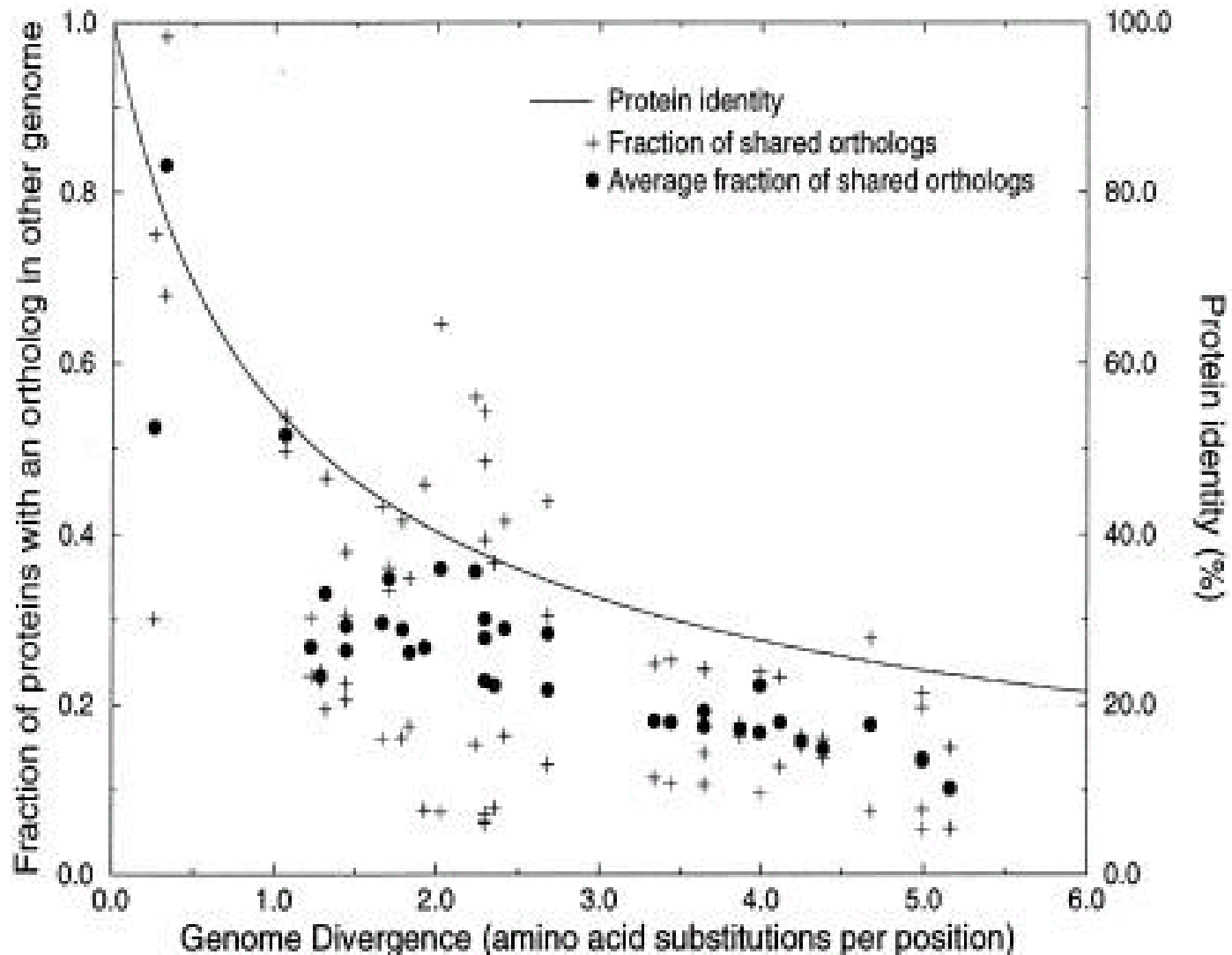
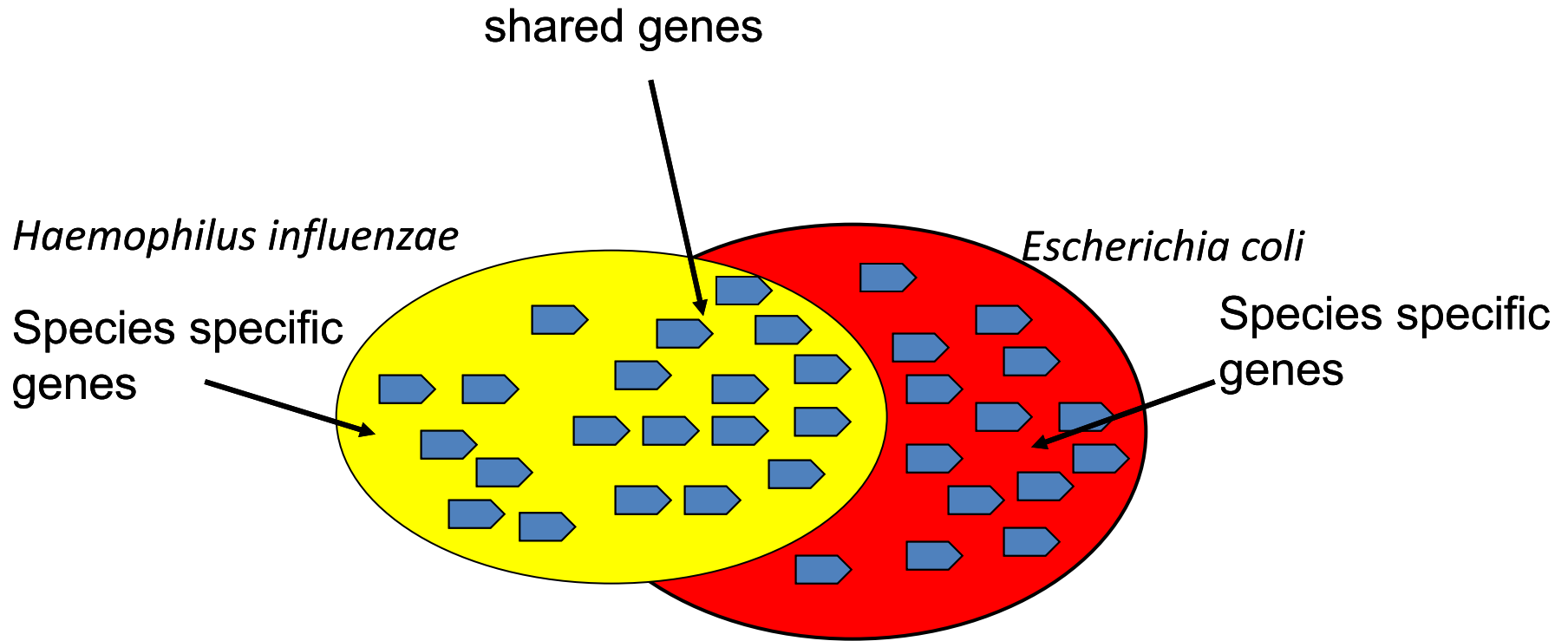


Bioinformatics and Evolutionary Genomics: Genome Evolution in terms of Gene Content

Gene Content Evolution



What about HGT / genome sizes?
Genome trees based on gene content:



Genome trees based on gene content

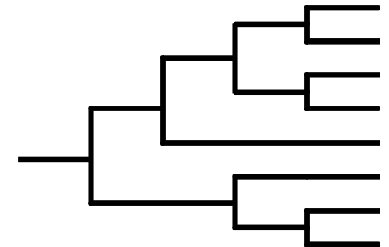
Presence / absence matrix:

	OG1	OG2	OG3	OG4	...
sp1	1	1	0	1	...
sp2	0	1	0	0	...
sp3	0	0	1	1	...
...

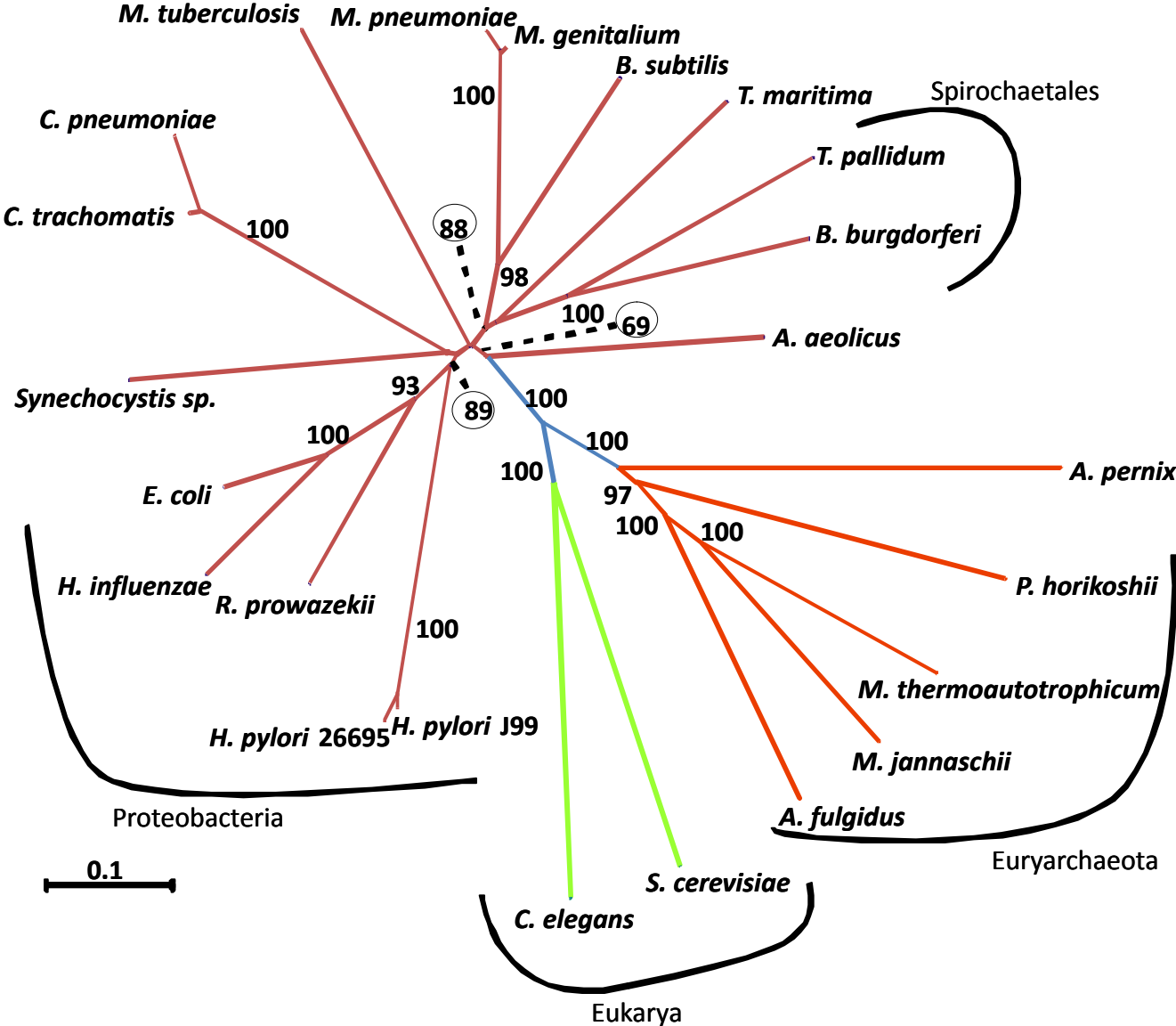
$$dist(spA, spB) = 1 - \left(\frac{\# \text{ shared OGs } (spA, spB)}{\text{Size of the smallest genome}} \right)$$

d\s	sp1	sp2	sp3	sp4	...
sp1	0\1	0.2	0.4	0.2	...
sp2	0.8	0\1	0.9	0.1	...
sp3	0.6	0.1	0\1	0.3	...
sp4	0.8	0.9	0.7	0\1	...
...

Neighbor joining

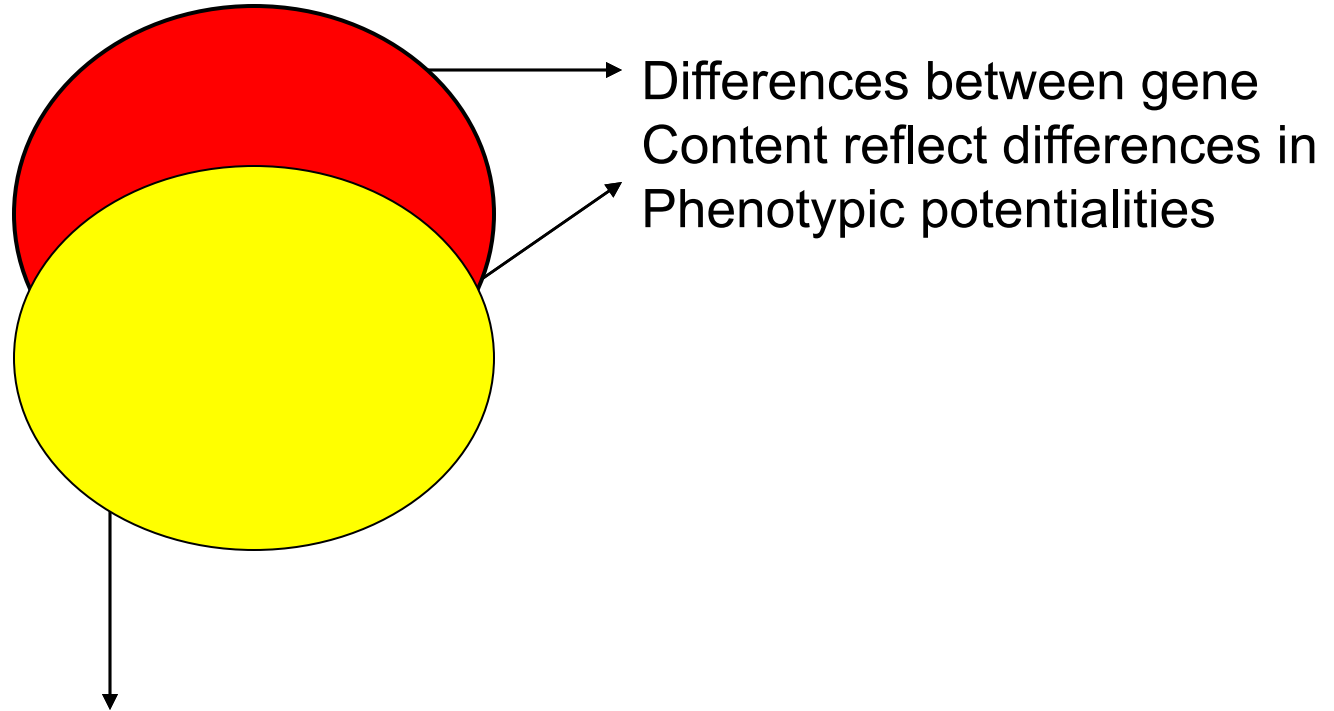


Genome trees based on gene content are remarkably similar to consensus on ToL



Presence / absence of genes

Gene content → co-evolution. (The easy case, few genomes.)



Genomes share genes for phenotypes they have in common

Three-way comparisons

Differential genome analysis

H. influenzae (1703)

132 metabolic enzymes

- sugar utilization
- pentose phosphate cycle
- fatty acid biosynthesis
- THF biosynthesis
- utilization/interconversion of nucleic acids/nucleotides
- 2 subunits of DNA polym. III

215 unknown

1 enzyme: dehydroquinase type II

12 host interaction factors
including virulence factors
such as urease

1 transporter

3 unknown

H. pylori (1577)

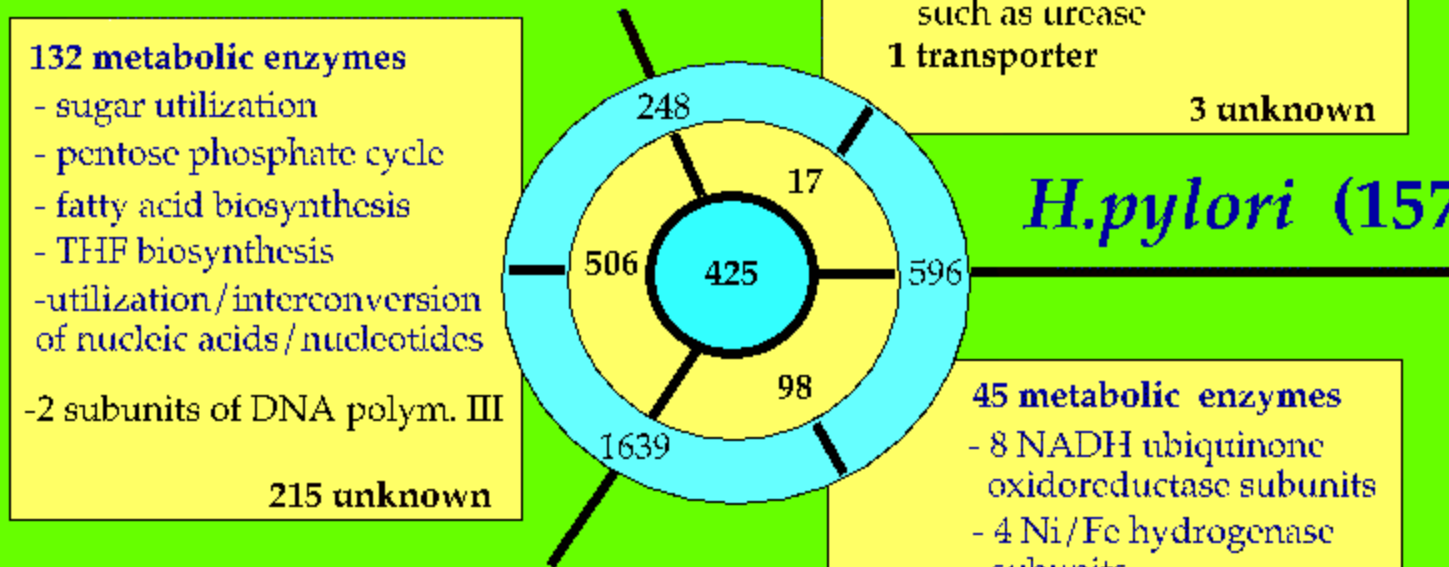
45 metabolic enzymes

- 8 NADH ubiquinone oxidoreductase subunits
- 4 Ni/Fe hydrogenase subunits

18 flagella-related genes

8 unknown

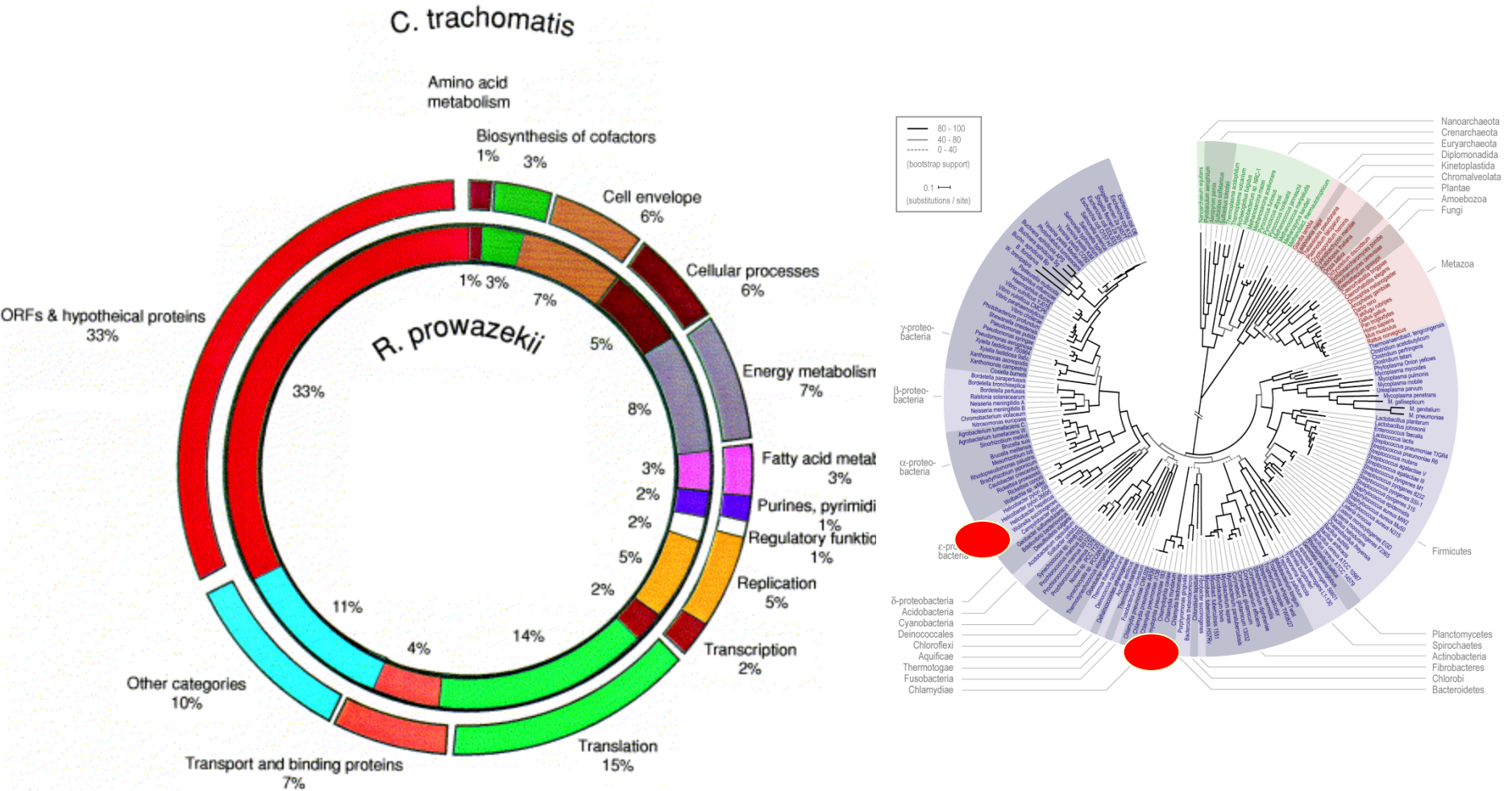
E. coli (4286)



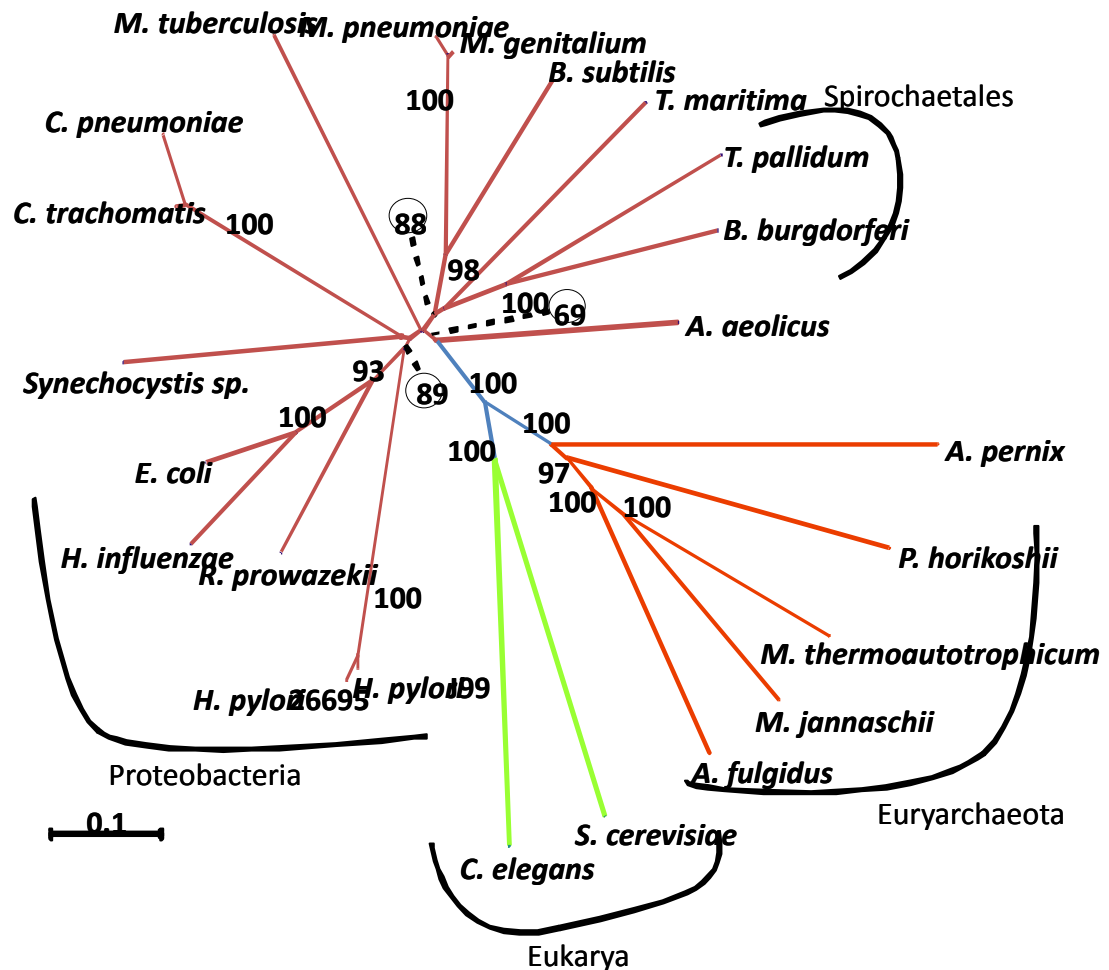
● orthologs shared by only two species

● species-specific orthologs

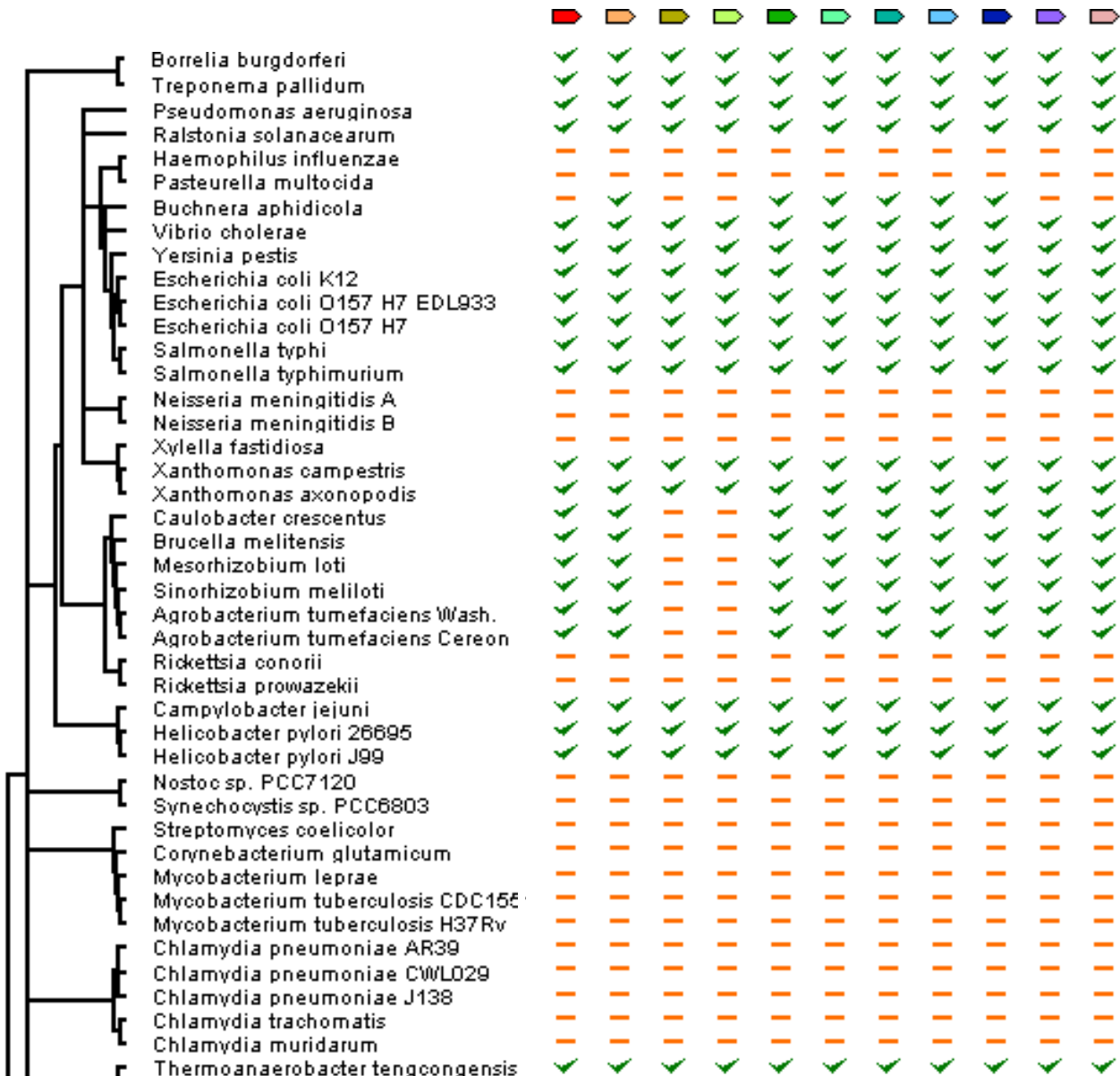
Convergence in functional classes of gene content in small intracellular bacterial parasites



Although we can, qualitatively, interpret the variations in shared gene content in terms of the phenotypes of the species, quantitatively they depend on the relative phylogenetic positions of the species. The closer two species are the larger fraction of their genes they share.



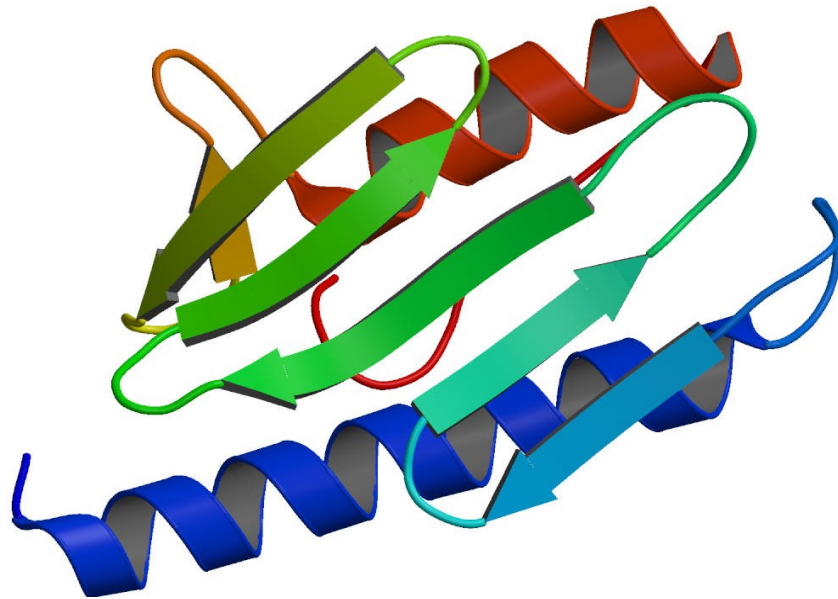
Co-occurrence of genes across genomes as prediction for interaction / association

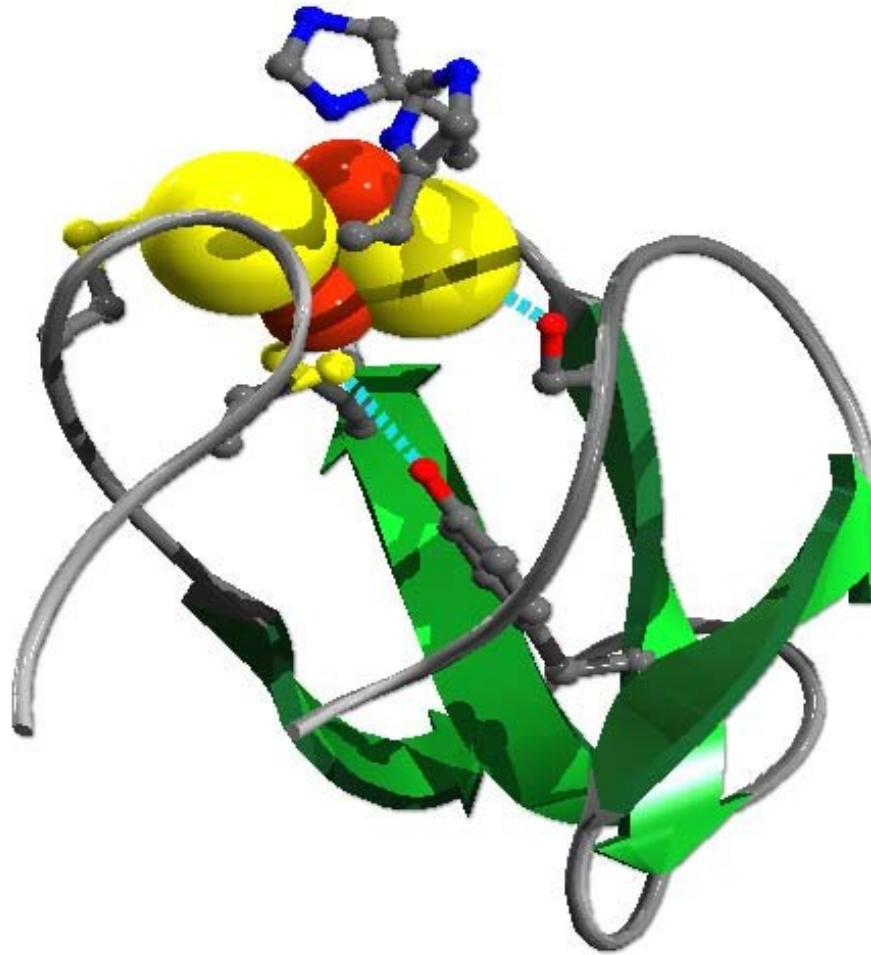


- i.e. two genes have the same presence/ absence pattern over multiple genomes:
- AKA phylogenetic profiles
- NB complete genomes absence -> needed for absence
- Correction for phylogenetic signal needed → events

Predicting function of a disease gene protein with unknown function, frataxin, using co-occurrence of genes across genomes / phylogenetic profiles

- Friedreich's ataxia
- No (homolog with) known function





Iron-Sulfur (2Fe-2S) cluster in the Rieske protein

Prediction:

© 2001 Oxford University Press

Human Molecular Genetics, 2001, Vol. 10, No. 21 2463–2468

The phylogenetic distribution of frataxin indicates a role in iron-sulfur cluster protein assembly

Martijn A. Huynen*, Berend Snel¹, Peer Bork and Toby J. Gibson¹

Biocomputing, EMBL/Max-Delbrueck-Center for molecular medicine, Berlin-Buch and ¹Biocomputing, EMBL, Meyerhofstrasse 1, 69117 Heidelberg, Germany

Received July 3, 2001; Revised and Accepted July 30, 2001

~Confirmation:

J|A|C|S
A R T I C L E S

Published on Web 04/26/2003

Iron–Sulfur Cluster Biosynthesis. Characterization of Frataxin as an Iron Donor for Assembly of [2Fe-2S] Clusters in ISU-Type Proteins

Taejin Yoon and J. A. Cowan*

*Contribution from Evans Laboratory of Chemistry, The Ohio State University,
100 West 18th Avenue, Columbus, Ohio 43210*

Received August 1, 2002; E-mail: cowan@chemistry.ohio-state.edu

STRING: functional protein association networks - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://string.embl.de/newstring.cgi/show_input_page.pl

Most Visited Google Gmail - Priority Inbox (...)

STRING: functional protein as...

Home · Download · Help/Info **STRING 8.3**

STRING - Known and Predicted Protein-Protein Interactions

search by name search by protein sequence multiple names multiple sequences

protein name: (examples: #1 #2 #3)

(STRING understands a variety of protein names and accessions; you can also try a [random entry](#))

organism: auto-detect

interactors wanted: COGs Proteins Reset GO!

please enter your protein of interest...

What it does ...

STRING is a database of known and predicted protein interactions. The interactions include direct (physical) and indirect (functional) associations; they are derived from four sources:

Genomic Context High-throughput Experiments (Conserved) Coexpression Previous Knowledge

STRING quantitatively integrates interaction data from these sources for a large number of organisms, and transfers information between these organisms where applicable. The database currently covers 2,590,259 proteins from 630 organisms.

More Info Funding / Support Acknowledgements Use Scenarios

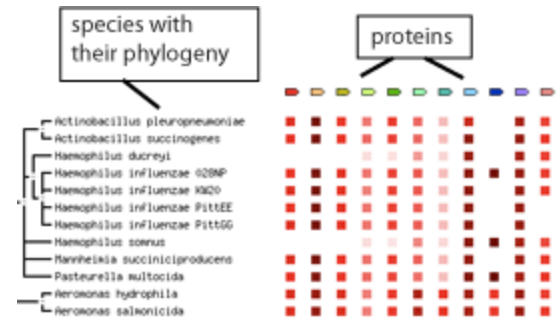
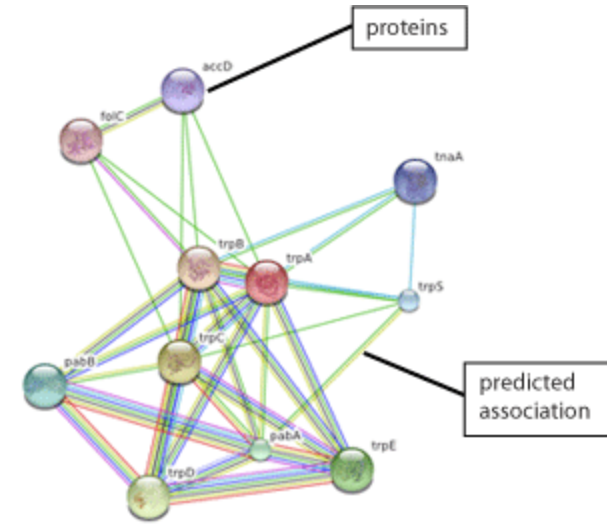
STRING ([Search Tool for the Retrieval of Interacting Genes/Proteins](#)) is being developed at CPR, EMBL, SIB, KU, IUD and UZH. STRING references: [Jensen et al. 2009](#) / [2007](#) / [2005](#) / [2003](#) / [Snel et al. 2000](#). Miscellaneous: [Access Statistics](#), [Robot Access Guide](#), [STRING/STITCH Blog](#), [Supported Browsers](#).

What's New? This is version 8.3 of STRING - June 2010: the latest interaction data, updated textmining, and bugfixes ...

Sister Projects: check out [STITCH](#) and [eggNOG](#) - two sister projects built on STRING data!

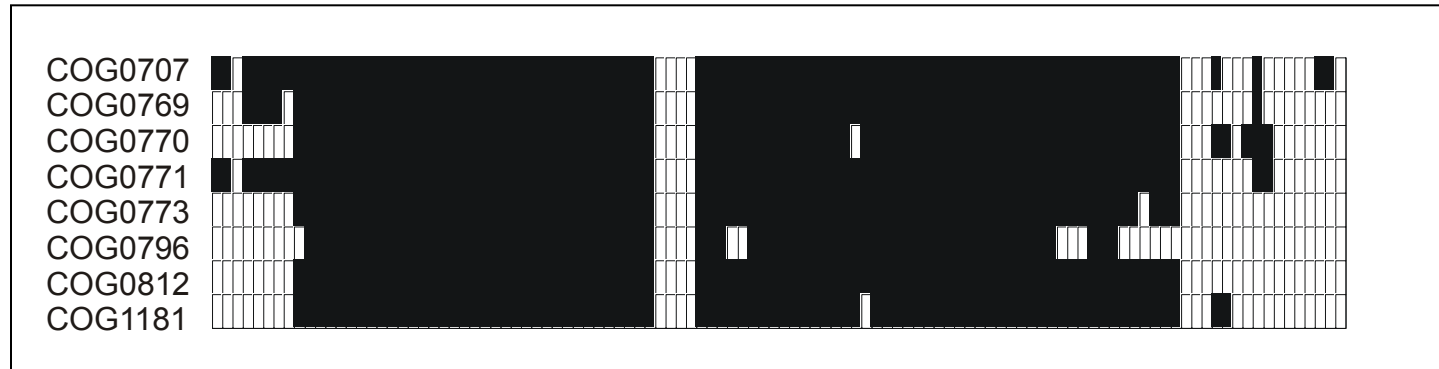
Previous Releases: Trying to reproduce an earlier finding? Confused? Refer to our [old releases](#).

xv controls UNREGISTERED Display Root



- If two genes have a “significantly” similar presence/absence pattern which is different from the phylogenetic signal, than their proteins are likely to interact / be in the same process/pathway
- This pattern can be created by independent loss and/or horizontal gene transfer (of operons)

However



peptidoglycan biosynthesis pathway (highly cohesiveness, far from perfect)



ribose phosphate metabolism (not cohesive at all)

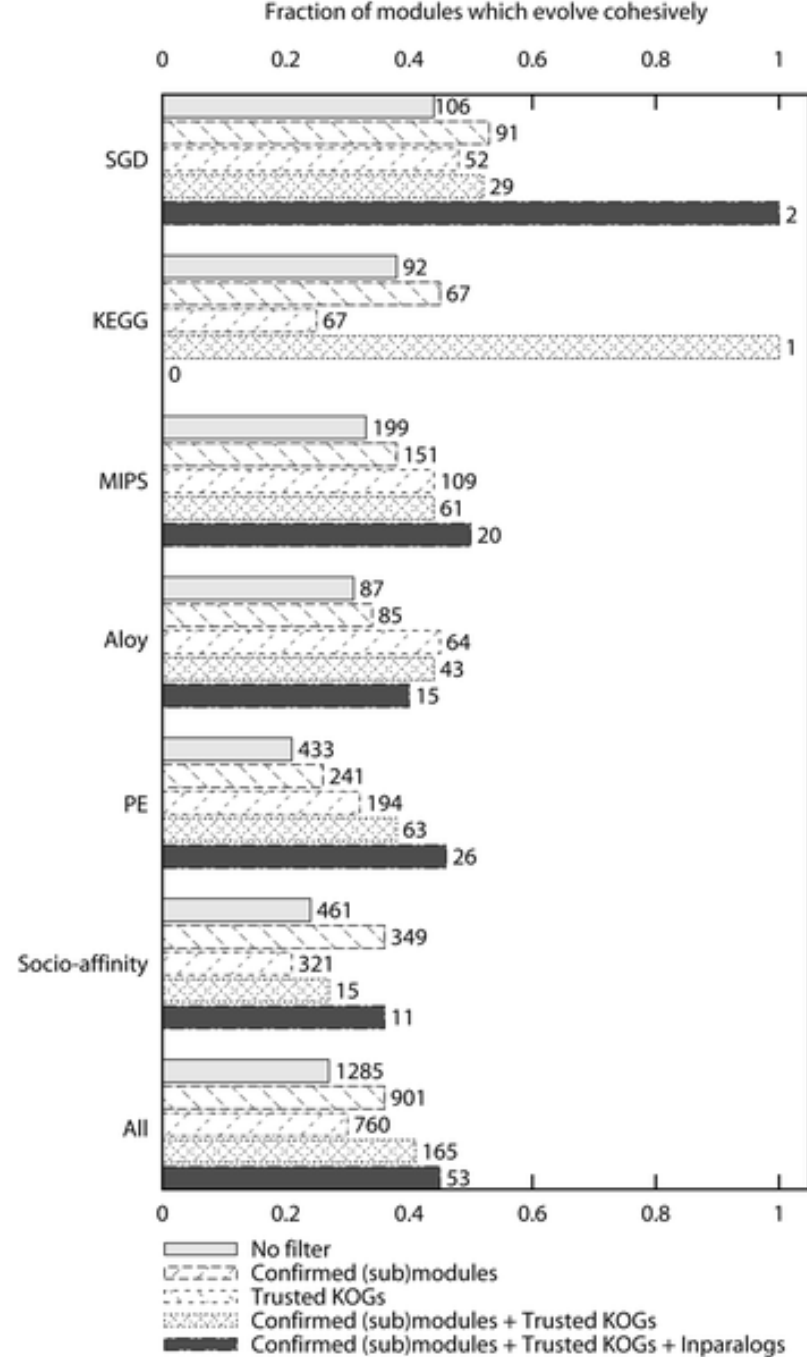
Very few functional modules are perfect; limited cohesiveness; functional units vs evolutionary units

Why?

If the phyletic patterns of two proteins are highly similar they tend to interact, but the reverse is not generally true!

~50% of modules (such as protein complexes) do not have highly similar phyletic patterns.

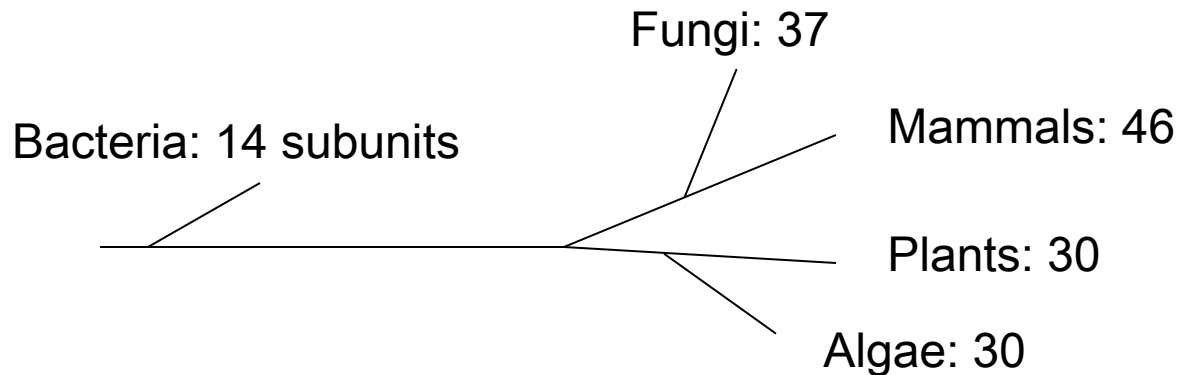
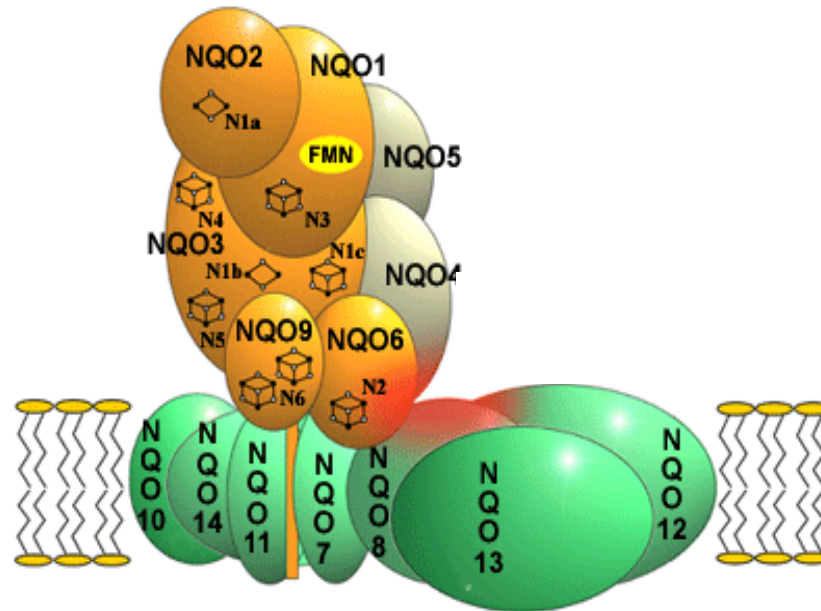
This seems to not only depend on dataset, noise in orthology detection, or noise in module definition

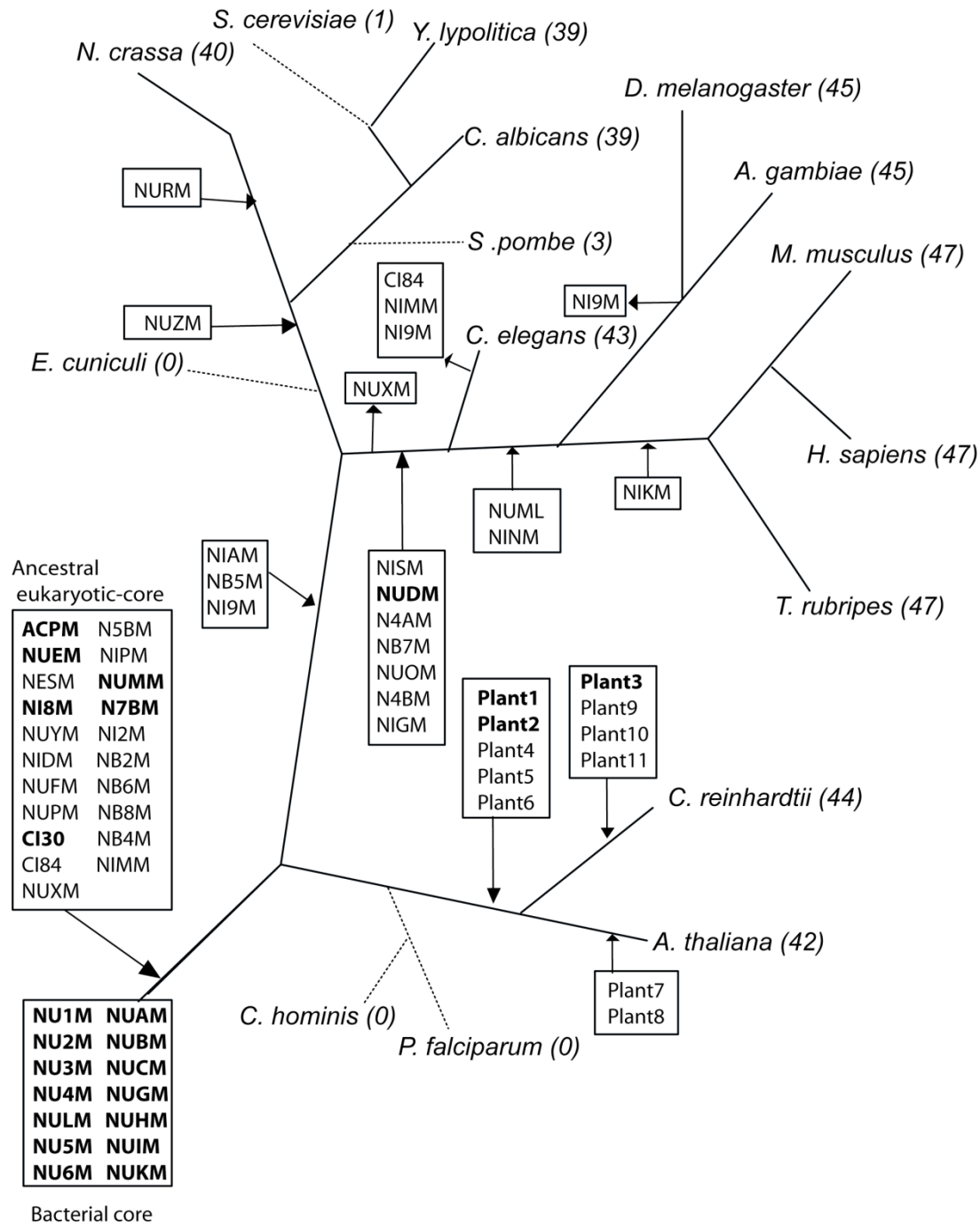


Explaining discordant phyletic patterns of proteins that interact

- Many cases stories and a few large scale studies
- (we could also just say that evolution is flexible and proteins change function; which I am not going to argue with but (A) conservation of interaction and (B) this is a “just so”, non testable explanation)

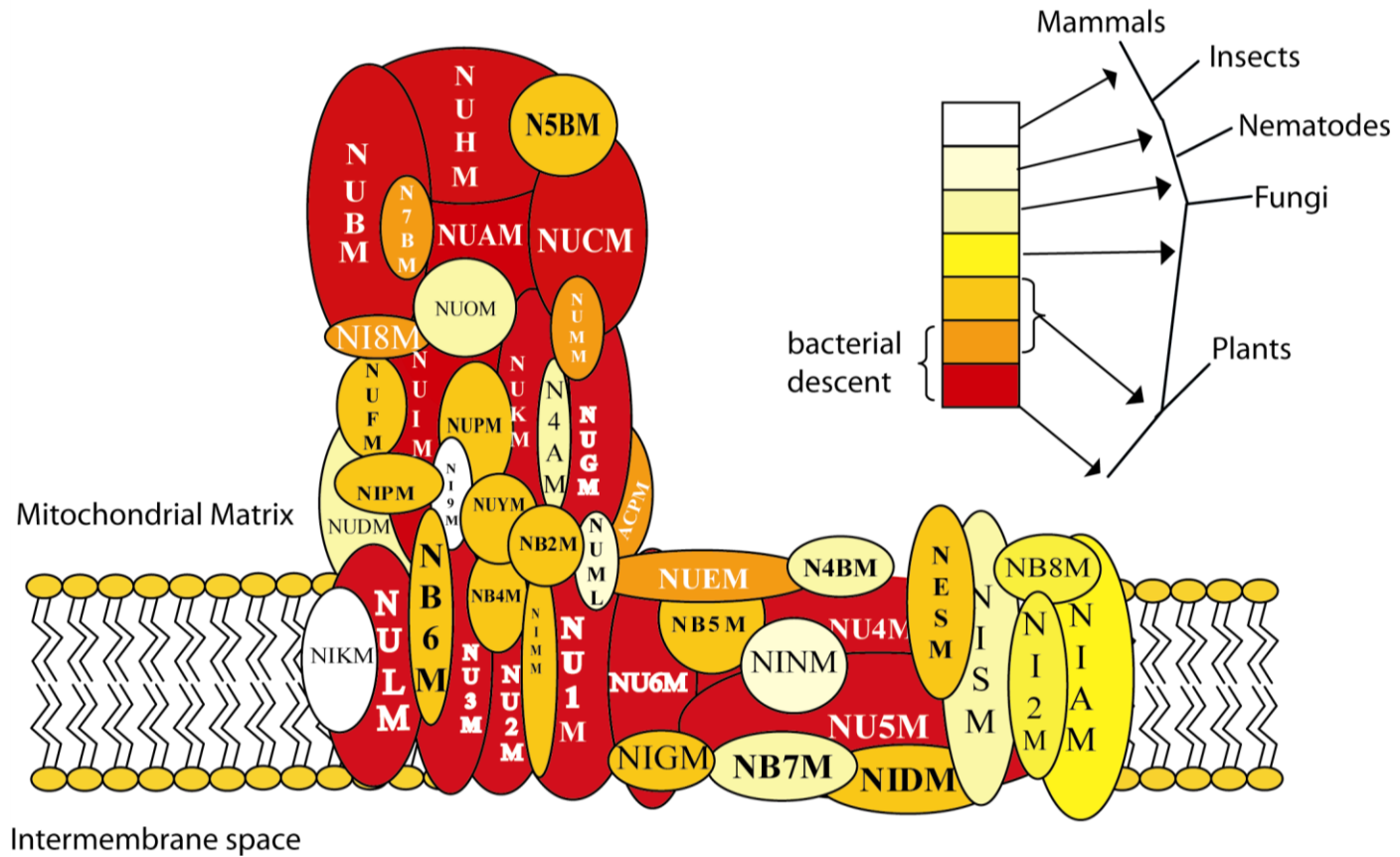
Tracing the evolution of NADH:ubiquinone oxidoreductase (Complex I of the oxidative phosphorylation), from 14 subunits (Bacteria) to 46 subunits (Mammals) by comparative genome analysis





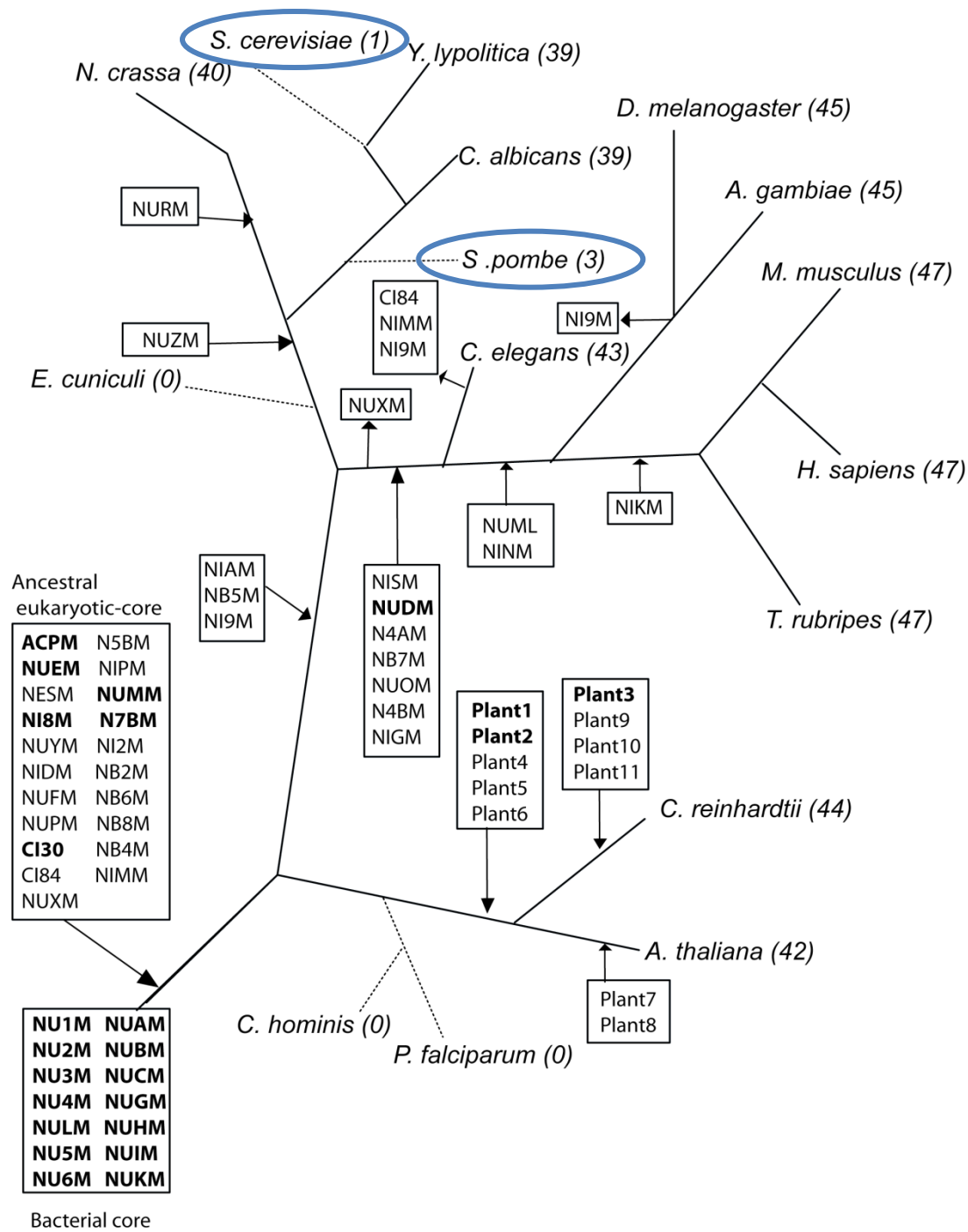
Reconstructing Complex I evolution by mapping the variation onto a phylogenetic tree. After an initial “surge” in complexity (from 14 to 35 subunits in early eukaryotic evolution) new subunits have been gradually added and incidentally lost. most other loss is large scale

In the eukaryotic evolution of Complex I, new subunits have been added “all over” the complex

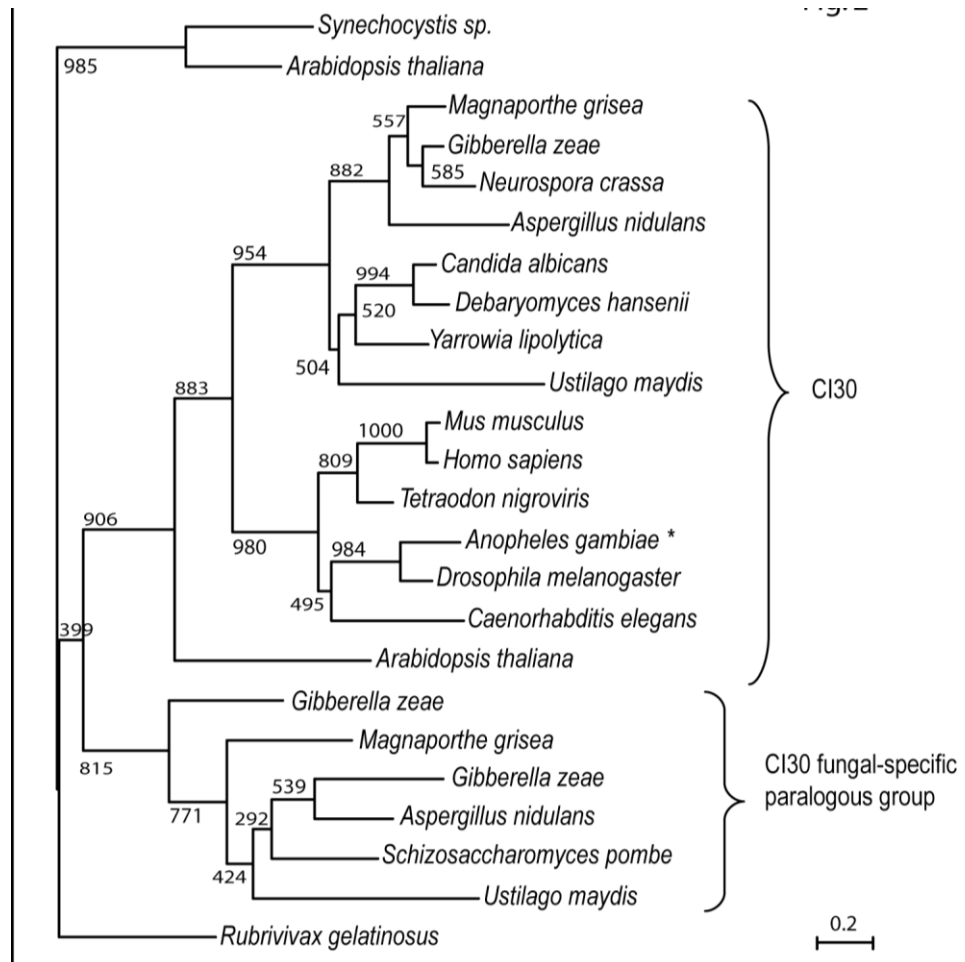


Reconstructing Complex I evolution by mapping the variation onto a phylogenetic tree. After an initial “surge” in complexity (from 14 to 35 subunits in early eukaryotic evolution) new subunits have been gradually added and incidentally lost., most other loss is large scale

Complex I loss is not always “complete”, *S.cerevisiae* and *S.pombe* have retained 1 and 3 proteins



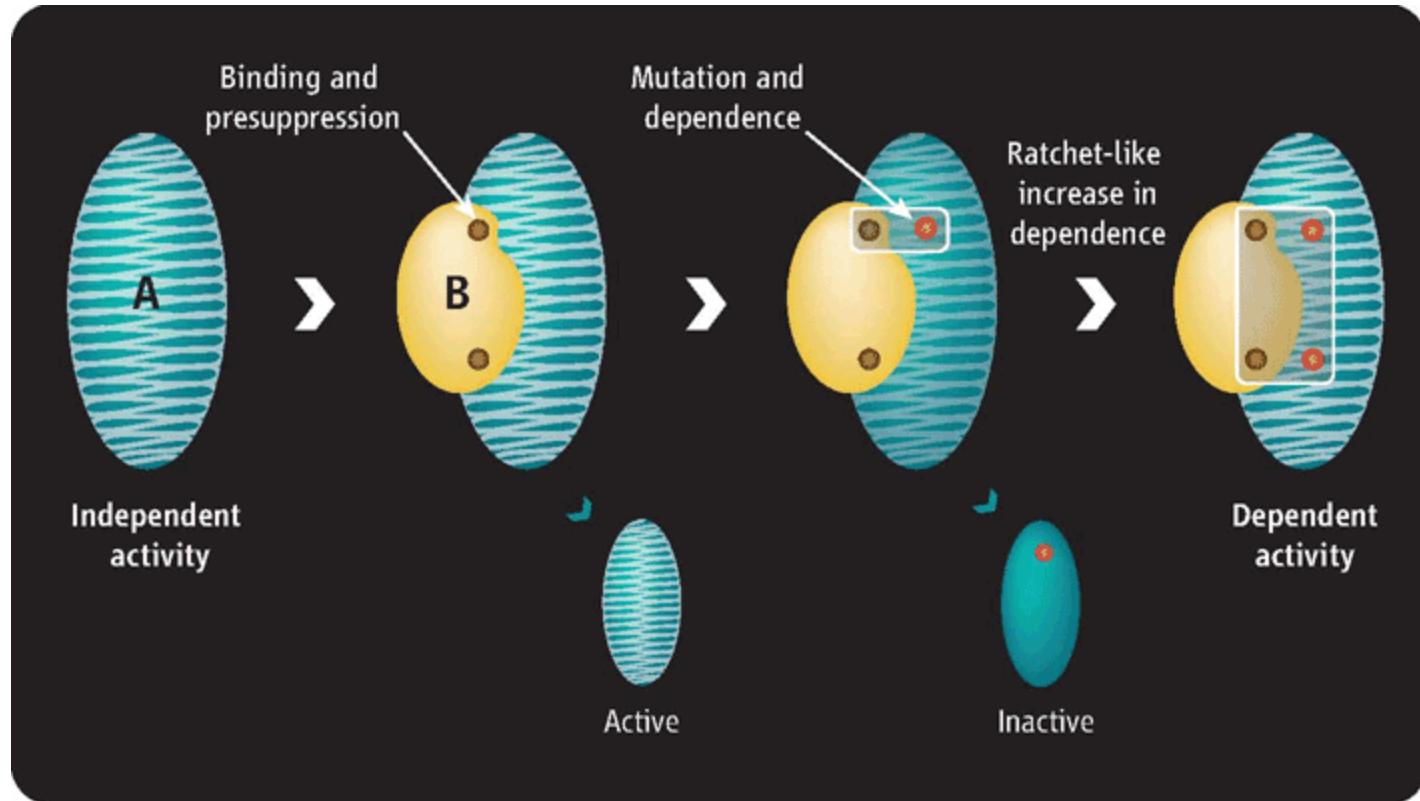
Phylogeny of a “remaining” complex I protein in pombe



?

“The Complex I assembly protein CI30 has been duplicated in the Fungi. This can explain the presence of a CIA30-homolog in Complex I-less *S.pombe*”

Why the accumulation: a neutral explanation



fixation of neutral or slightly deleterious features as a general and unavoidable source of complexity in taxa with small populations

Science. 2010 Nov 12;330(6006):920-1.

Cell biology. Irremediable complexity?

[Gray MW](#), [Lukes J](#), [Archibald JM](#), [Keeling PJ](#), [Doolittle WF](#).

How to falsify?

e.g. *Neurospora* mito-TyrRS

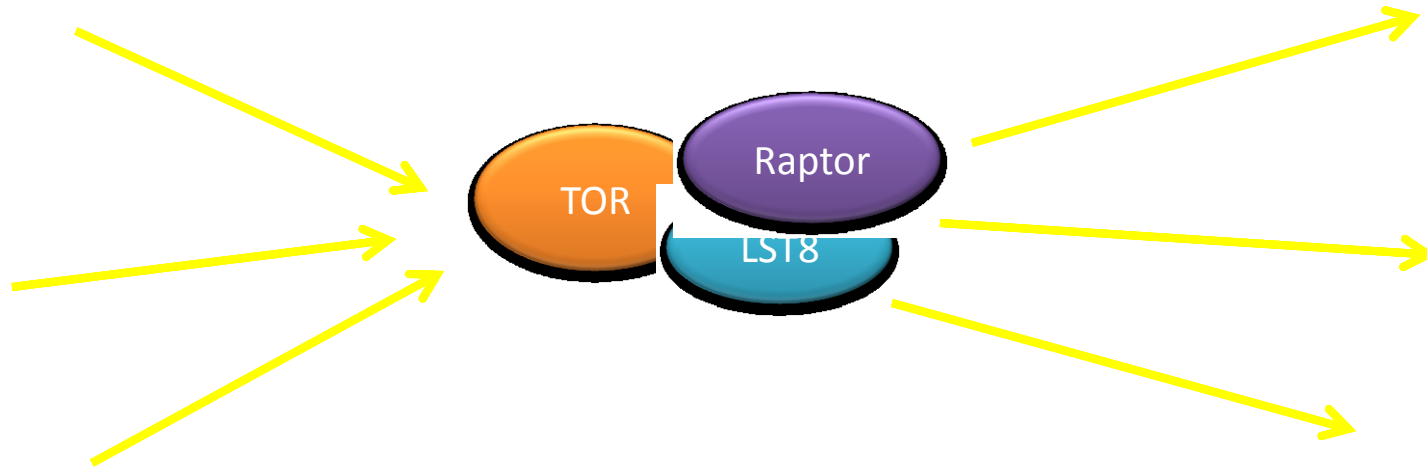
- *Neurospora* mitochondrial genome encodes several introns which require a tyrosyl tRNA synthetase (TyrRS) to splice.
- “to compensate for structural defects acquired by the intron sequences “
- BUT Introns with defects arising -> negative selection
- ? Reverse: first binding (fortuitously or for reason unrelated to splicing)—> accumulation of mutations in the intron that inactivate splicing, if TyrRS not bound.
- Because the compensatory / suppressive activity exists before mutation “presuppression,”
- the protein dependence by the intron could be selectively neutral (or slightly disadvantageous

“Constructive neutral evolution”

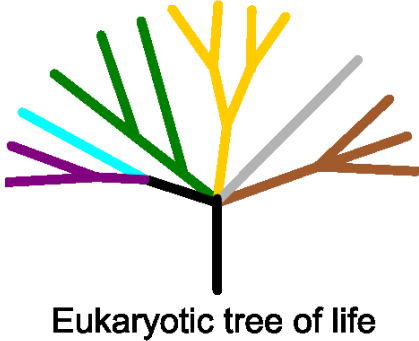
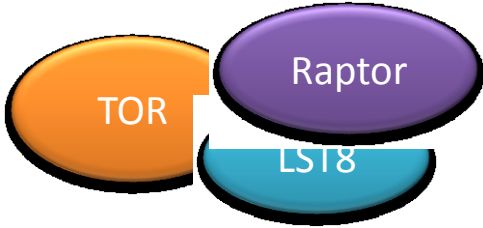
- Suggested that many taxon specific subunits (taxon specific proteins that are a subunit in a complex) are regulatory subunits
- Hypothesis: neutrally added but necessary subunits could have been appropriated as regulatory subunits?

TOR1 complex

- Kinase
- Regulates growth
- Mutations of TOR1 components involved in Cancer



Evolution of TOR

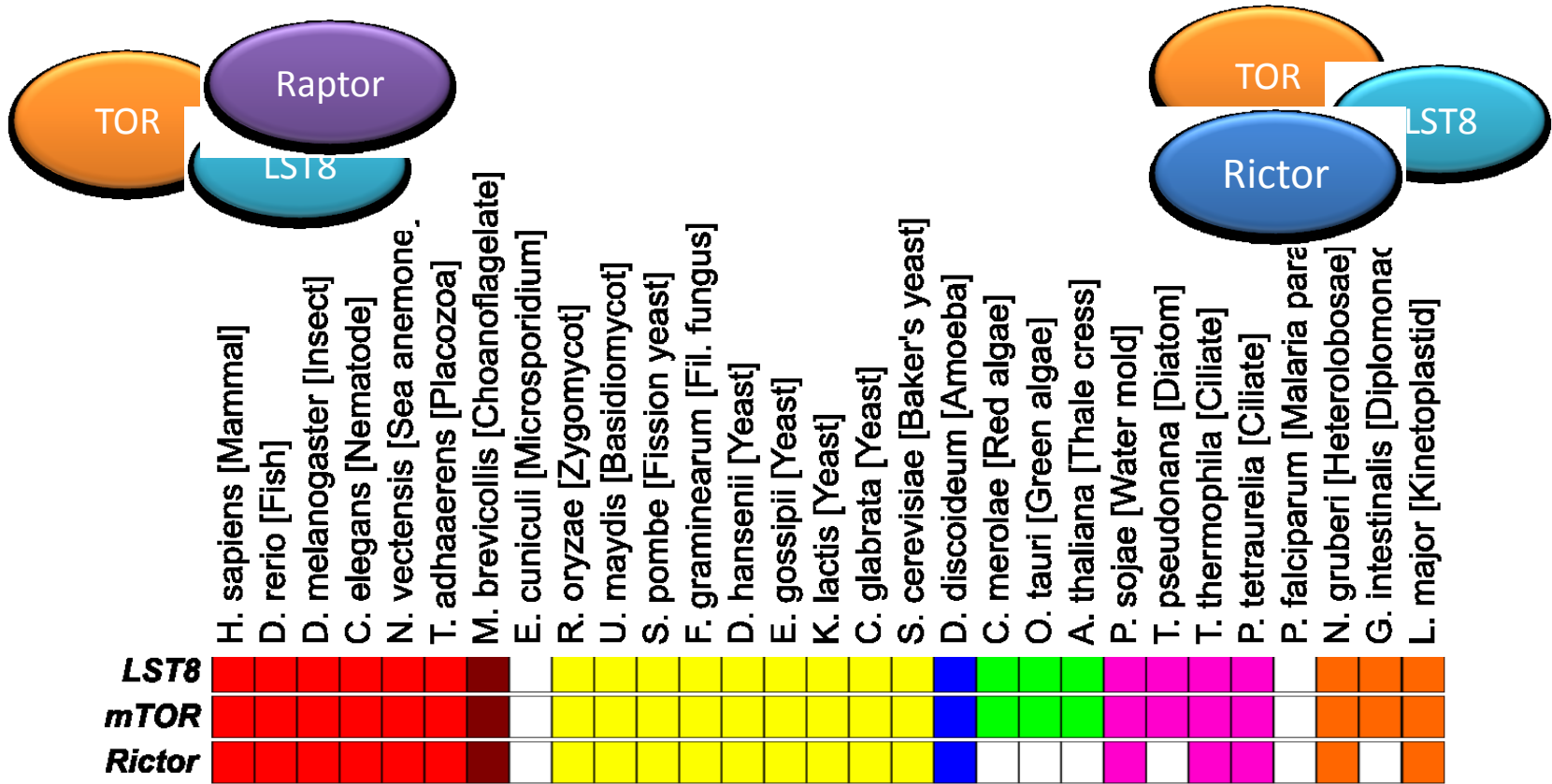


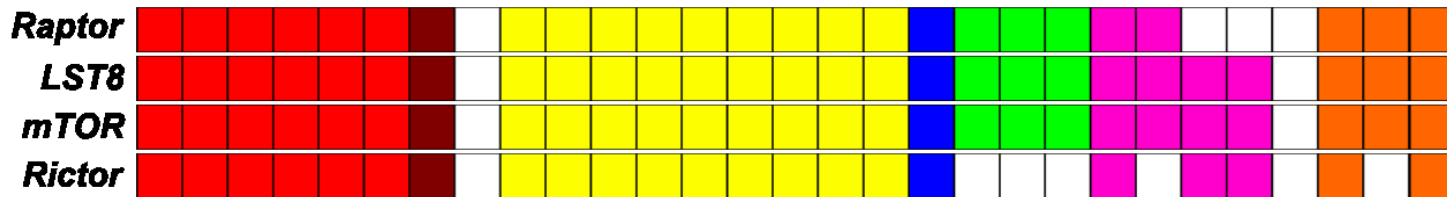
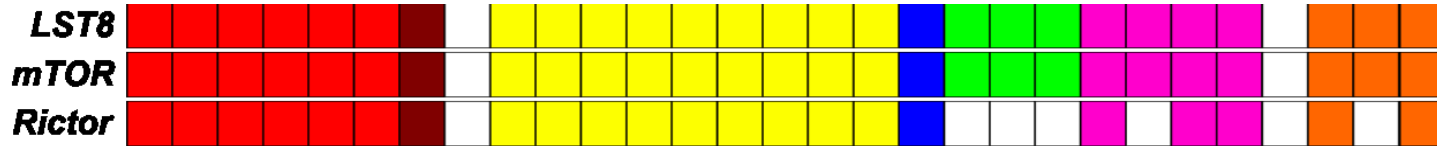
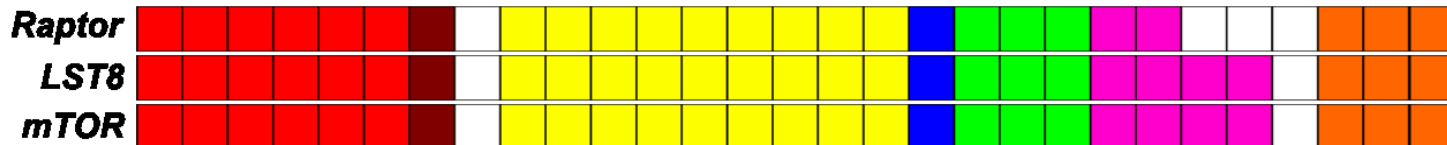
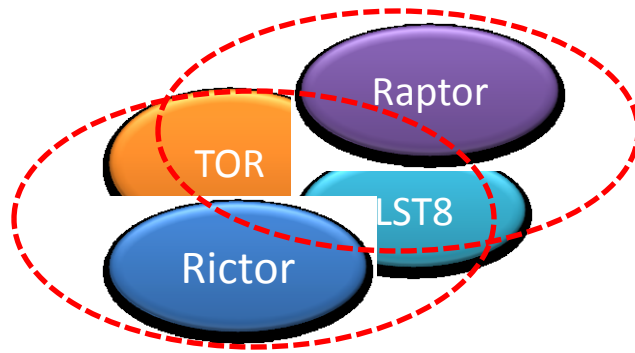
	1. sapiens [Mammal]	2. rerio [Fish]	3. melanogaster [Insect]	4. elegans [Nematode]	5. vectensis [Sea anemone]	6. adhaerens [Placozoa]	7. brevicollis [Choanoflagelate]	8. cuniculi [Microsporidium]	9. oryzae [Zygomycot]	10. maydis [Basidiomycot]	11. pombe [Fission yeast]	12. graminearum [Fil. fungus]	13. hansenii [Yeast]	14. gossipii [Yeast]	15. lactis [Yeast]	16. glabrata [Yeast]	17. cerevisiae [Baker's yeast]	18. discoideum [Amoeba]	19. merolae [Red algae]	20. tauri [Green algae]	21. thaliana [Thale cress]	22. sojae [Water mold]	23. pseudonana [Diatom]	24. thermophila [Ciliate]	25. tetraurelia [Ciliate]	26. falciparum [Malaria parasite]	27. gruberi [Heterobosae]	28. intestinalis [Diplomonad]	29. major [Kinetoplastid]					
Raptor	Red	Red	Red	Red	Red	Red	Red	White	Yellow	Yellow	Yellow	Yellow	Yellow	Yellow	Yellow	Blue	Blue	Green	Green	Green	Green	Green	Pink	Pink	Pink	White	White	White	White	Orange	Orange	Orange		
LST8	Red	Red	Red	Red	Red	Red	Red	White	Yellow	Yellow	Yellow	Yellow	Yellow	Yellow	Yellow	Blue	Blue	Green	Green	Green	Green	Green	Pink	Pink	Pink	Pink	Pink	White	White	White	White	Orange	Orange	Orange
mTOR	Red	Red	Red	Red	Red	Red	Red	White	Yellow	Yellow	Yellow	Yellow	Yellow	Yellow	Yellow	Blue	Blue	Green	Green	Green	Green	Green	Pink	Pink	Pink	Pink	Pink	White	White	White	White	Orange	Orange	Orange

Does evolution of TOR make more sense if we consider the whole network of interactions:

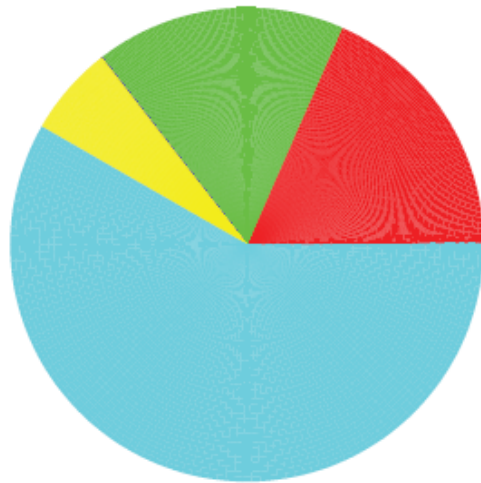
TOR2 complex

- TOR2 is involved in rearrangement of cytoskeleton

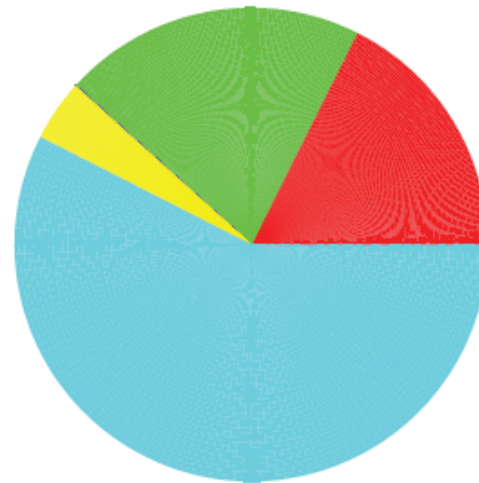




A



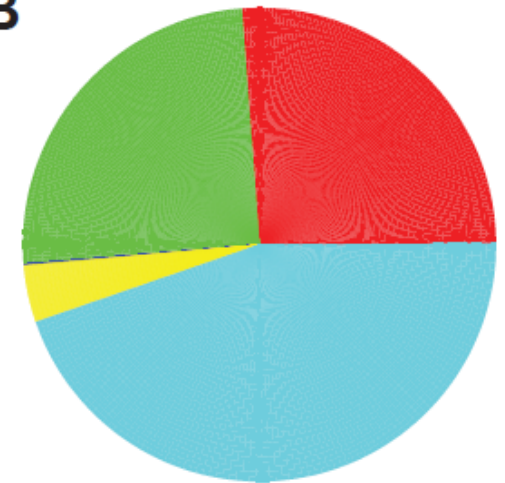
yeast data



human data

protein complex

B



yeast data

protein-protein
interaction (BioGRID)

? Taxonomic subunit constructive neutral
evolution ?

Considering triangles of interactions instead of pairs: a complementarity score

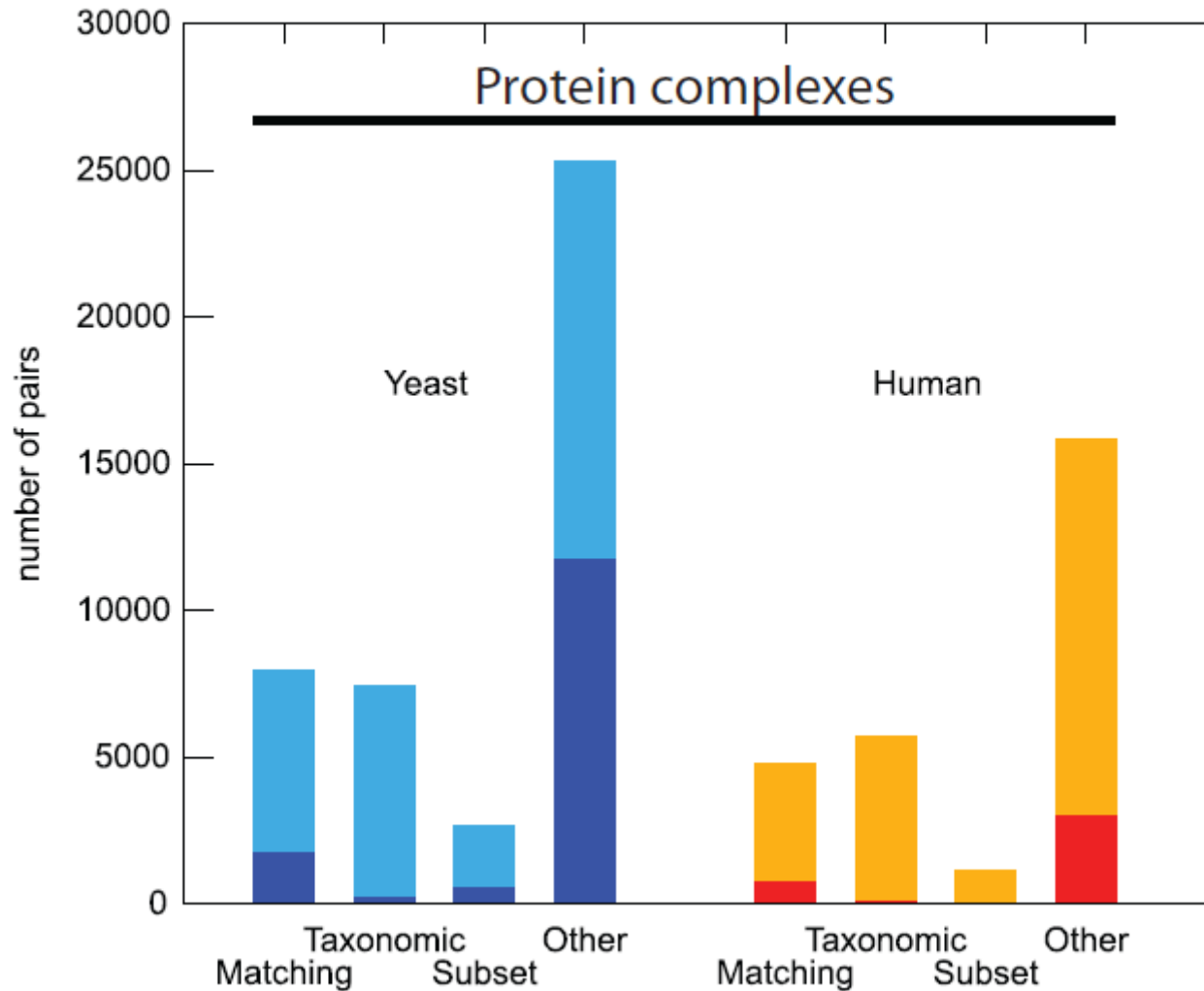
C	Genes		
	A	B	C
	O	O	O
	O	O	X
	O	X	O
	O	X	X
	X	O	O
	X	O	X
	X	X	O
	X	X	X

smaller number = nGood

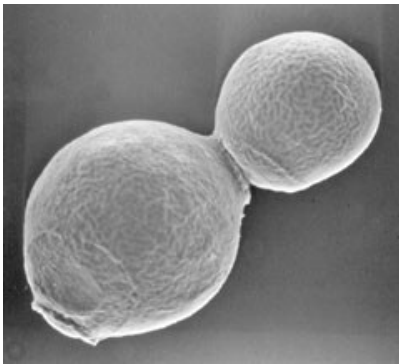
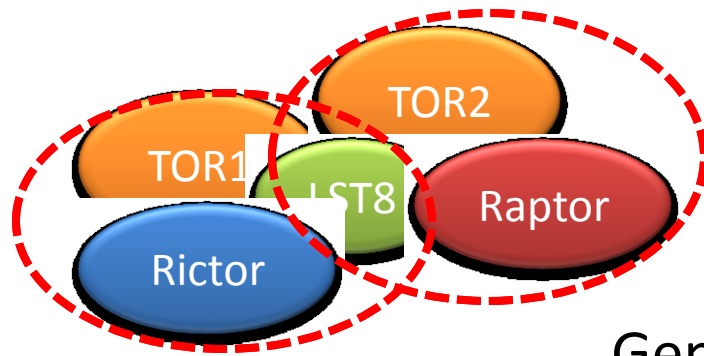
sum of these = nBad

$$\text{Complementarity} = n\text{Good} - n\text{Bad}$$

A substantial fraction of the “others” has a high complementarity score



Gene duplications are important in evolution and function of TOR complexes



Gene duplication of the tor kinase in *Sacceromyces cerevisiae*, chytrid fungi, oomycetes, poplar.

In yeast molecular biology demonstrated specialization

Why in some species this happens and others not?

Other forms of co-evolution

- Speed/rate of evolution
- Acceleration after loss of binding partner
- Compensatory mutations, co-evolving residues: old problem, never solved, now maybe possibly in reach thanks to evoFold (?), not applied to study evolution yet

Non-orthologous gene displacement/analogous proteins

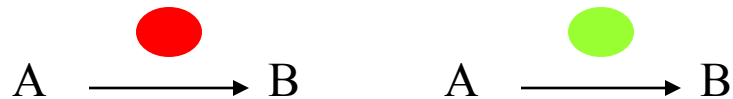
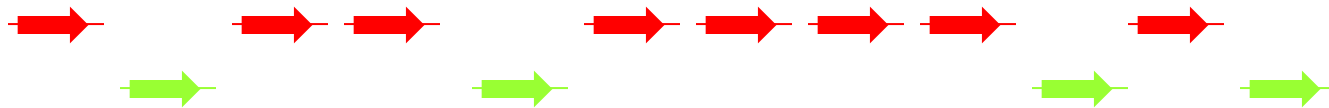
- First systematic analysis on *M.genitalium*
(Koonin et al., Trends Genet. 1997)

TABLE 1. Non-orthologous genes coding for the same function in *Mycoplasma genitalium* and *Haemophilus influenzae*

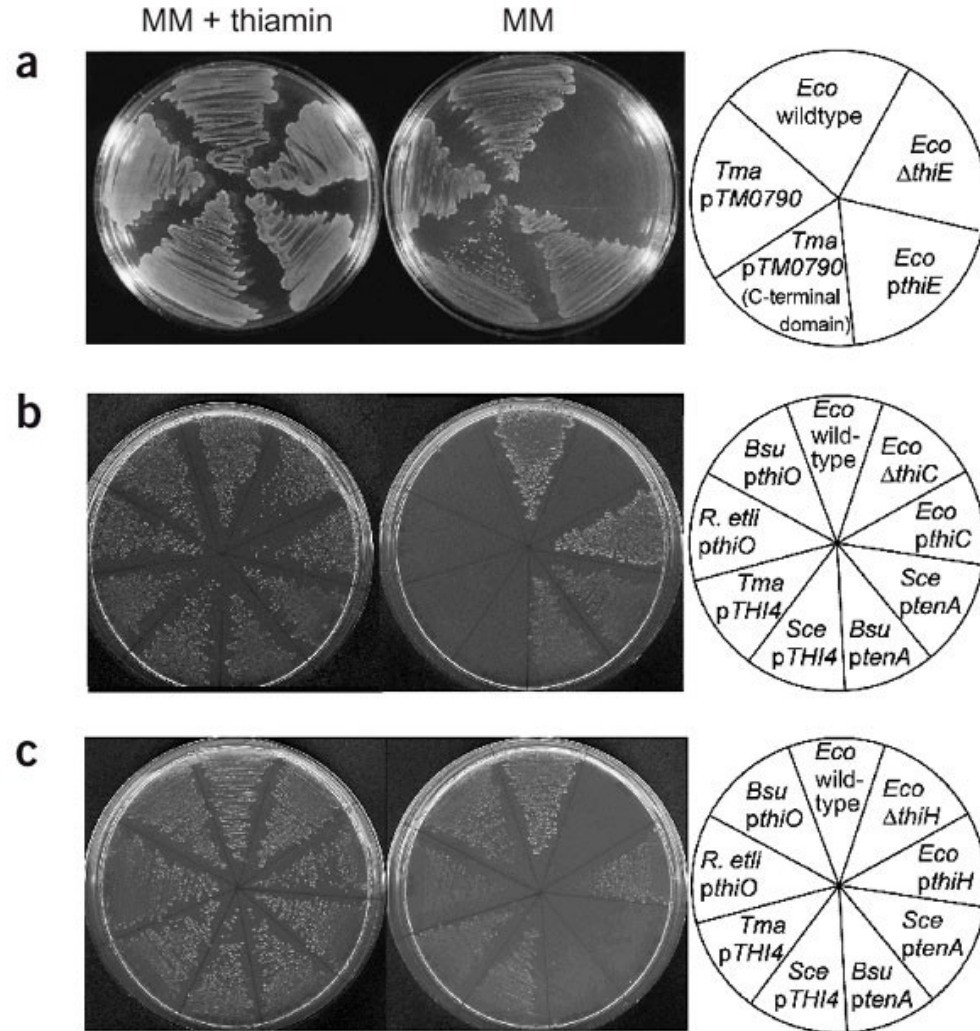
Enzyme	<i>M. genitalium</i>		<i>H. influenzae</i>		Comment ^c
	Gene ^a	Orthologs ^b	Gene ^a	Orthologs ^b	
No sequence similarity between <i>M. genitalium</i> and <i>H. influenzae</i> proteins					
Phosphoglycerate mutase	MG430 (<i>yibO</i>)	PMGI_BACSU PMGI_ECOLI PMGI_MAIZE	HI0757 (<i>gpmA</i>)	PMG1_ECOLI PMGM_HUMAN not in G(+)	<i>Escherichia coli</i> encodes both types of enzymes
L-lactate dehydrogenase	MG460	LDH_BACSU LDHM_HUMAN	HI1739B (<i>lctD</i> or <i>lldD</i>)	LLDD_ECOLI G(+)	The HI enzyme is distantly related to eukaryotic cytochrome B2
Lipoate-protein ligase	MG270	LPLA_ECOLI SCYJL046W_1	HI0027 (<i>lipB</i>)	LIPB_ECOLI S51458 (yeast)	<i>E. coli</i> and yeast encode both types of enzymes
Nucleoside diphosphate kinase	MG264 ^d MG268 ^d	None	HI0876 (<i>ndk</i>)	NDK_ECOLI NDKB_HUMAN	The two predicted kinases in MG are candidates for this indispensable activity
DNA polymerase, repair	MG261 (<i>dnaE</i>)	DP3A_HAEIN DP3A_ECOLI	HI0856 (<i>polA</i>)	DPO1_ECOLI DPO1_MYCTU	MG encodes two homologs of DNA polymerase III. MG261 is the likely repair polymerase as it belongs to a putative repair operon ⁶
RNase H	MG262 ^d	DPO1_BACCA DPO1_HAEIN	HI0138 (<i>rnbA</i>); HI1059 (<i>rnbB</i>)	RNH_ECOLI RNH1_YEAST RNH2_ECOLI MC326_1 (<i>M. capricol.</i>) SC23CDS_13 (yeast)	MG262 is homologous to the 5'-3' exonuclease domain of DNA polymerase I. It is predicted to replace the two unrelated RNases H of HI in primer removal during DNA replication
Glycyl-tRNA synthetase	MG251	SYG_HUMAN	HI0927 (<i>glyQ</i>) HI0924 (<i>glyS</i>)	SYGA_ECOLI SYGB_ECOLI CTU20547_1 (<i>Chlamydia</i>) G(-)	The MG enzyme contains one subunit, the HI counterpart two
Paralogs in <i>M. genitalium</i> and <i>H. influenzae</i>					
Prolyl-tRNA synthetase	MG283	YHI0_YEAST	HI0729 (<i>proS</i>)	SYP_ECOLI YER7_YEAST	Yeast encodes both types of enzymes
Cytidine deaminase	MG052	CDD_BACSU CDD_HUMAN	HI1350 (<i>cdd</i>)	CDD_ECOLI	The MG cytidine deaminase is more closely related to eukaryotic enzymes than to those from G(+) bacteria

The opposite of co-occurrence: anti-correlation / complementary patterns: predicting analogous enzymes

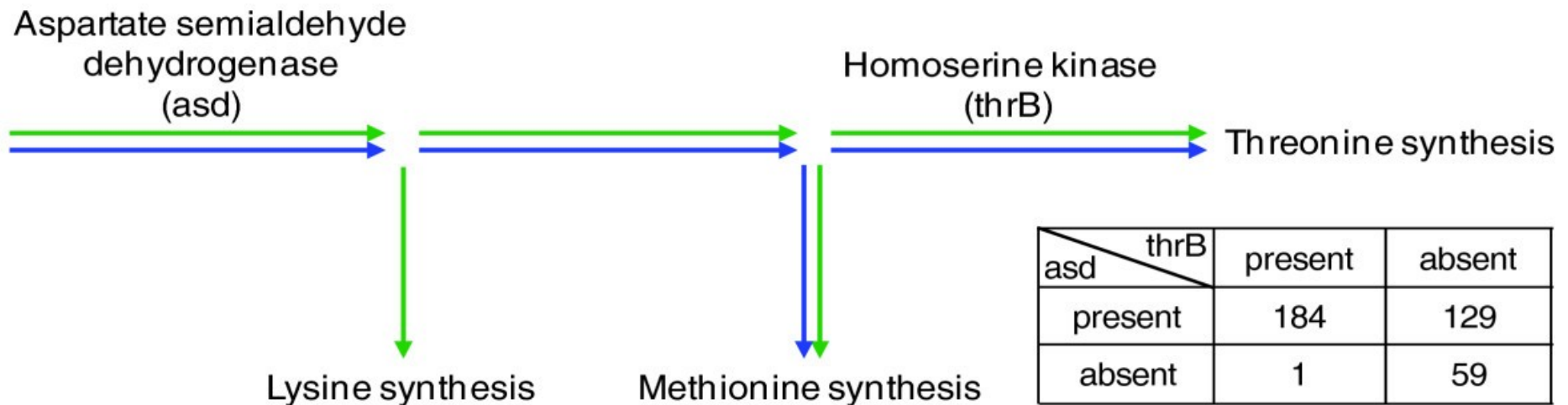
Genes with complementary phylogenetic profiles could have a similar biochemical function.



Prediction of analogous enzymes is confirmed

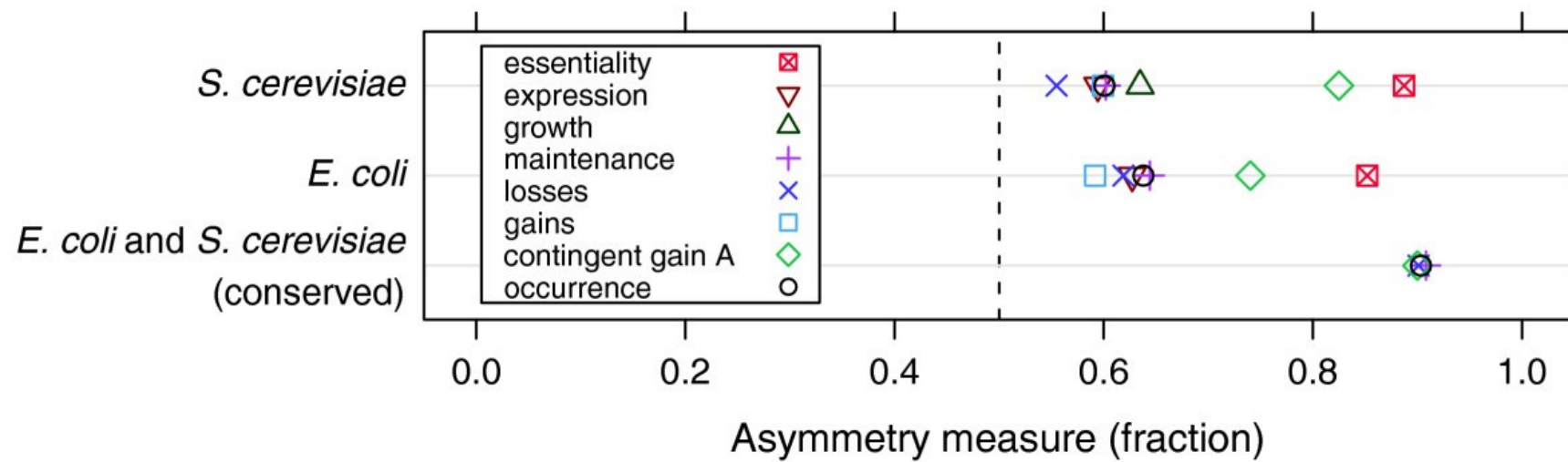
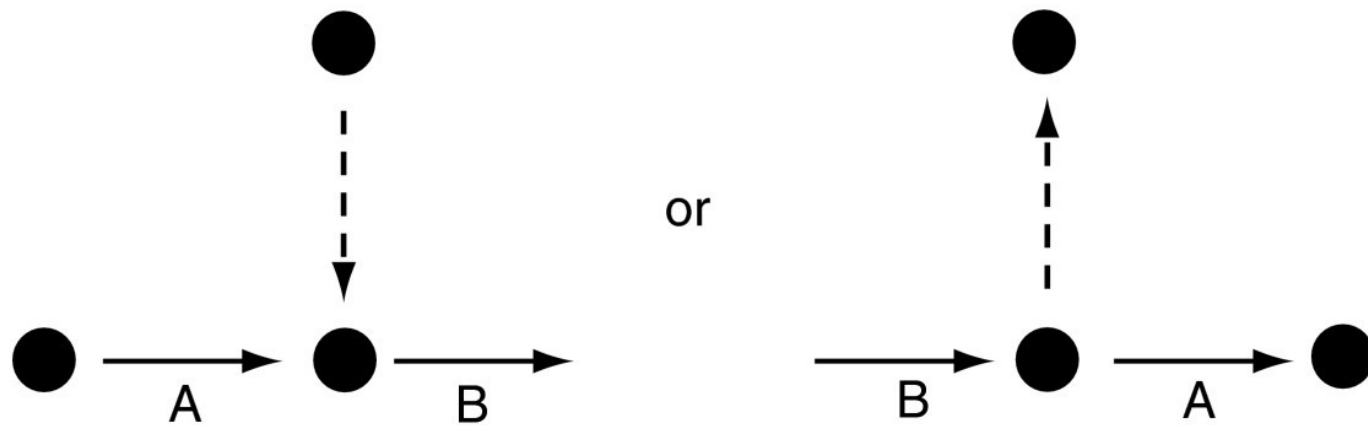


Asymmetric functional/metabolic relations explain non-similar presence absence patterns

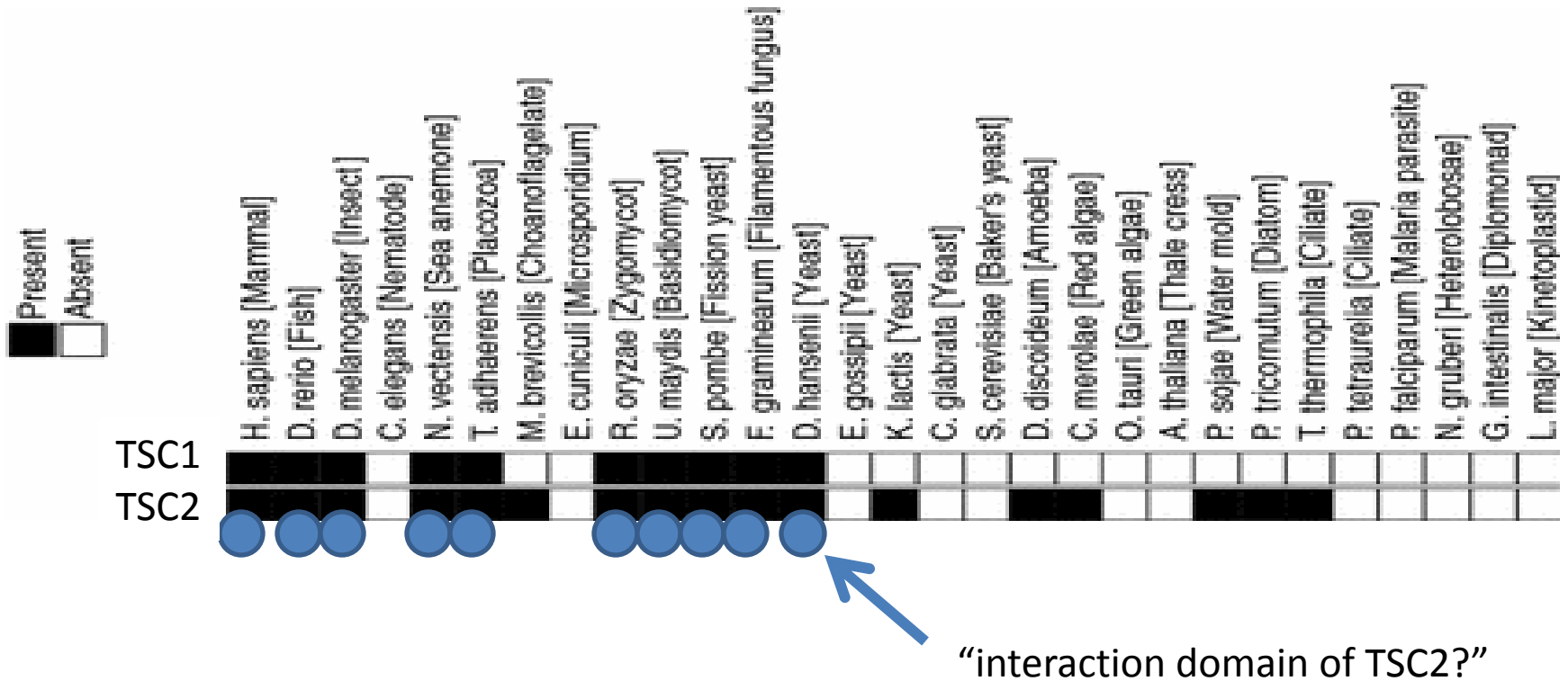


[Asymmetric relationships between proteins shape genome evolution.](#)

Notebaart RA, Kensche PR, **Huynen MA**, Dutilh BE.
Genome Biol. 2009 Feb 12;10(2):R19.



Domains vs proteins?



But why the innovation? Regulatory subunit / neutral accumulation of taxon specific subunits / constructive neutral evolution?

- Explaining the evolution of genes stimulates a better / more focused discussion on what we mean by gene function(al) relationship
- The more/better HTP functional data, the better for studying genome evolution
- Many different plausible, interlocking, reasons for disrupted co-occurrence across genomes of interacting proteins; (role of duplication least systematically researched?)