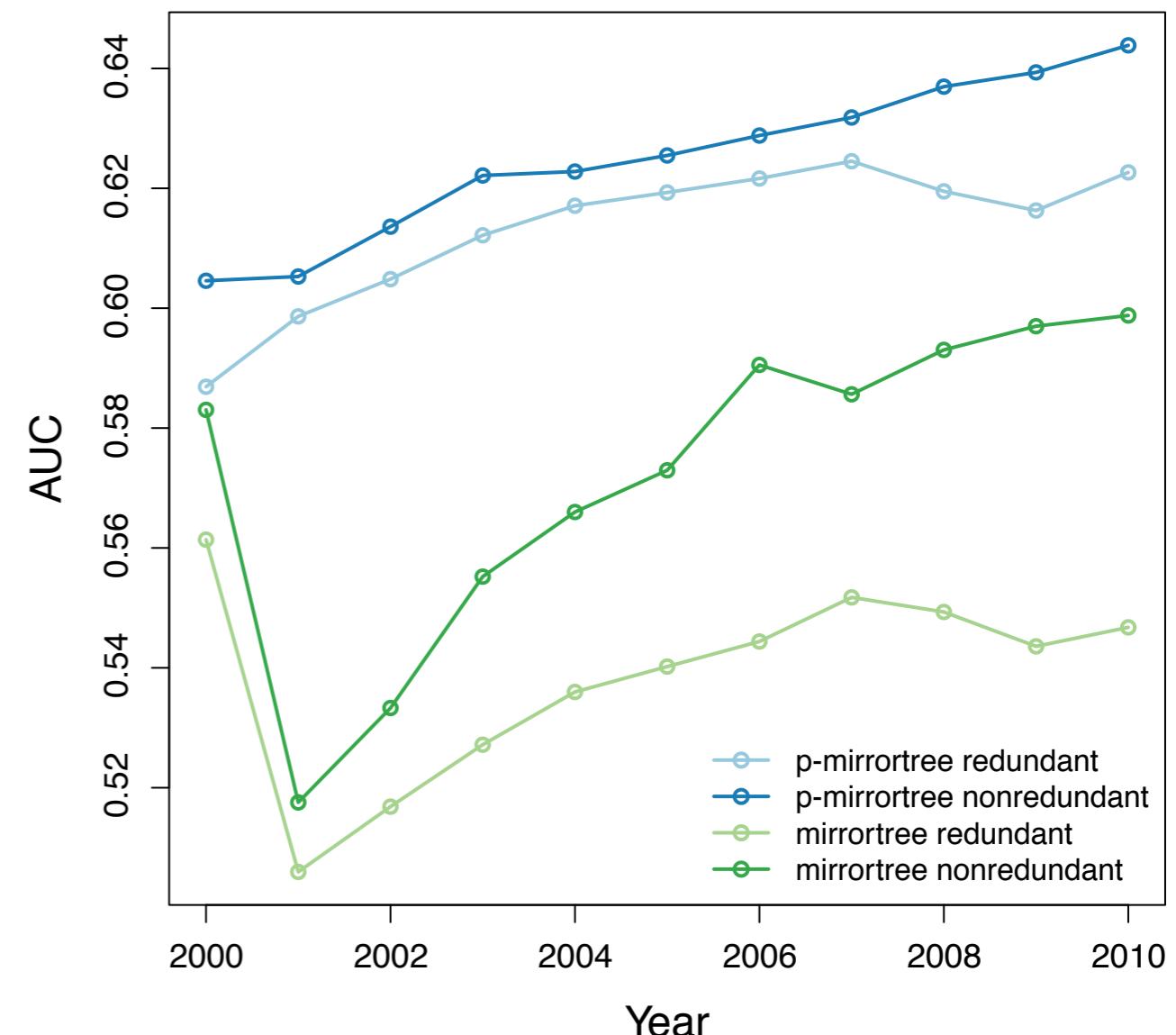
d**e****f****g****h****i**



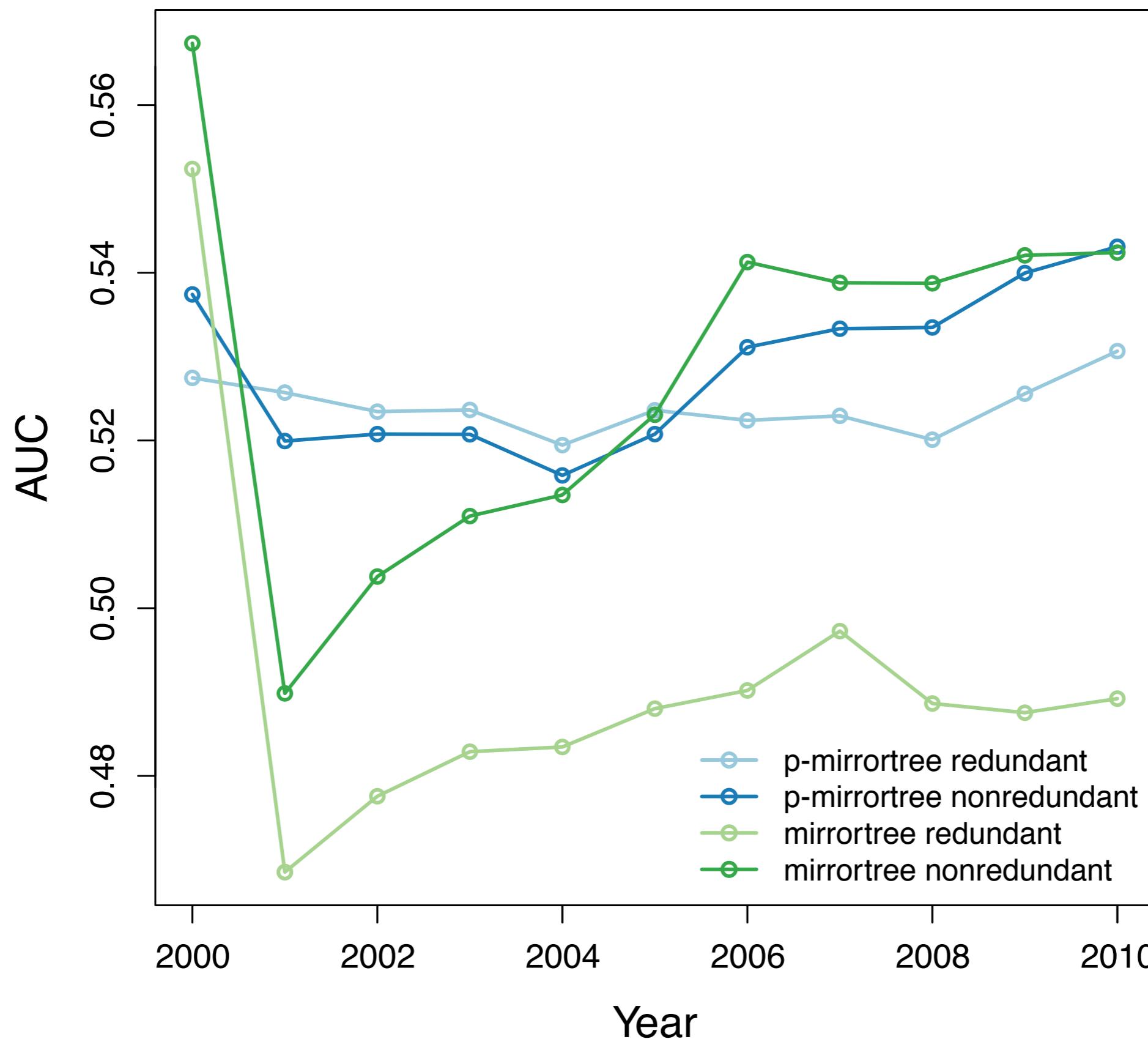
Complexes



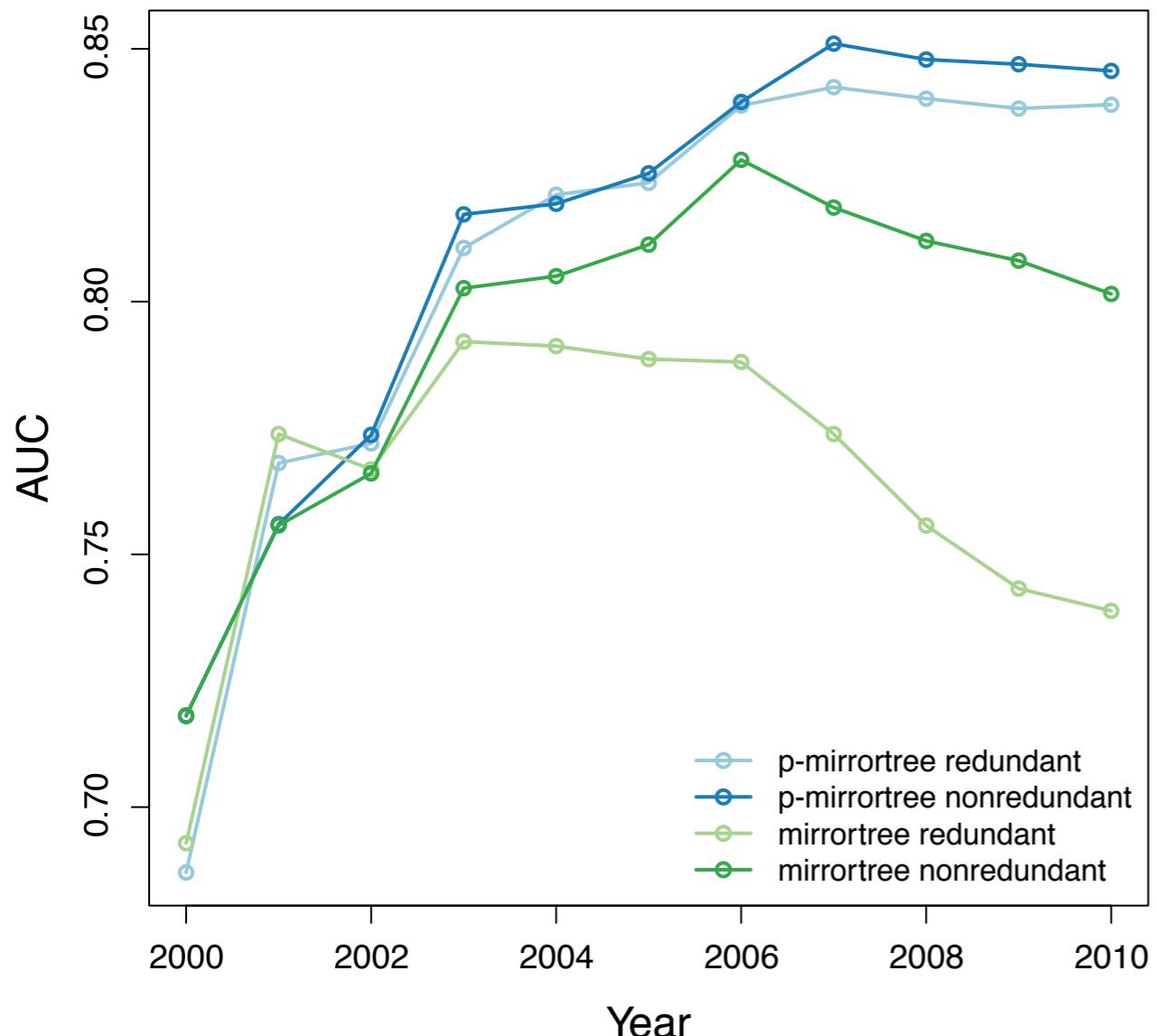
Binary Physical



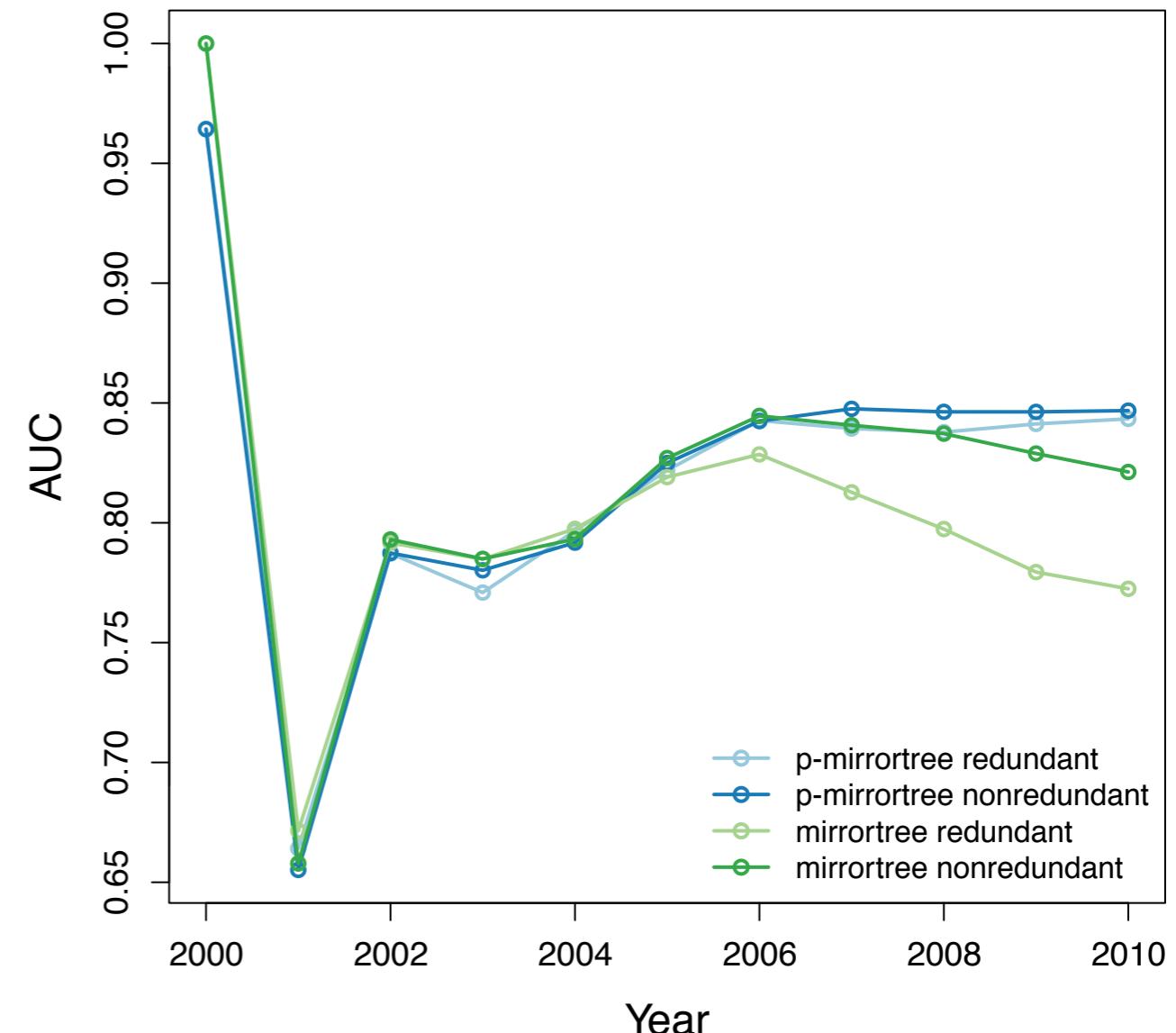
Pathways



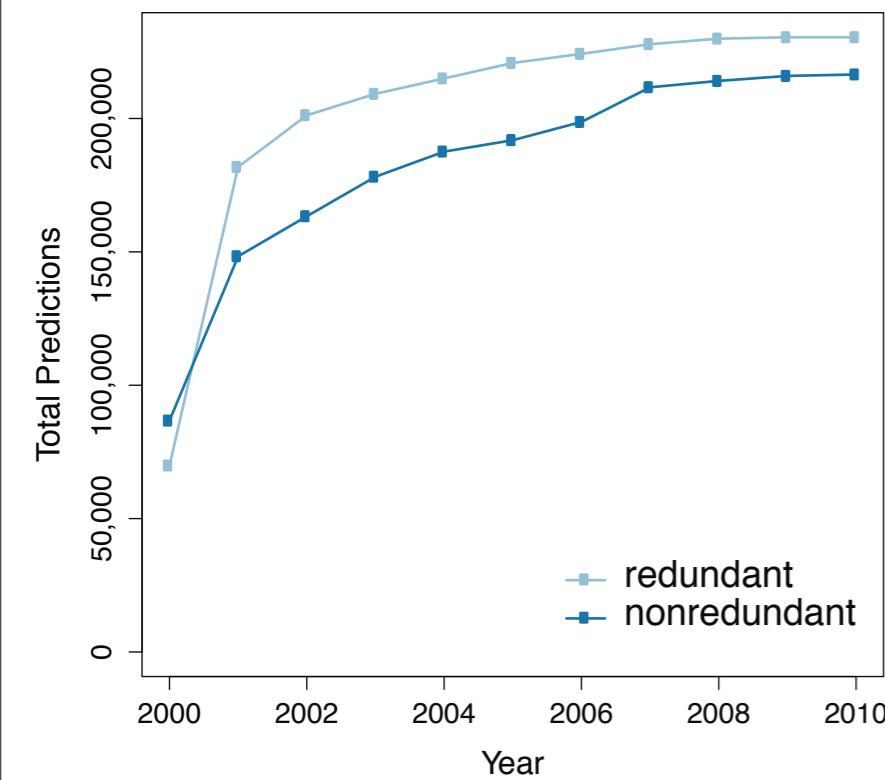
$n > 15$



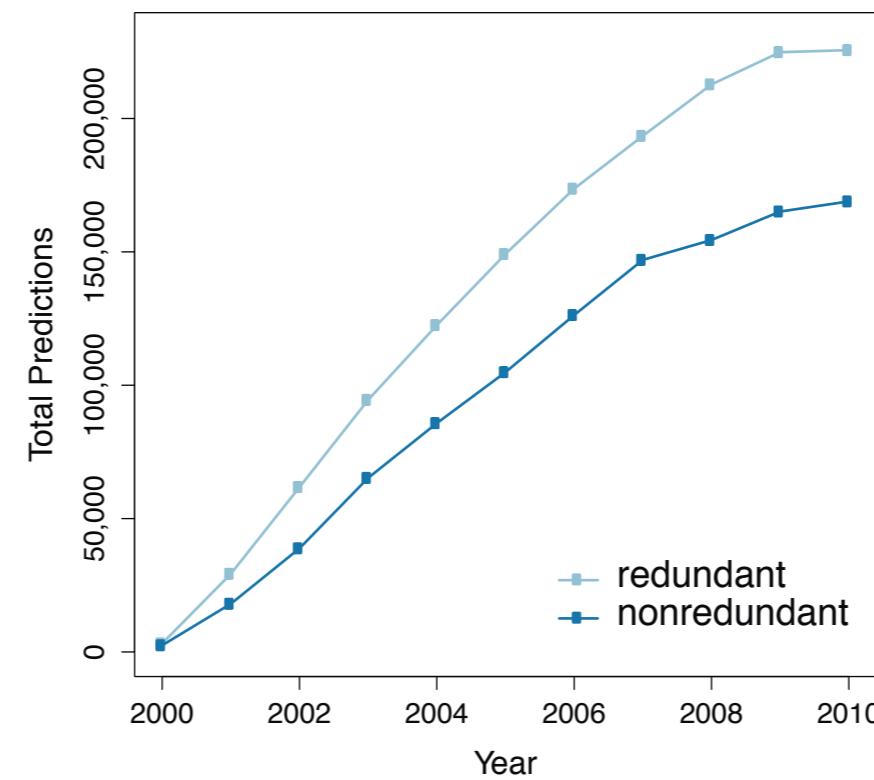
$n > 30$



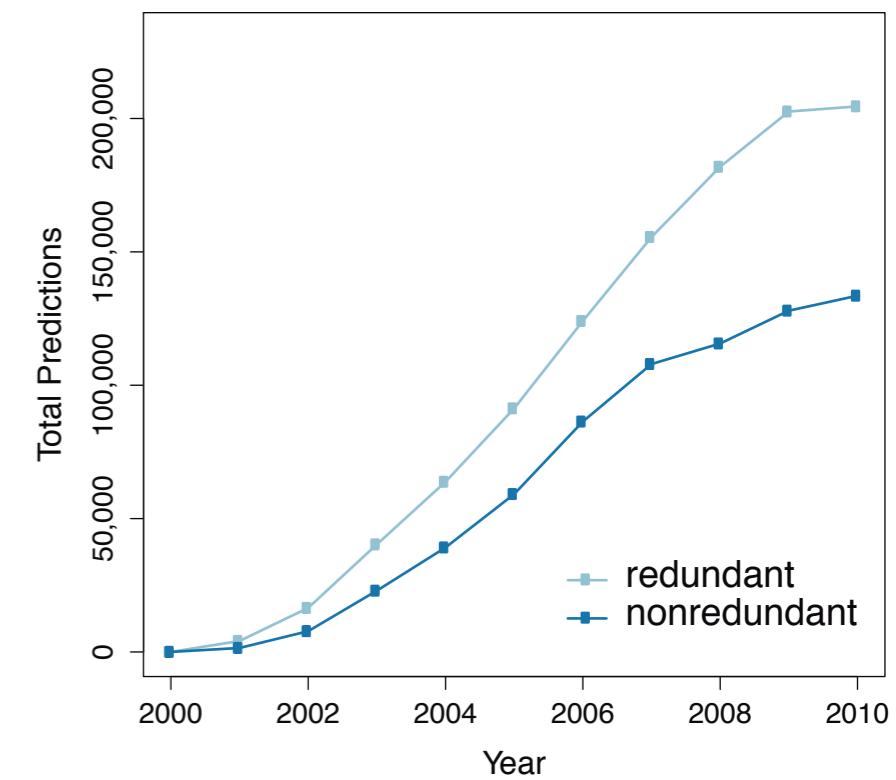
All



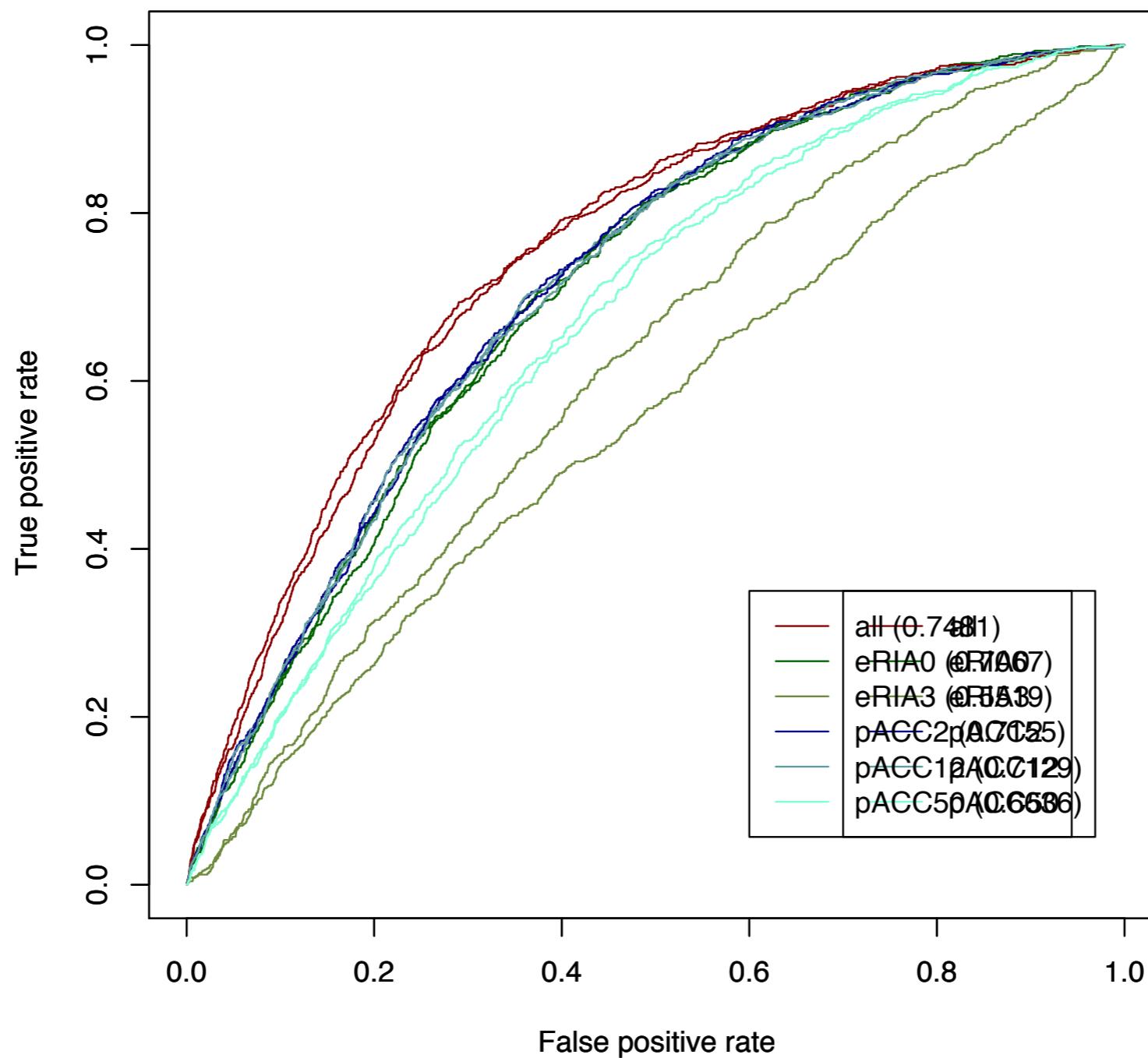
$n > 15$



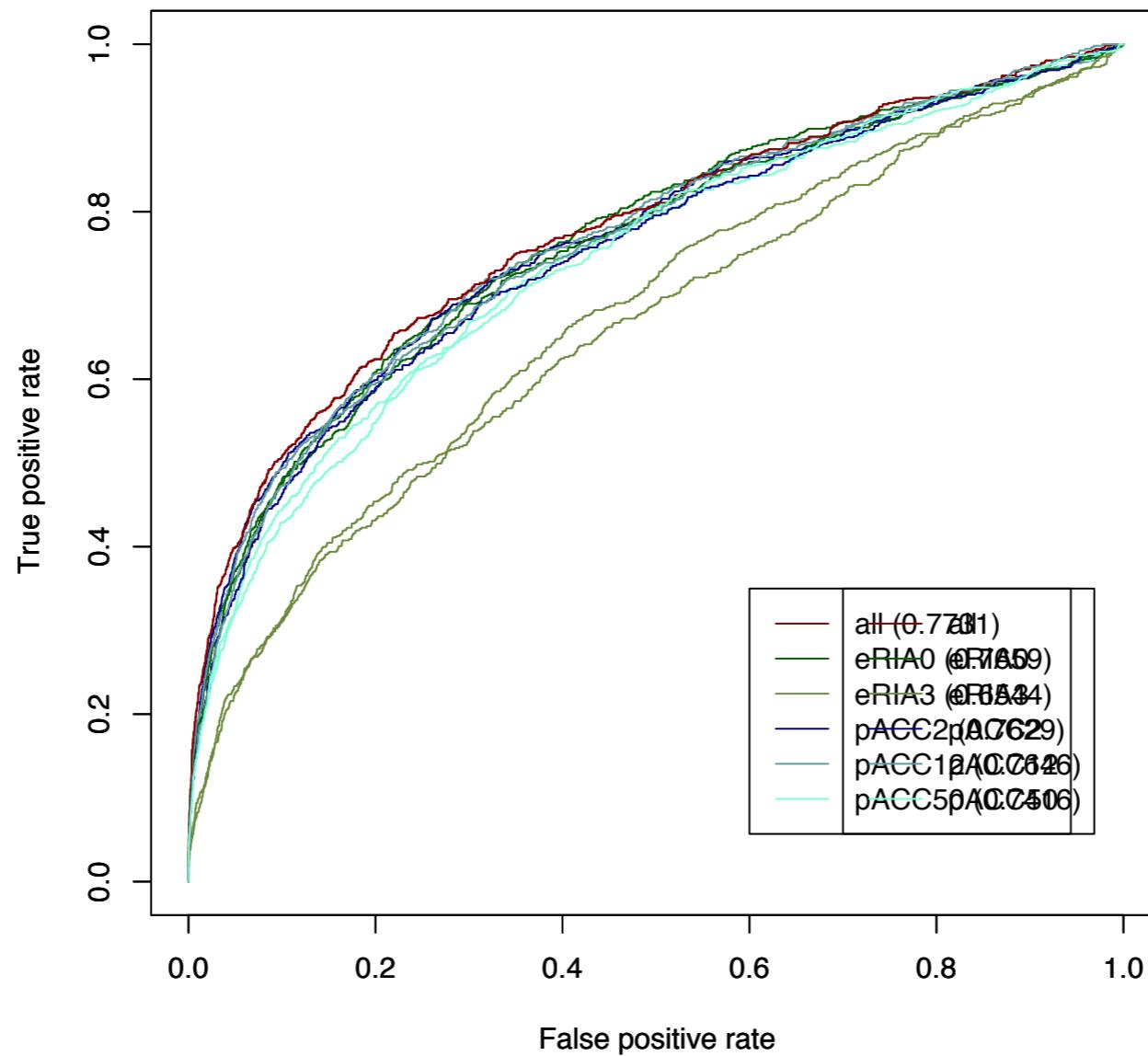
$n > 30$



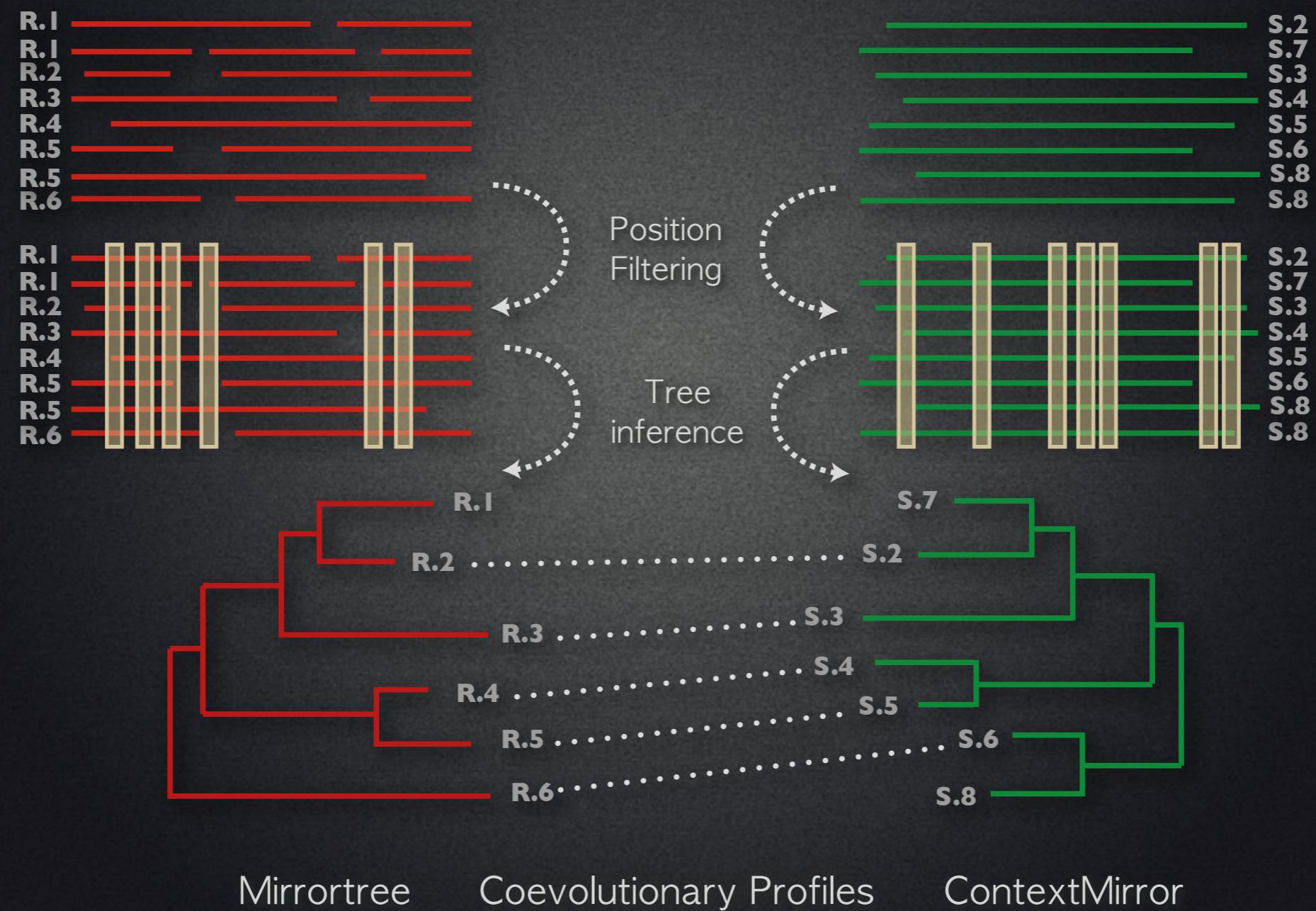
mt_COMPLEX



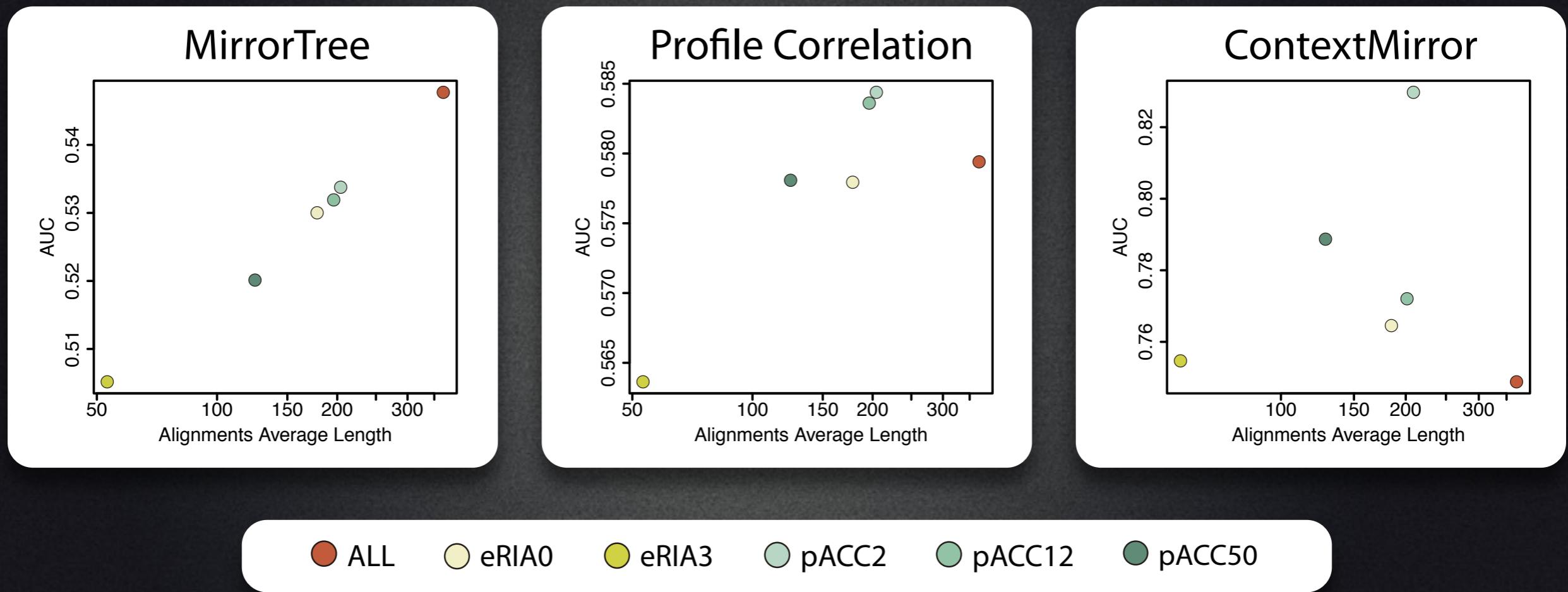
profiles_COMPLEX



Adapted workflow



Complexes



Ochoa, D., García-Gutiérrez, P., Juan, D., Valencia, A. & Pazos, F. Incorporating information on predicted solvent accessibility to the co-evolution-based study of protein interactions. Molecular bioSystems 9, 70–76 (2013).

Distances matrix

	S2	S3	S4	S5	S6	S7	S8
S2							
S3							
S4							
S5							
S6							
S7							
S8							

	S2	S3	S4	S5	S6
S2					
S3					
S4					
S5					
S6					

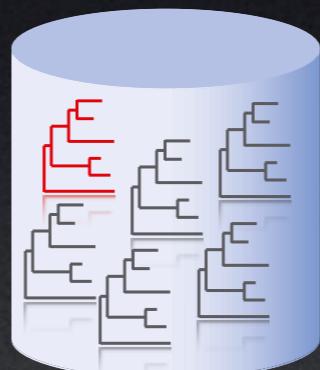
	S2	S3	S4	S5	S6
S2					
S3					
S4					
S5					
S6					

Distances matrix

	R1	R2	R3	R4	R5	R6
R1						
R2						
R3						
R4						
R5						
R6						

	R2	R3	R4	R5	R6
R2					
R3					
R4					
R5					
R6					

	R2	R3	R4	R5	R6
R2					
R3					
R4					
R5					
R6					

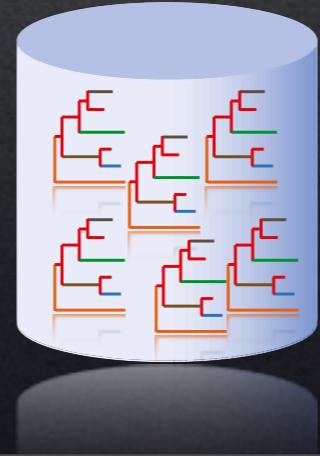


E.coli

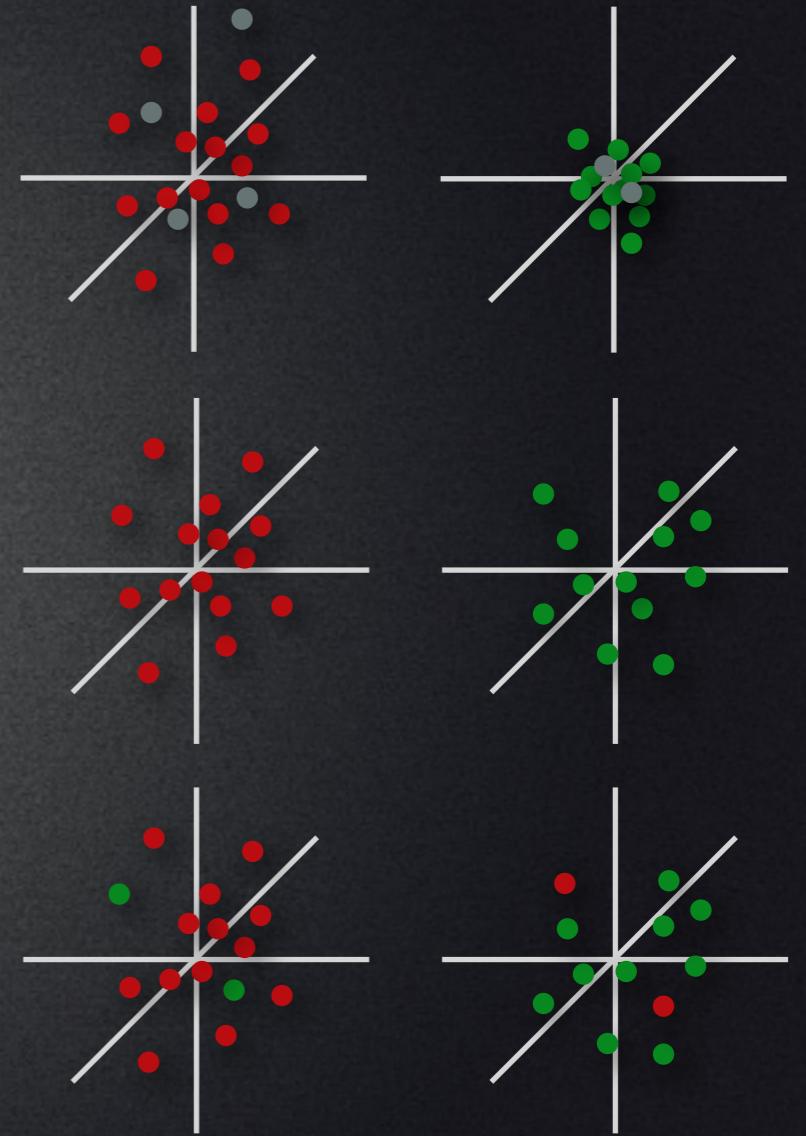
Randomly
select pair

Normalization

Switch
distances



Null E.coli



S5							
S6							
S7							
S8							

R4							
R5							
R6							

	S2	S3	S4	S5	S6		
S2							
S3							
S4							
S5							
S6							

	R2	R3	R4	R5	R6		
R2							
R3							
R4							
R5							
R6							

	S2	S3	S4	S5	S6		
S2							
S3							
S4							
S5							
S6							

	R2	R3	R4	R5	R6		
R2							
R3							
R4							
R5							
R6							

	S2	S3	S4	S5	S6		
S2							
S3							
S4							
S5							
S6							

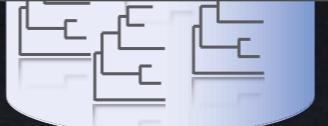
	R2	R3	R4	R5	R6		
R2							
R3							
R4							
R5							
R6							

	S2	S3	S4	S5	S6	S7	S8
S2							
S3							
S4							
S5							
S6							
S7							
S8							

	R1	R2	R3	R4	R5	R6	
R1							
R2							
R3							
R4							
R5							
R6							

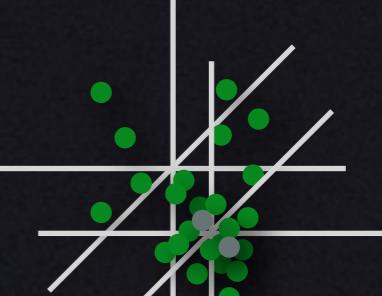
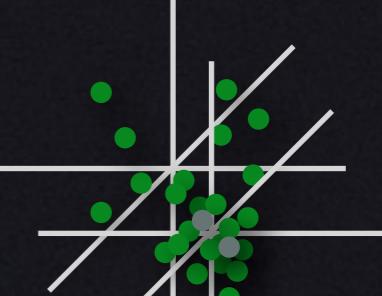
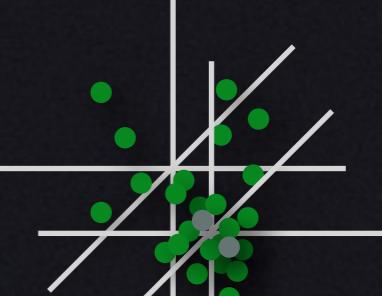
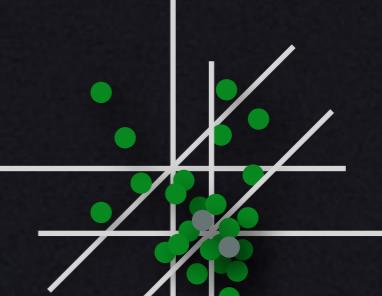
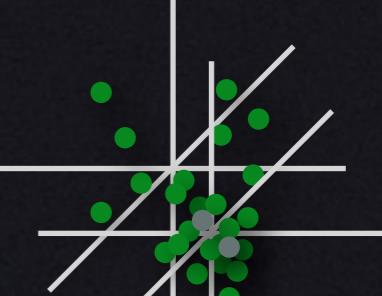
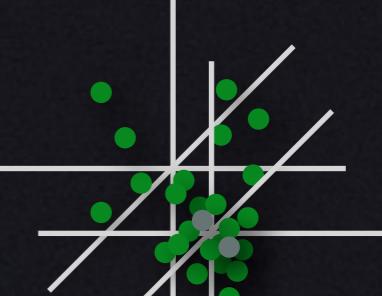
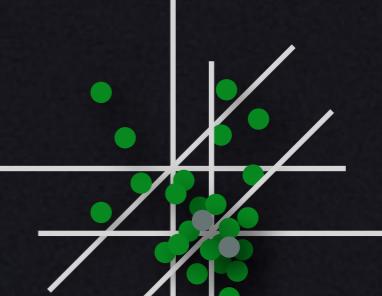
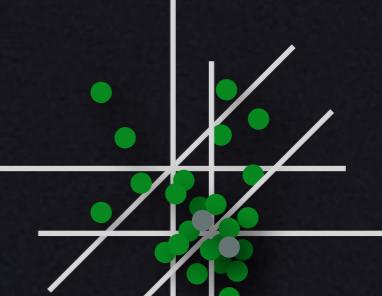
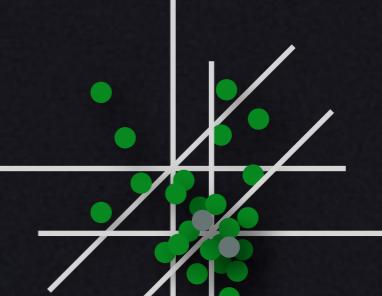
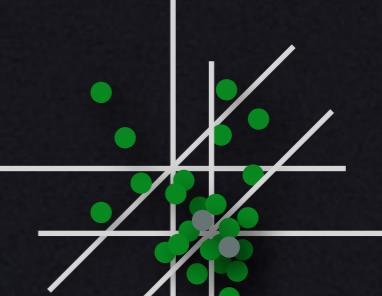
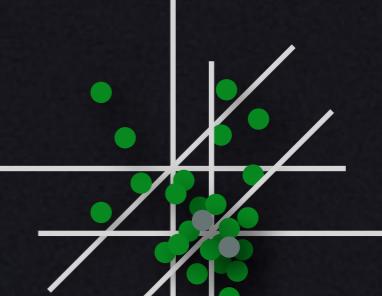
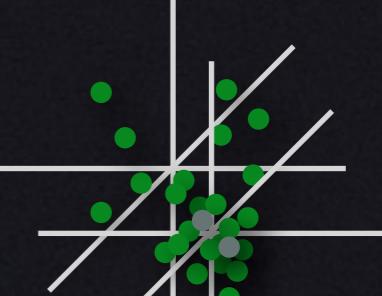
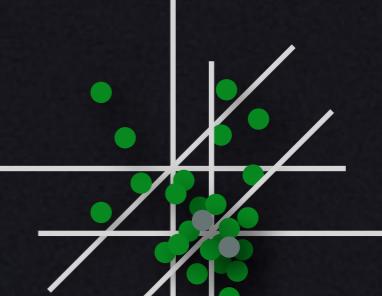
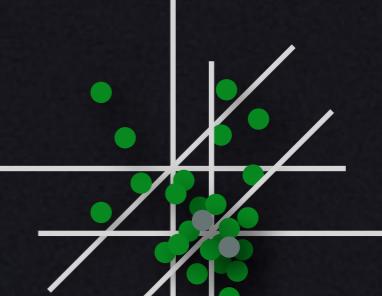
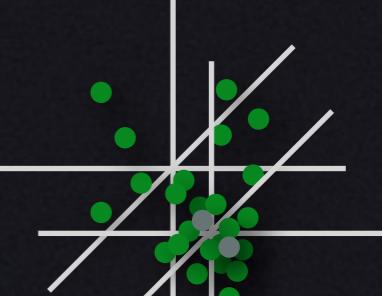
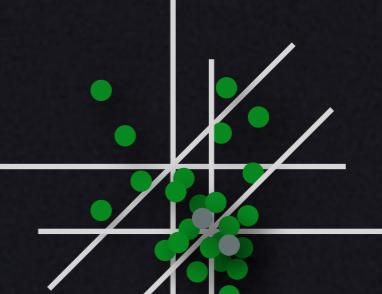
	S2	S3	S4	S5	S6	S7	S8
S2							
S3							
S4							
S5							
S6							
S7							
S8							

	R1	R2	R3	R4	R5	R6	
R1							
R2							
R3							
R4							
R5							
R6							

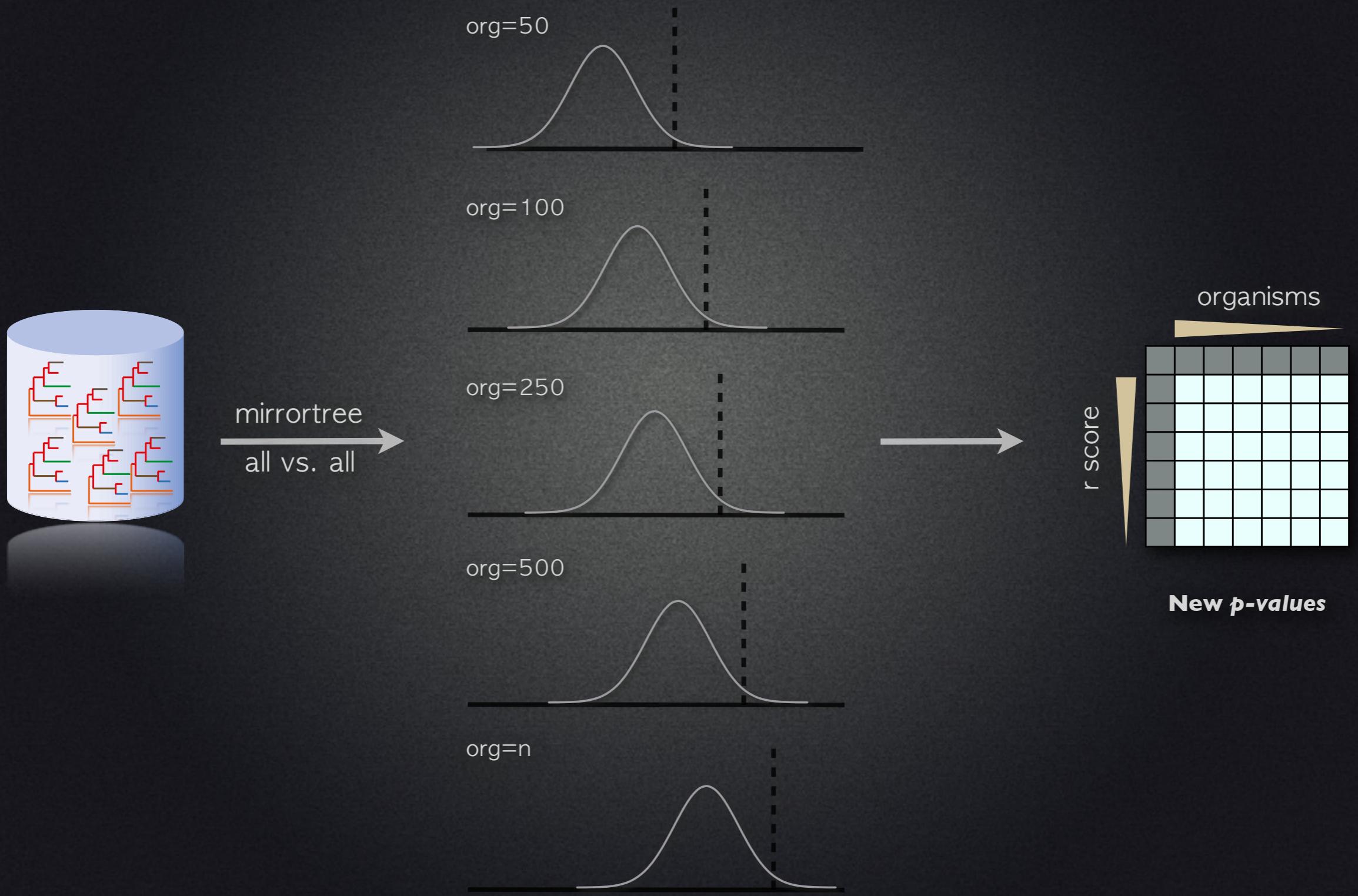


Randomly select pair

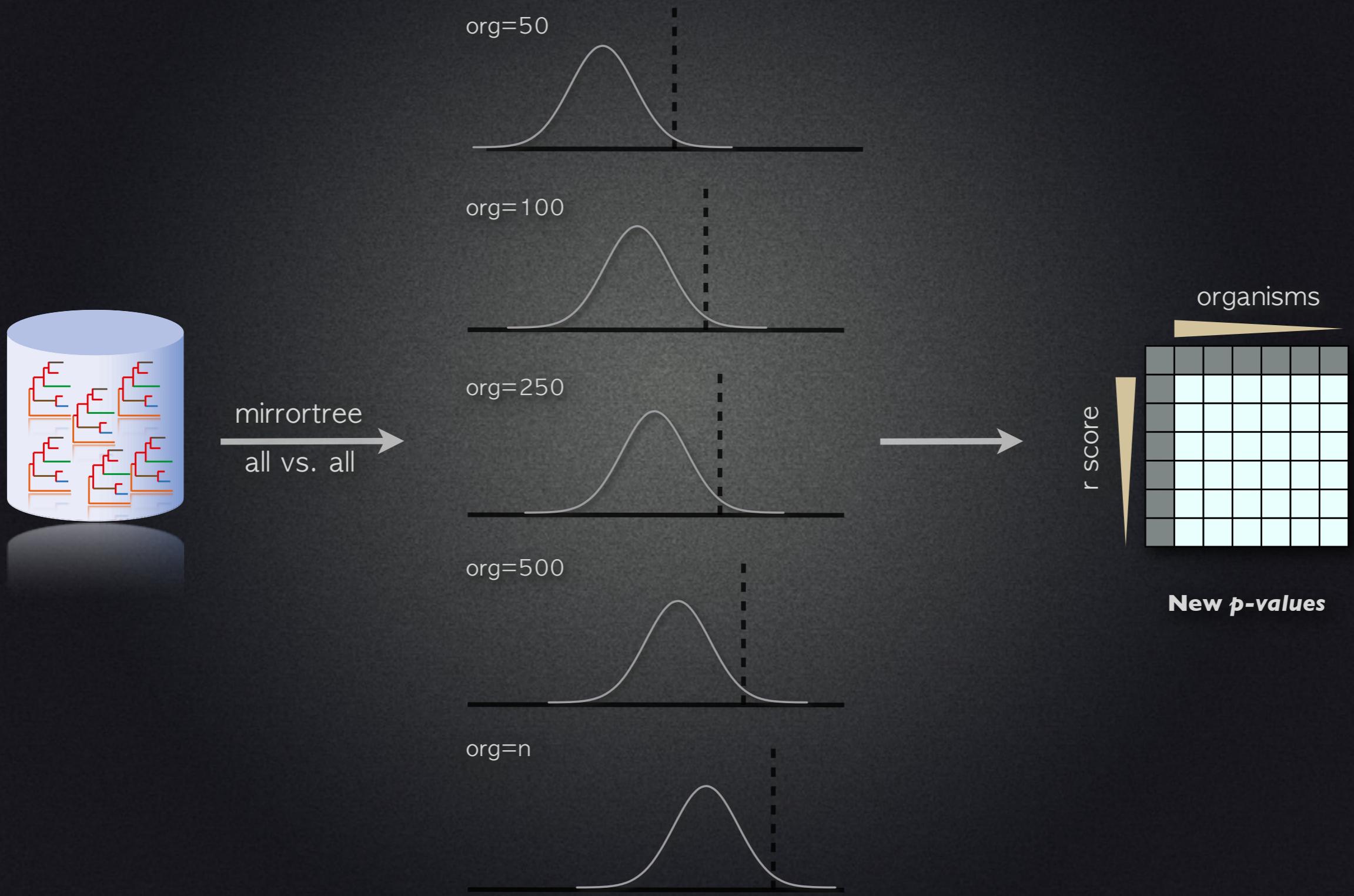
Standardization



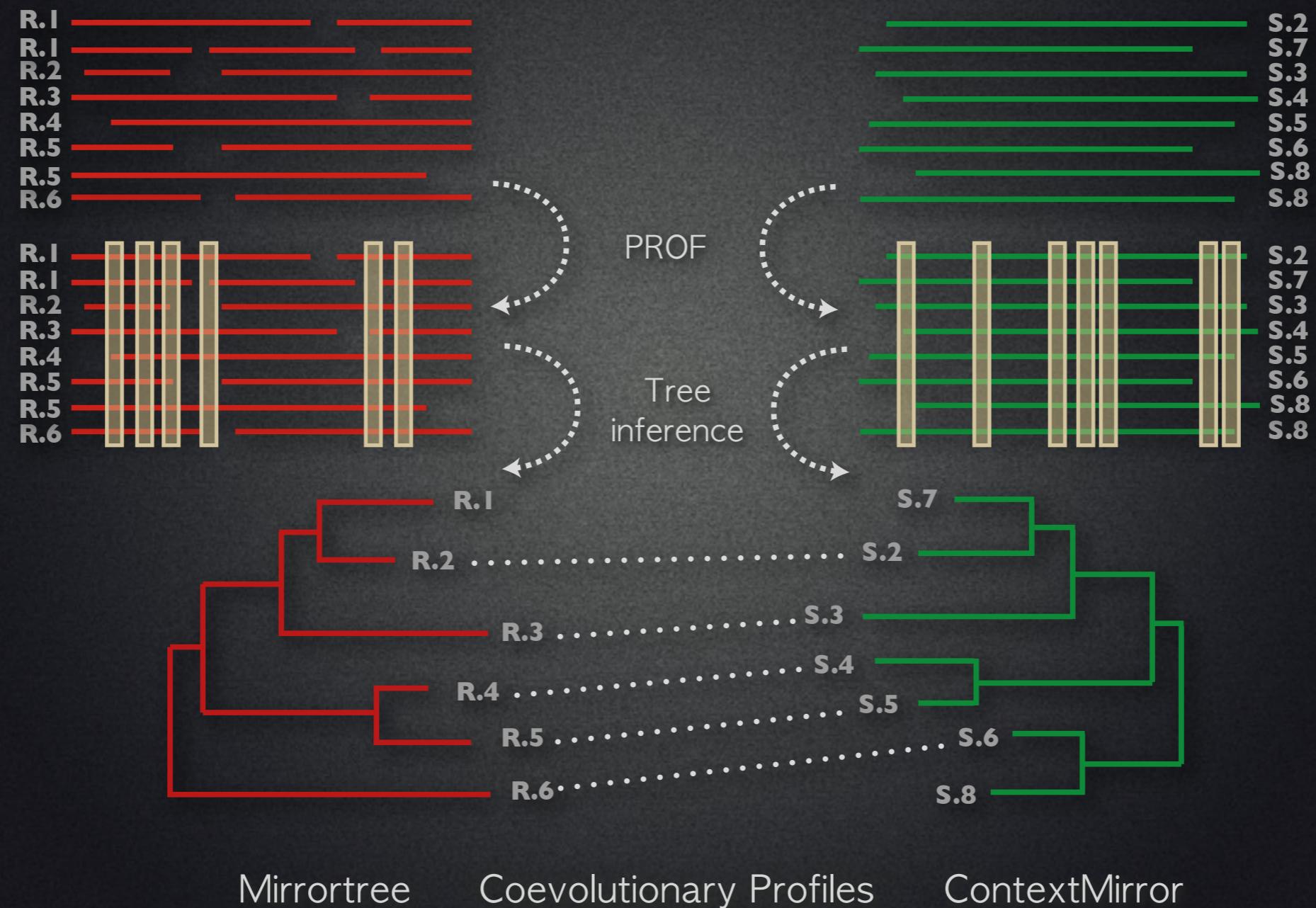
Corrected by interaction age



null distribution



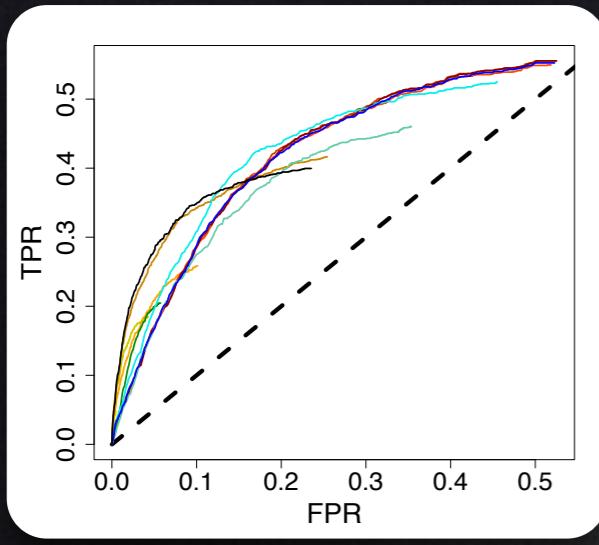
Adapted workflow



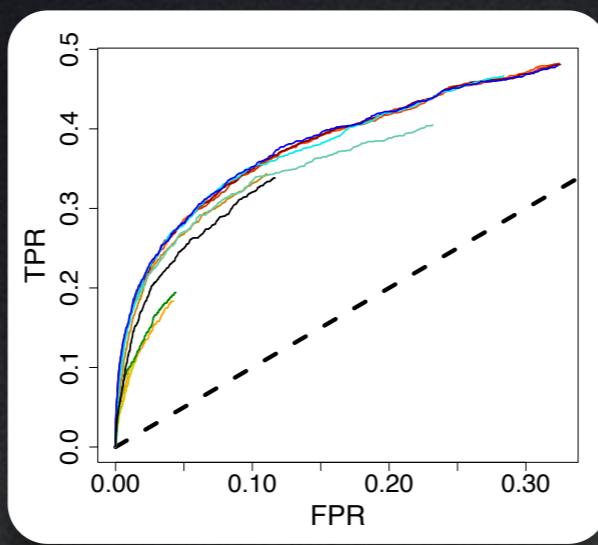
Rost, B. & Sander, C. Conservation and prediction of solvent accessibility in protein families. *Proteins* 20, 216–226 (1994).

Complexes

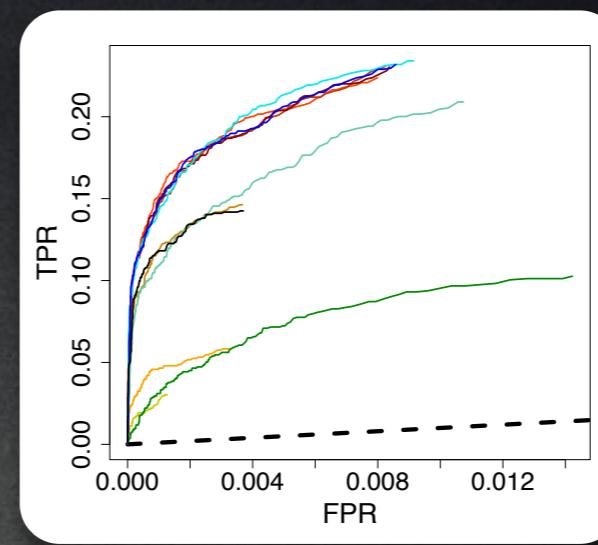
Mirrortree



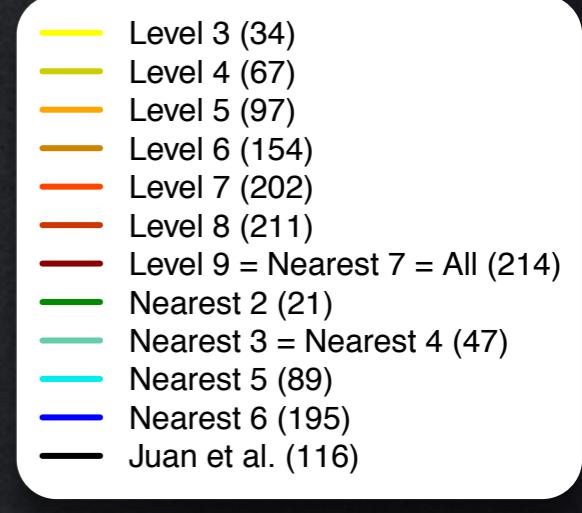
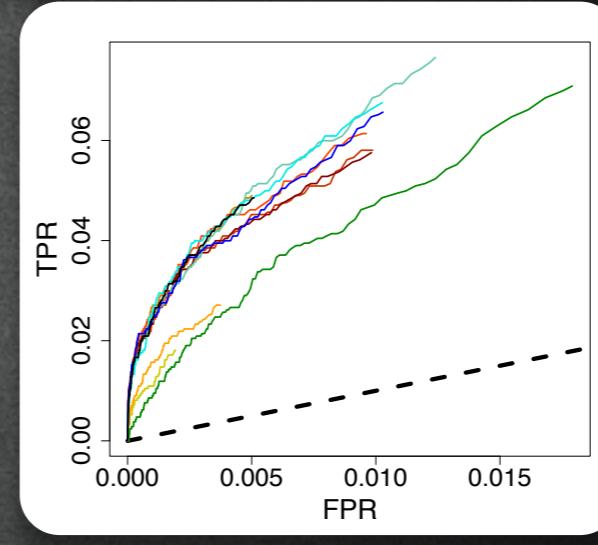
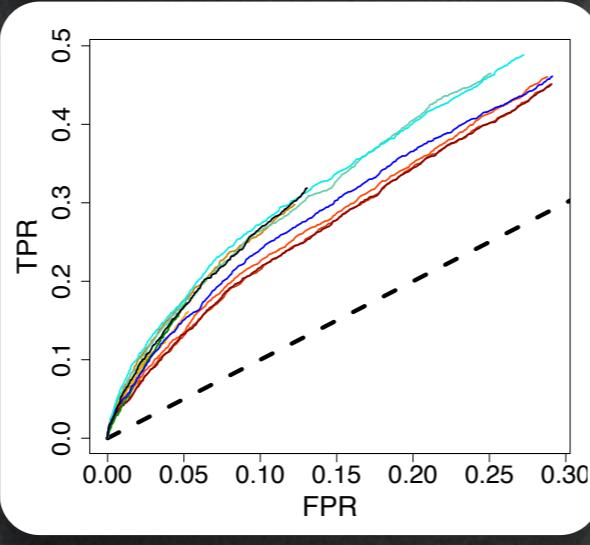
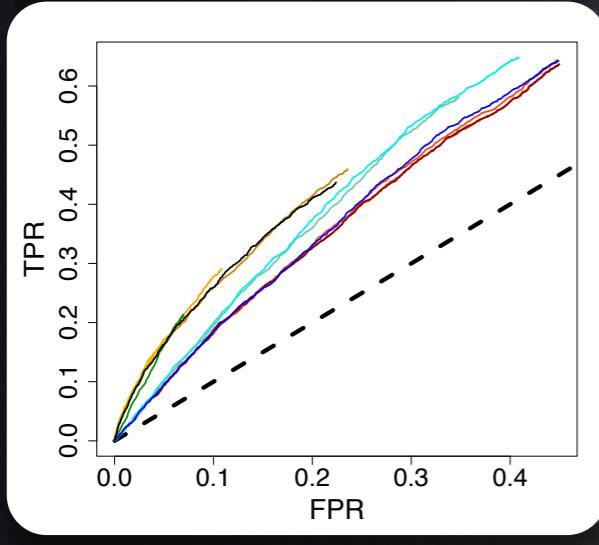
Profile Correlation



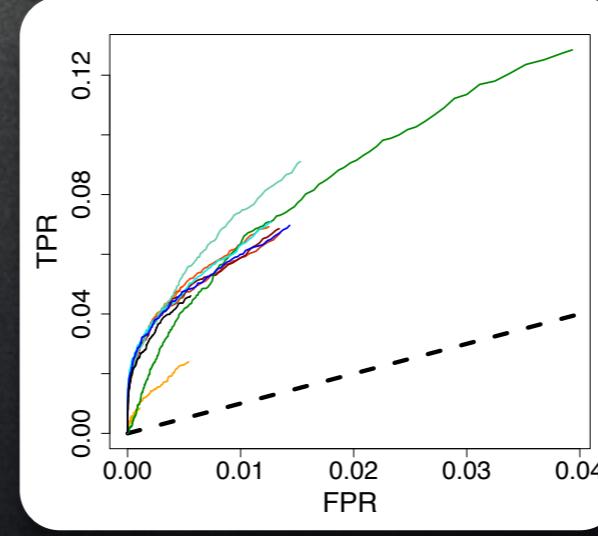
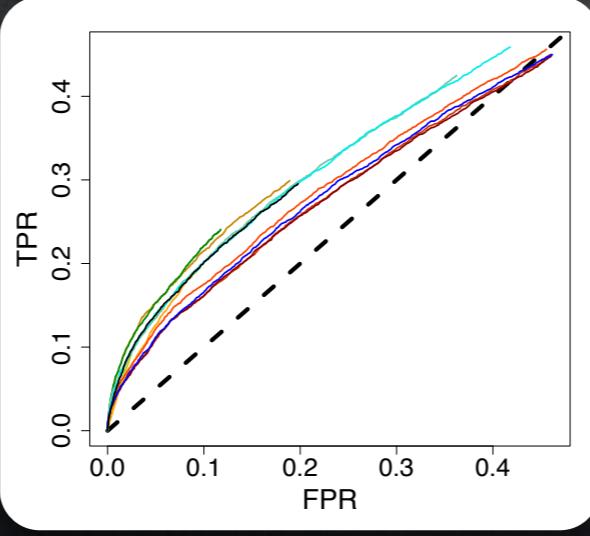
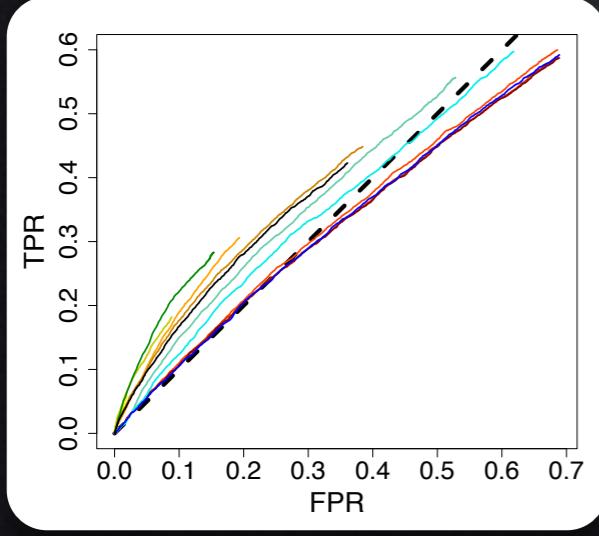
ContextMirror



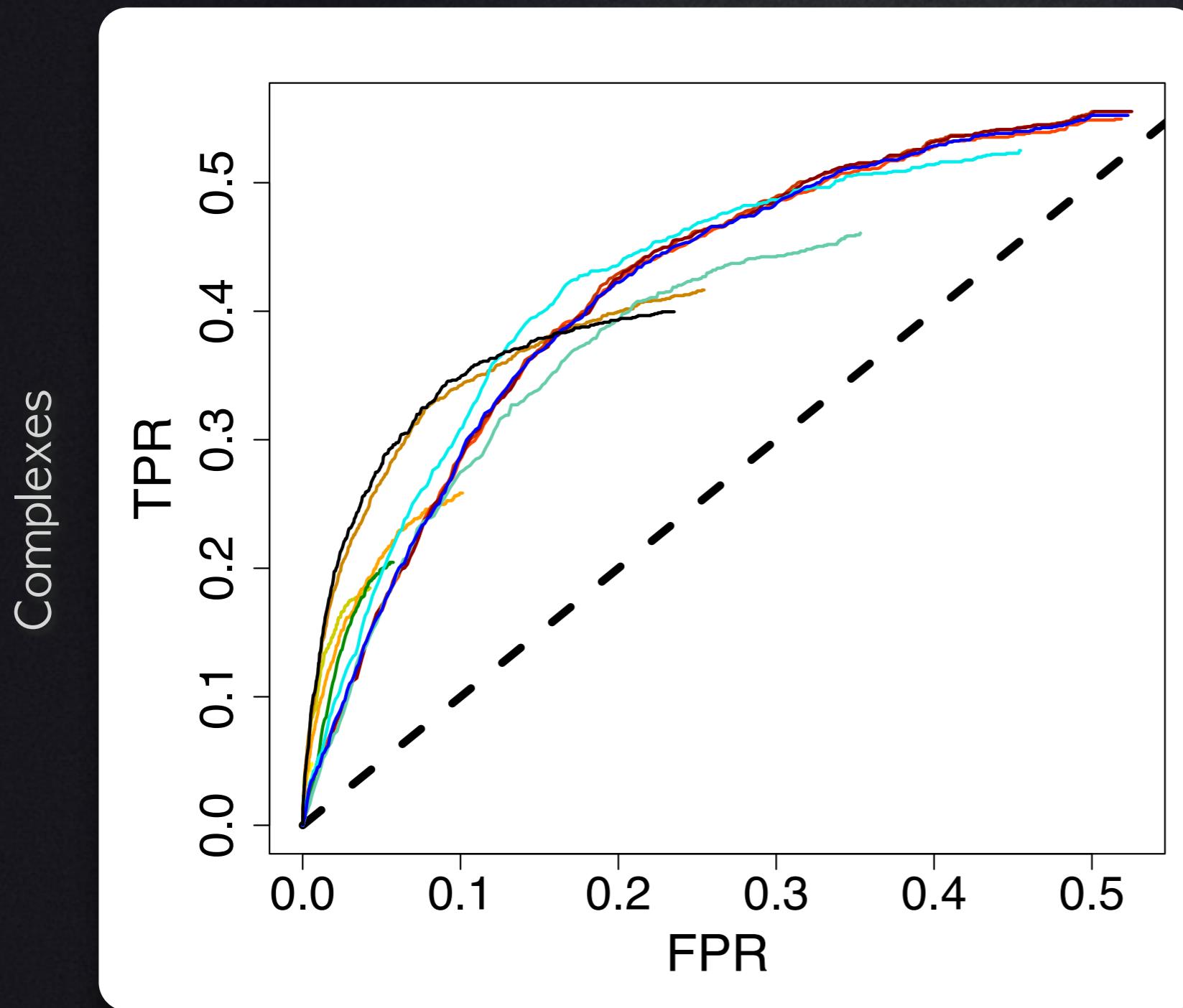
Binary Physical



Pathways

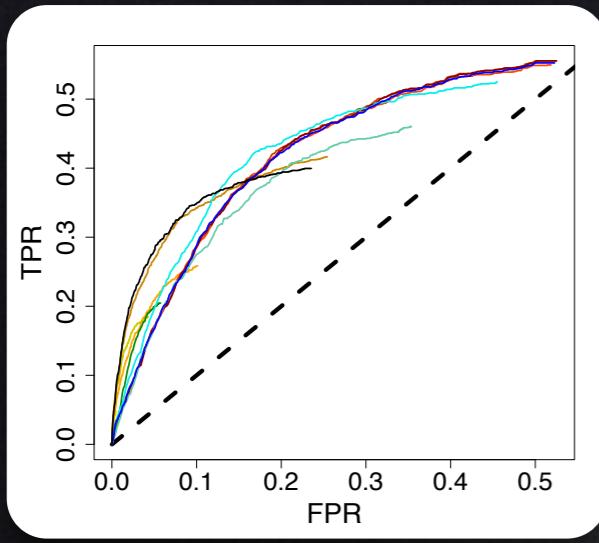


Mirrortree

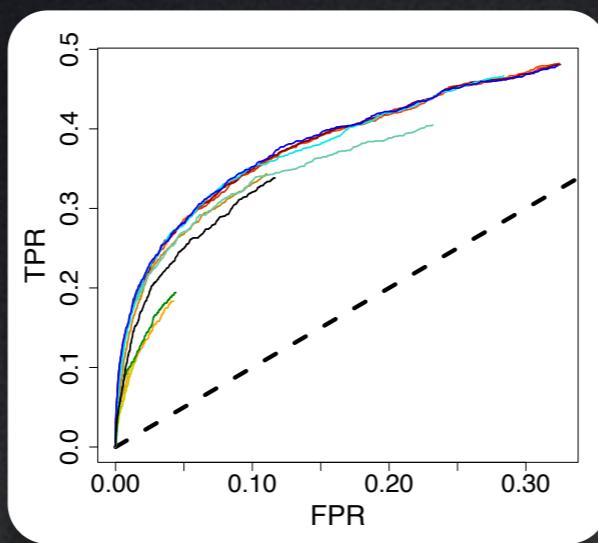


Complexes

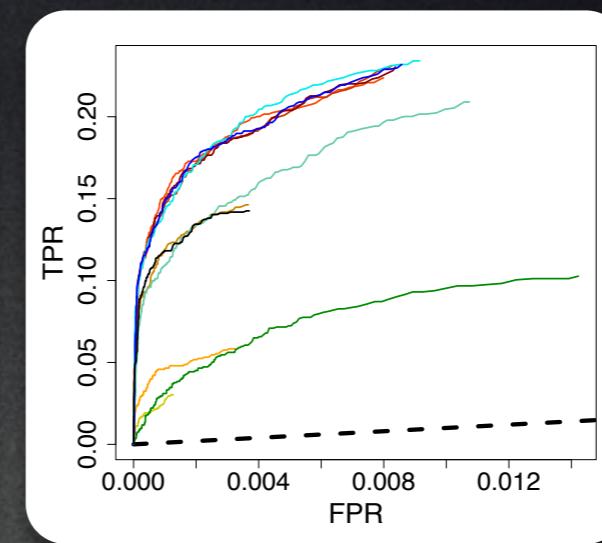
Mirrortree



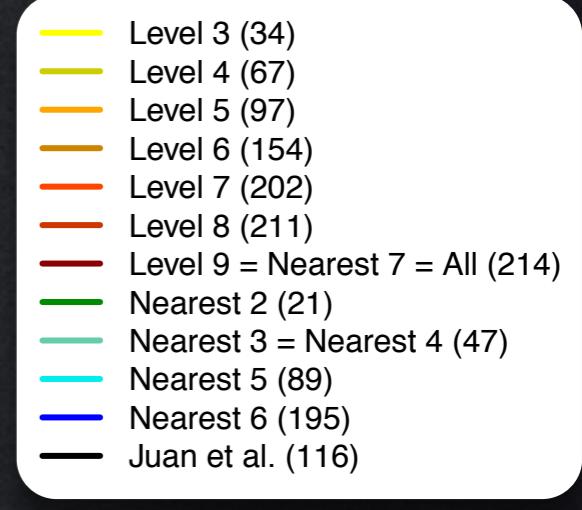
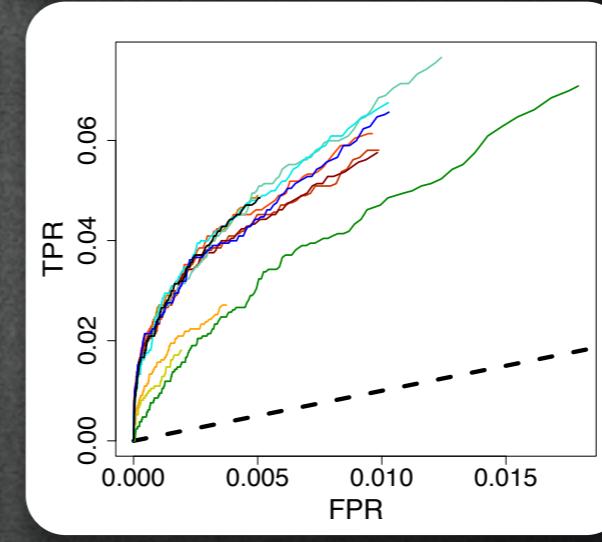
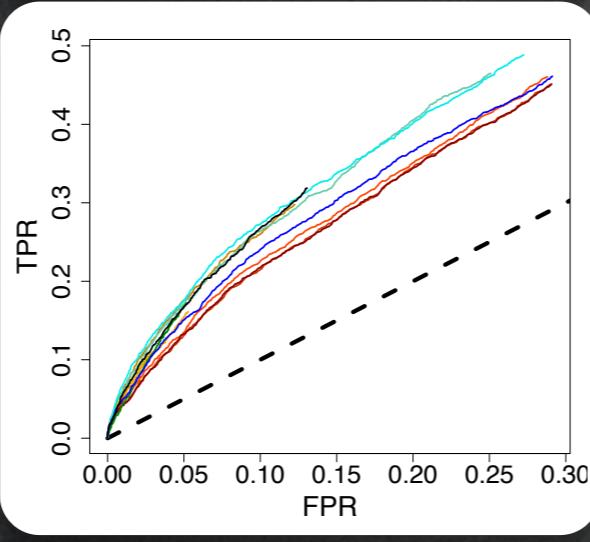
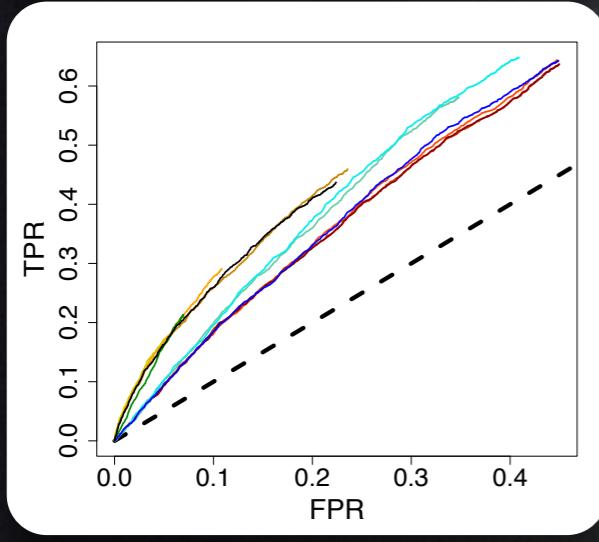
Profile Correlation



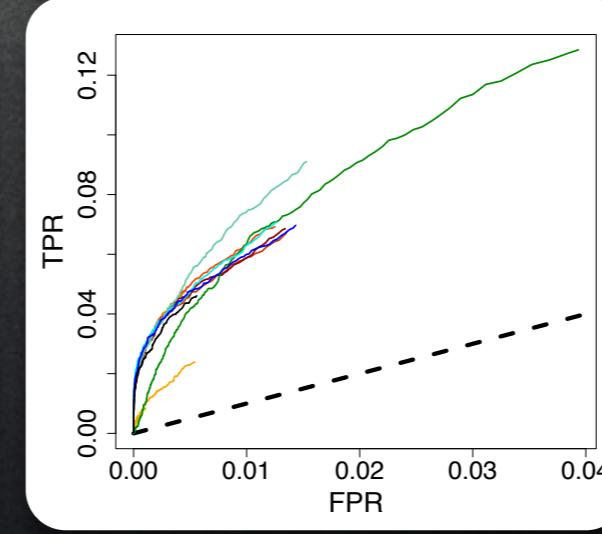
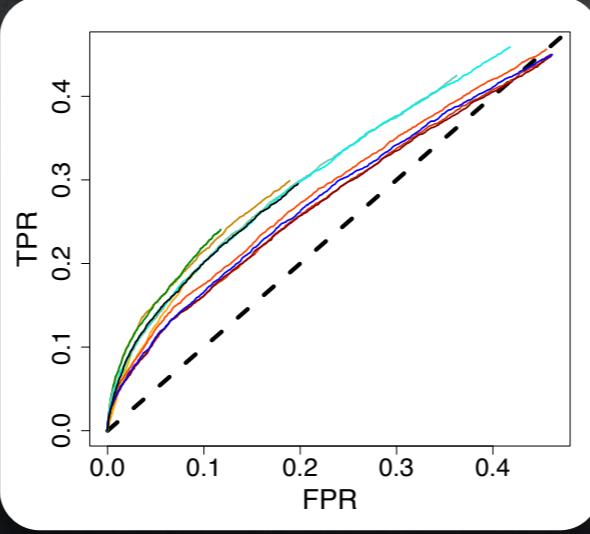
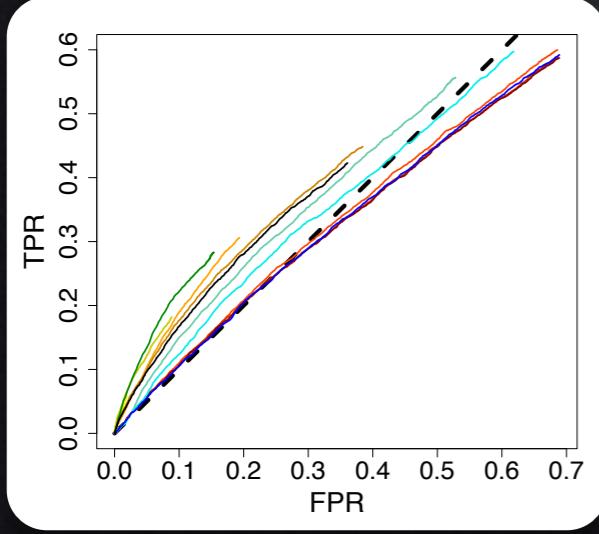
ContextMirror



Binary Physical

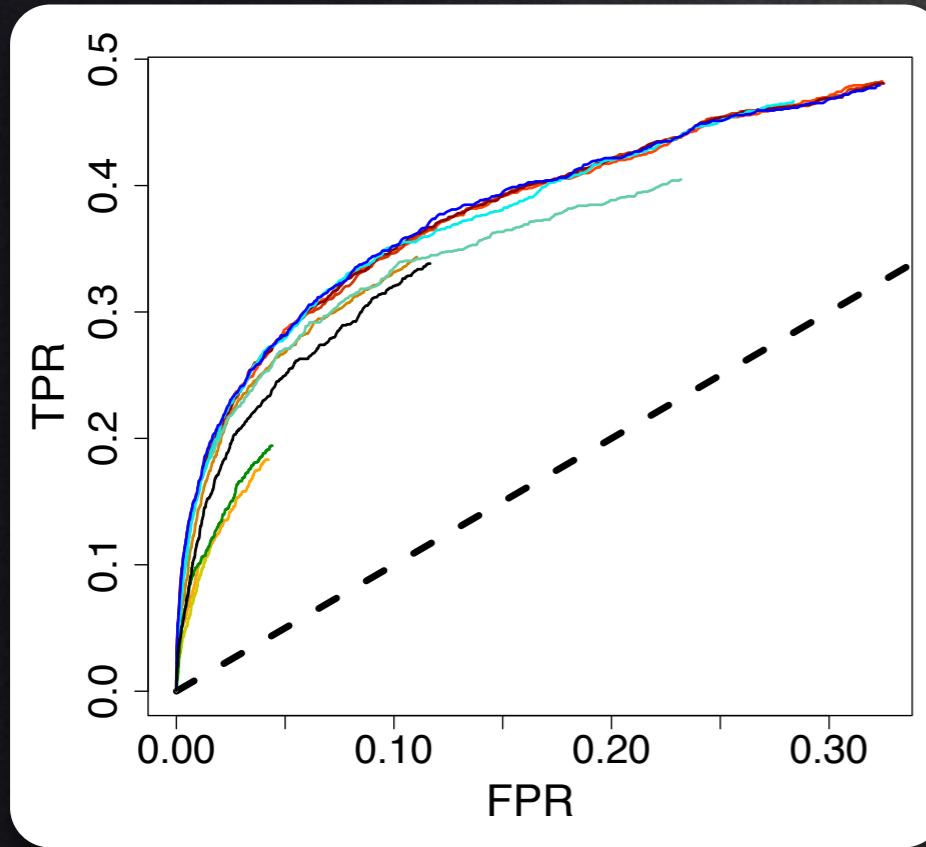


Pathways

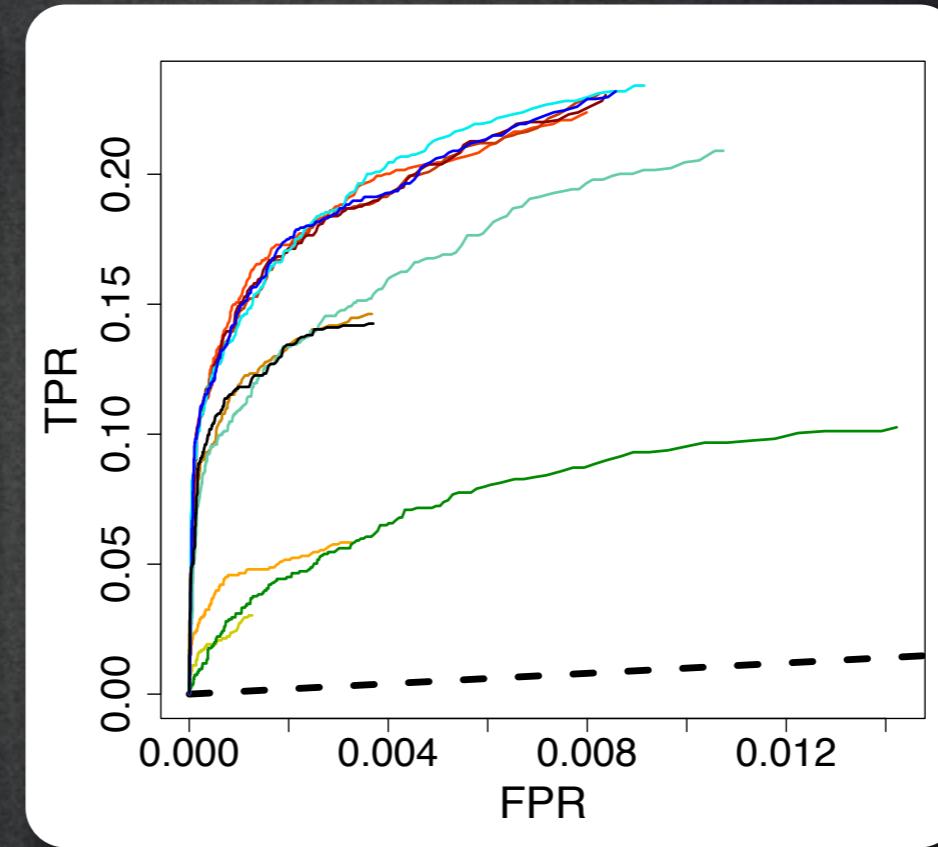


Complexes

Profile Correlation



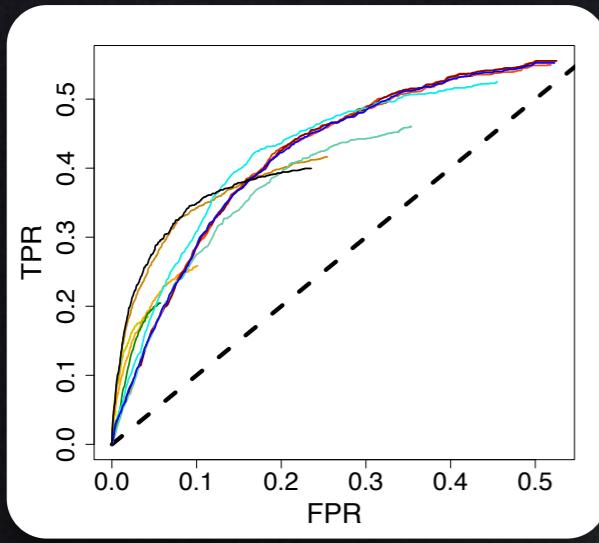
ContextMirror



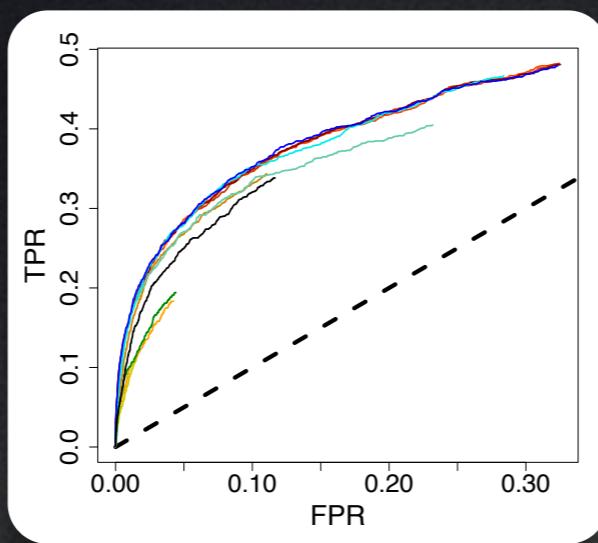
- Level 3 (34)
- Level 4 (67)
- Level 5 (97)
- Level 6 (154)
- Level 7 (202)
- Level 8 (211)
- Level 9 = Nearest 7 = All (214)
- Nearest 2 (21)
- Nearest 3 = Nearest 4 (47)
- Nearest 5 (89)
- Nearest 6 (195)
- Juan et al. (116)

Complexes

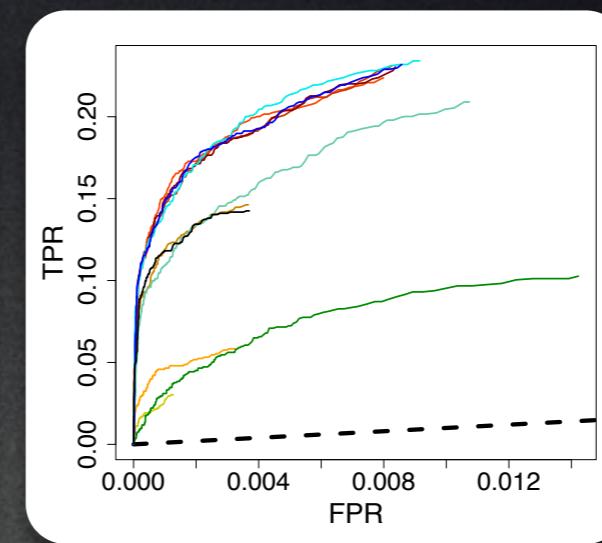
Mirrortree



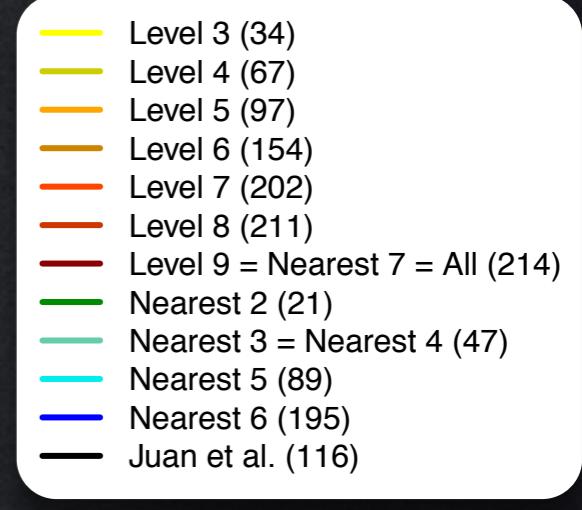
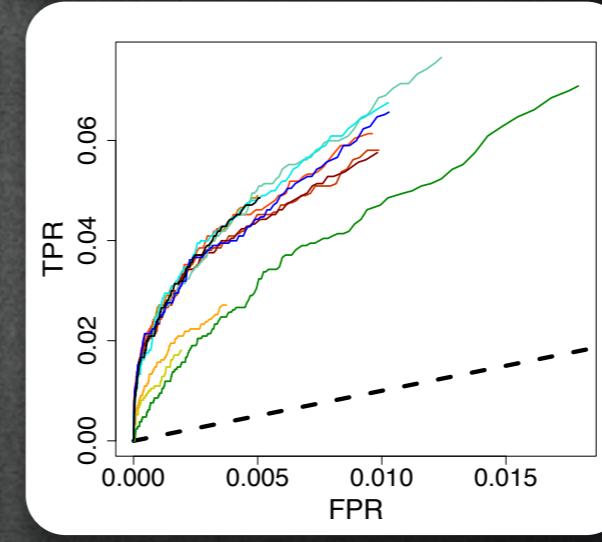
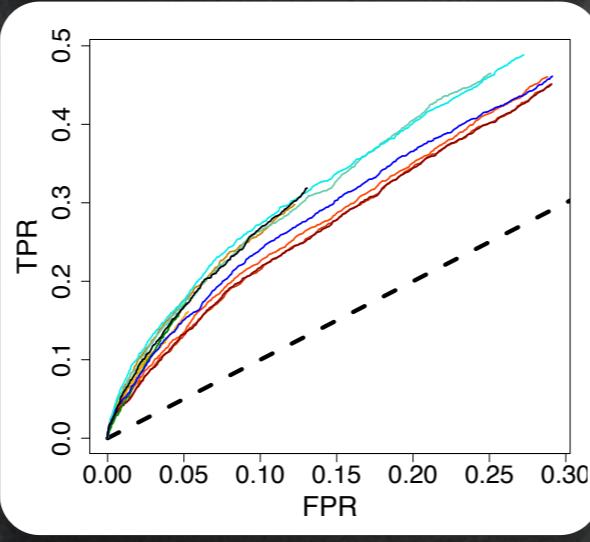
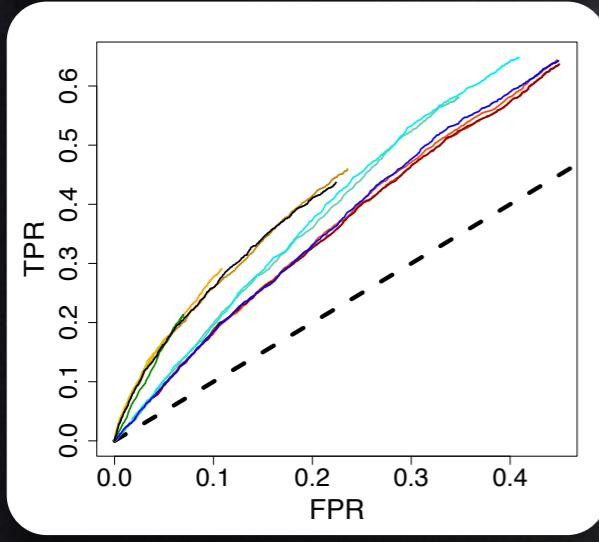
Profile Correlation



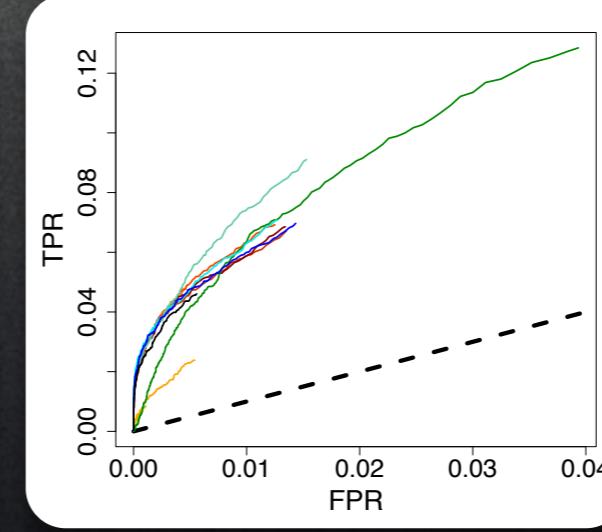
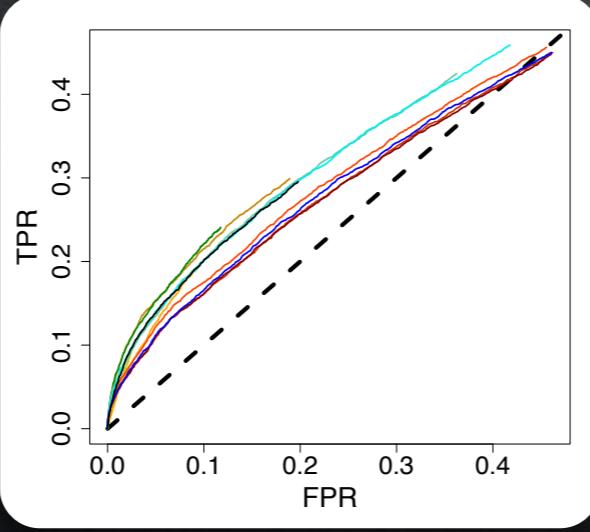
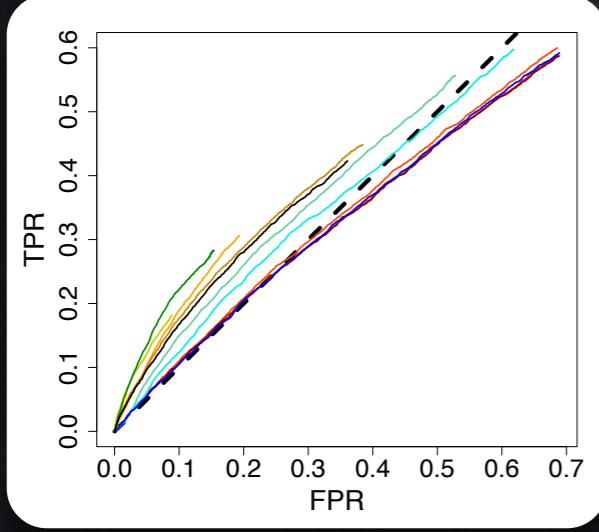
ContextMirror



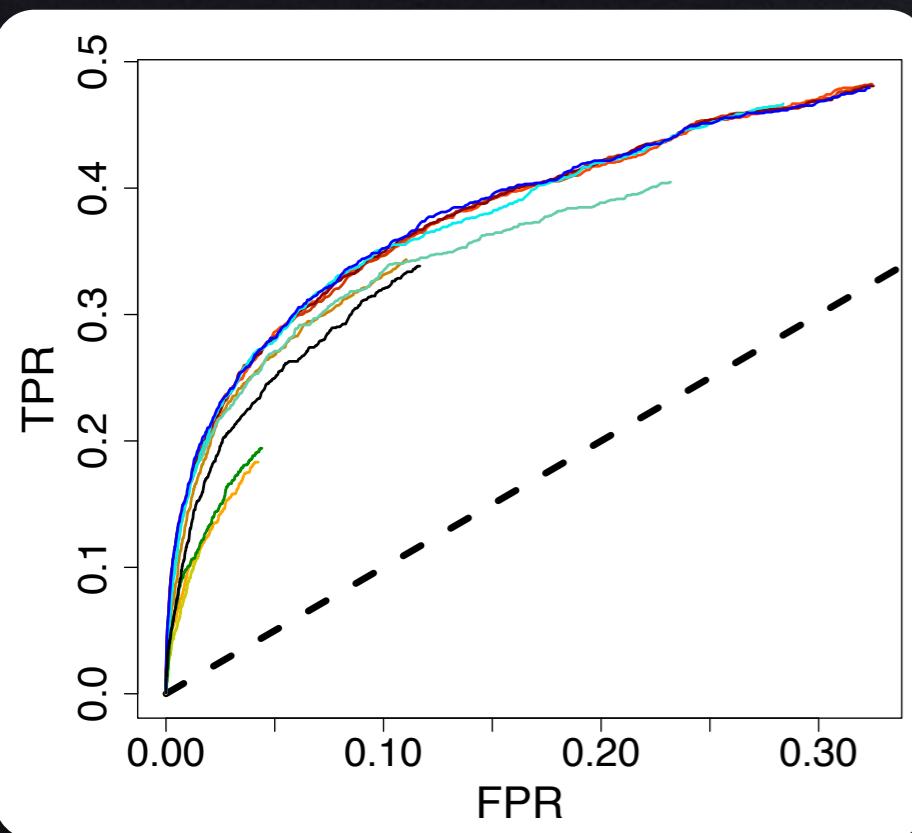
Binary Physical



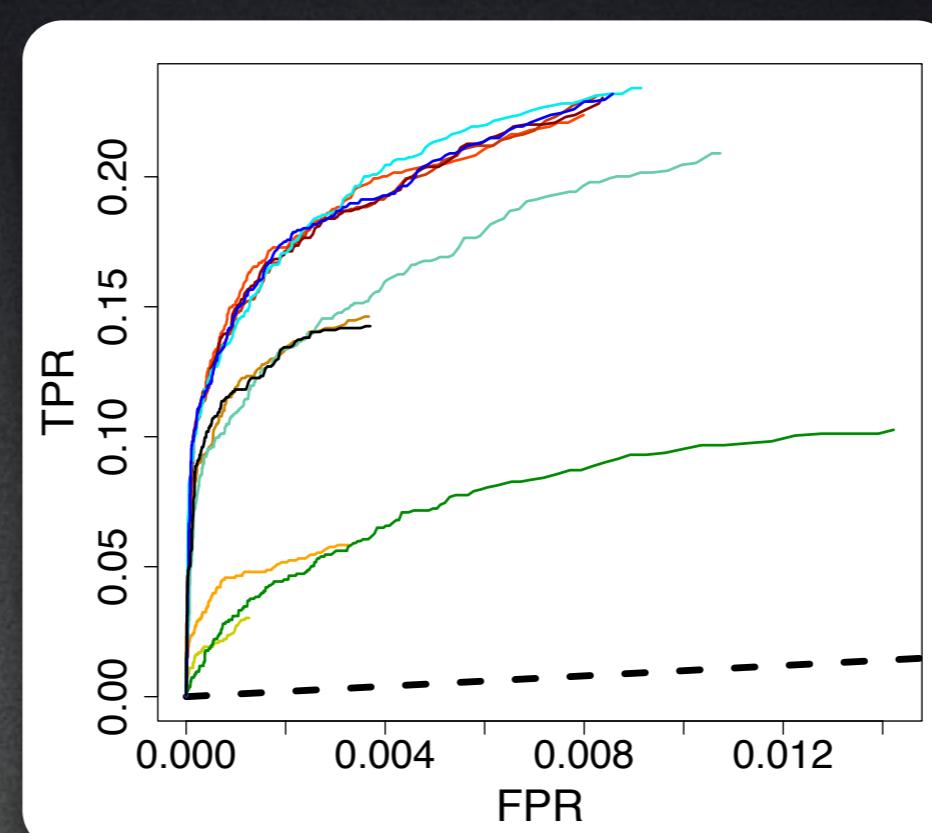
Pathways



Profile Correlation

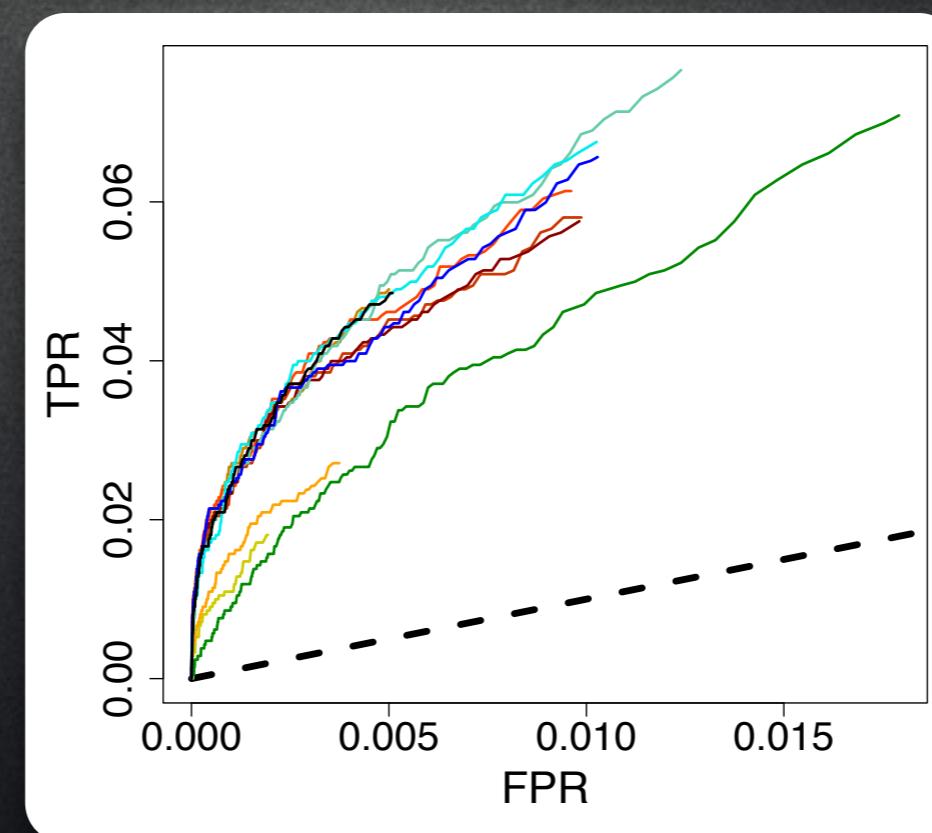
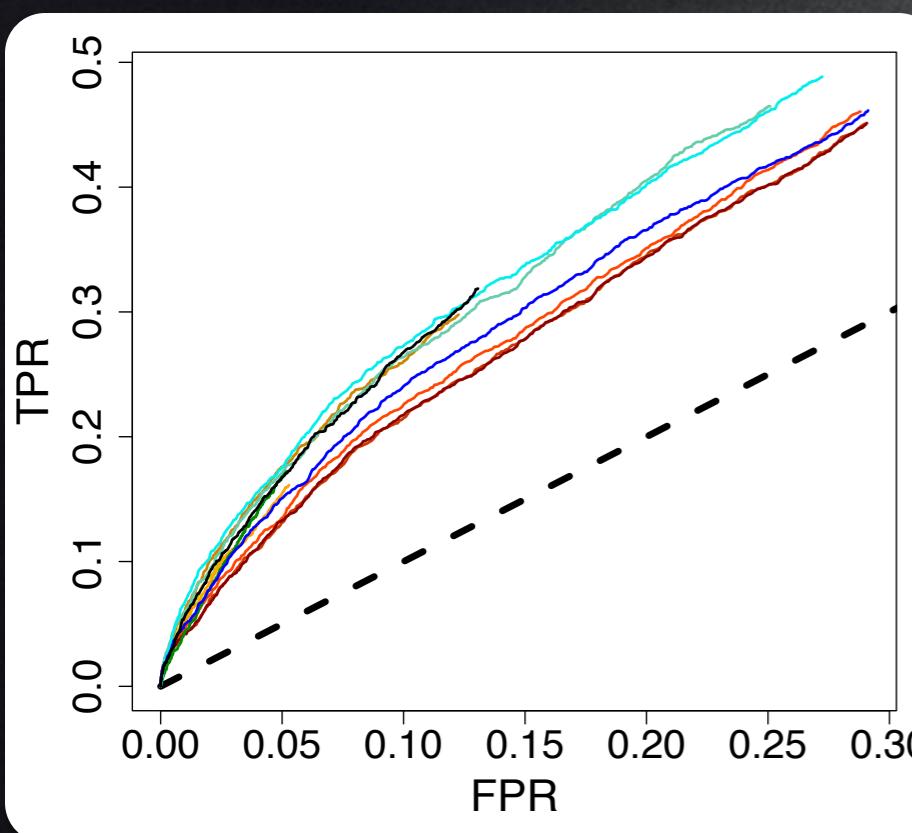


ContextMirror



Complexes

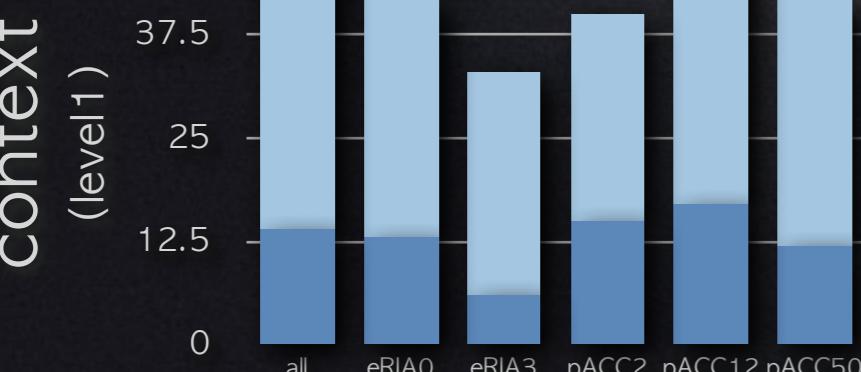
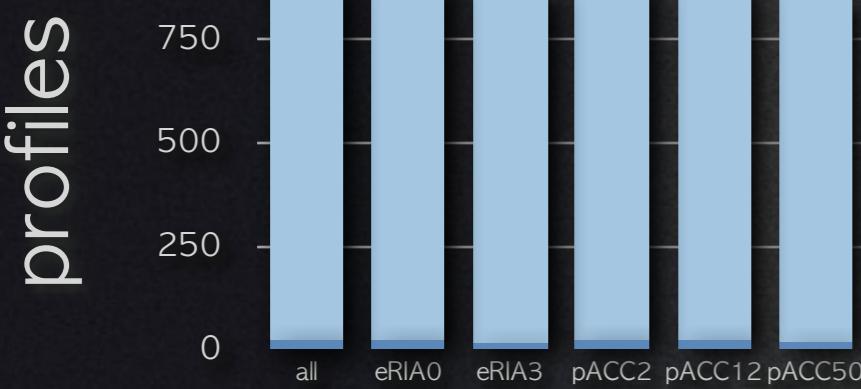
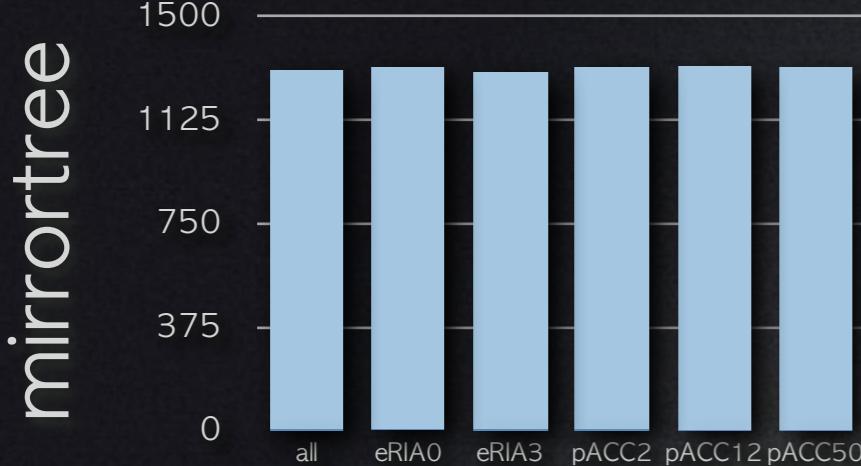
- Level 3 (34)
- Level 4 (67)
- Level 5 (97)
- Level 6 (154)
- Level 7 (202)
- Level 8 (211)
- Level 9 = Nearest 7 = All (214)
- Nearest 2 (21)
- Nearest 3 = Nearest 4 (47)
- Nearest 5 (89)
- Nearest 6 (195)
- Juan et al. (116)



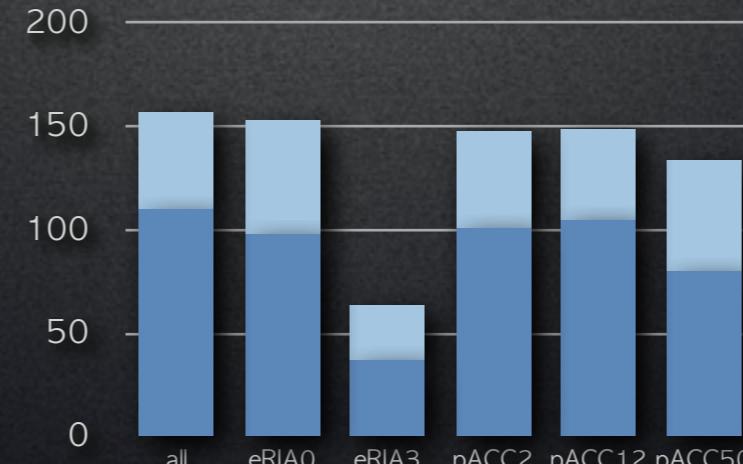
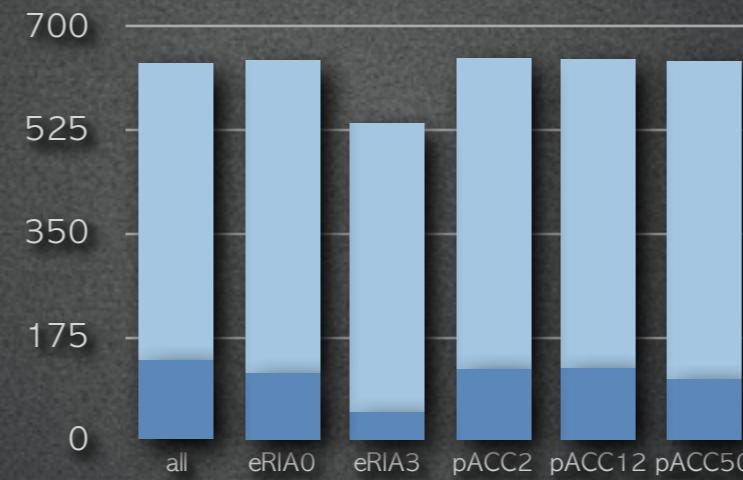
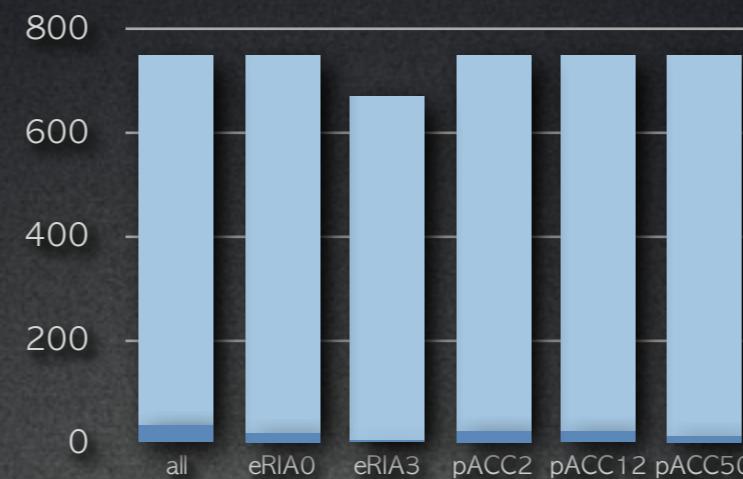
Binary Physical



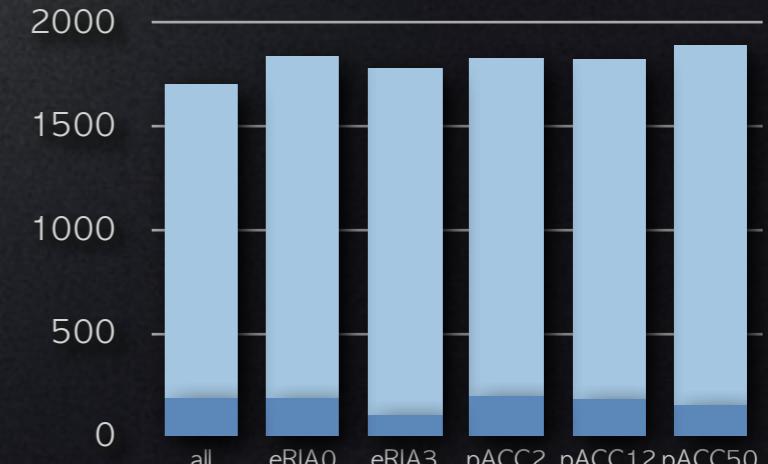
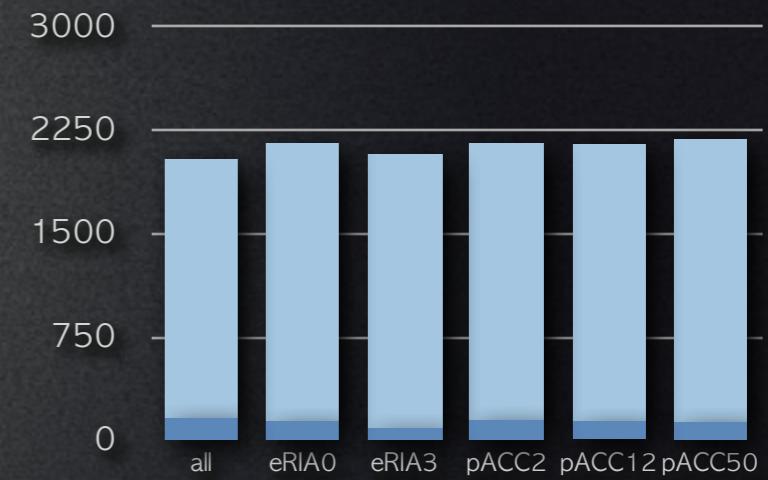
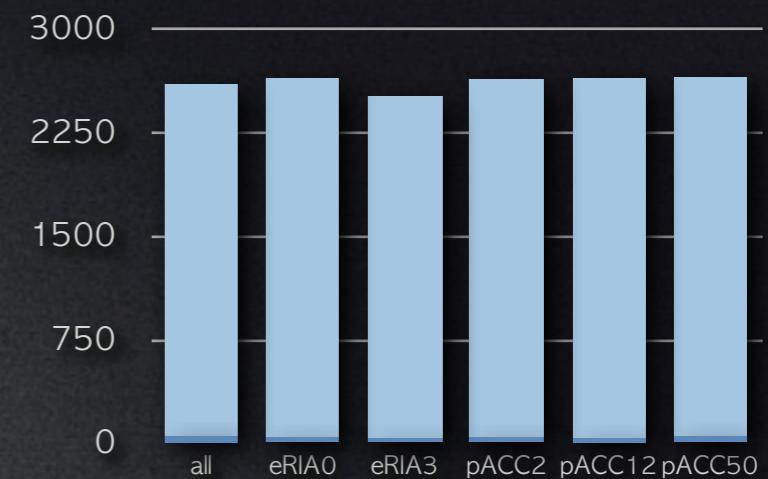
Binary Physical



Complexes



Pathways



ALL/215-275 SLGKALRQLPLTDEQKKRIPEDASTLHRL_LGAQPGSQRLRHAGNPLHLDVLVVDEASMID
eRIA0/103-130 --GK--RQ-PLTDEQK-R--ED-----G-QPGSQR-RHHAG-P-H-----
eRIA3/25-27 ---K--RQ-----
pACC2/113-148 --GK--RQLPLTDEQKKR--ED-----R--GAQPGSQRLRHAGNPLH-D-----
pACC12/107-138 --GK--RQ-PLTDEQKKR--ED-----R--GAQPGSQR-RHHAGNP-H-----
pACC50/58-78 ---K--RQ-PLT-EQK---E-----QPGSQR-R--AG-P-H-----

