

A Comparative Genomic Investigation of Fungal Genome Evolution

Jason Stajich
Duke University
University Program in Genetics & Genomics

Evolutionary genomics

Evolutionary & Organismal Biology

Phylogeny
Population genetics
and structure
Phenotype
Ecological adaptations

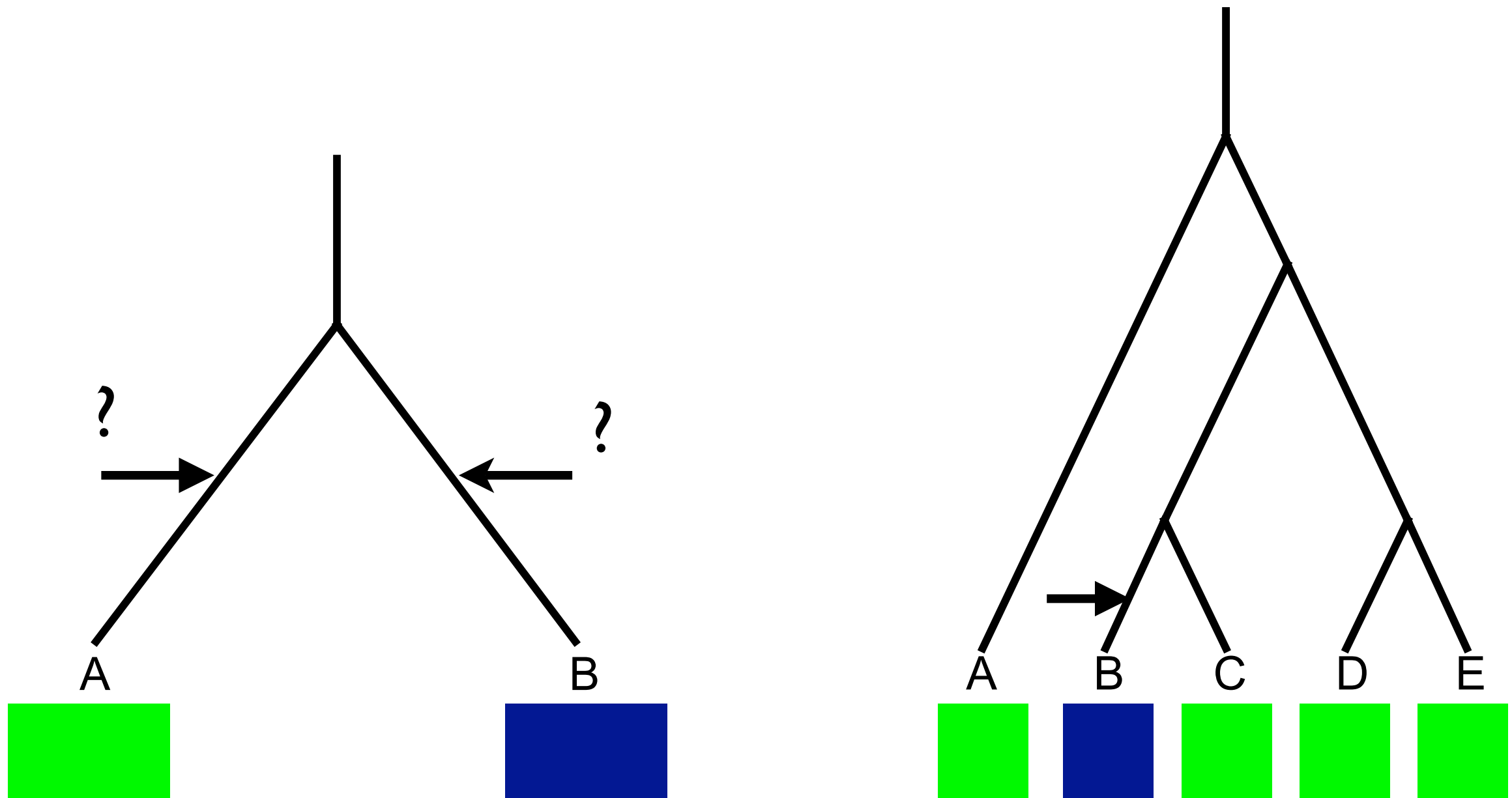
Comparative Genomics

Molecular evolution
Gene order
Gene families
Gene and genome
structure
Gene content
Conserved elements
Rates of molecular
evolution
Gene function
inference

Model Systems

Genetic tools
Gene function &
expression
Regulatory networks
Pathways
Molecular & cellular
biology
Disease models

Power of the comparative approach



Industrial uses of fungi

- Bread, beer, wine - *Saccharomyces cerevisiae*
- Sake and soy sauce - *Aspergillus oryzae*
- Dairy - *Penicillium roqueforti*, *Kluyveromyces lactis*
- Citric acid - *Aspergillus niger*
- Riboflavin - *Ashbya gossypii*
- Stonewashed jeans - *Trichoderma reesei*
- Penicillin antibiotic - *Penicillium notatum*
- Button Mushrooms - *Agaricus bisporus*

Agricultural impact of fungi

Most of plant disease is caused by



USDA



A.G. Bölker

deposit mycotoxins - e.g. ergot

Some fungi provide nutrient
and nitrogen fixation

Impact of fungi on human health

- Mostly immunocompromised individuals are at risk of life-threatening infections
- Primary pathogens
 - *Histoplasma, Coccidioides, Cryptococcus gattii*
- Opportunistic pathogens
 - *Candida albicans, Aspergillus fumigatus, Cryptococcus neoformans, Rhizopus oryzae*

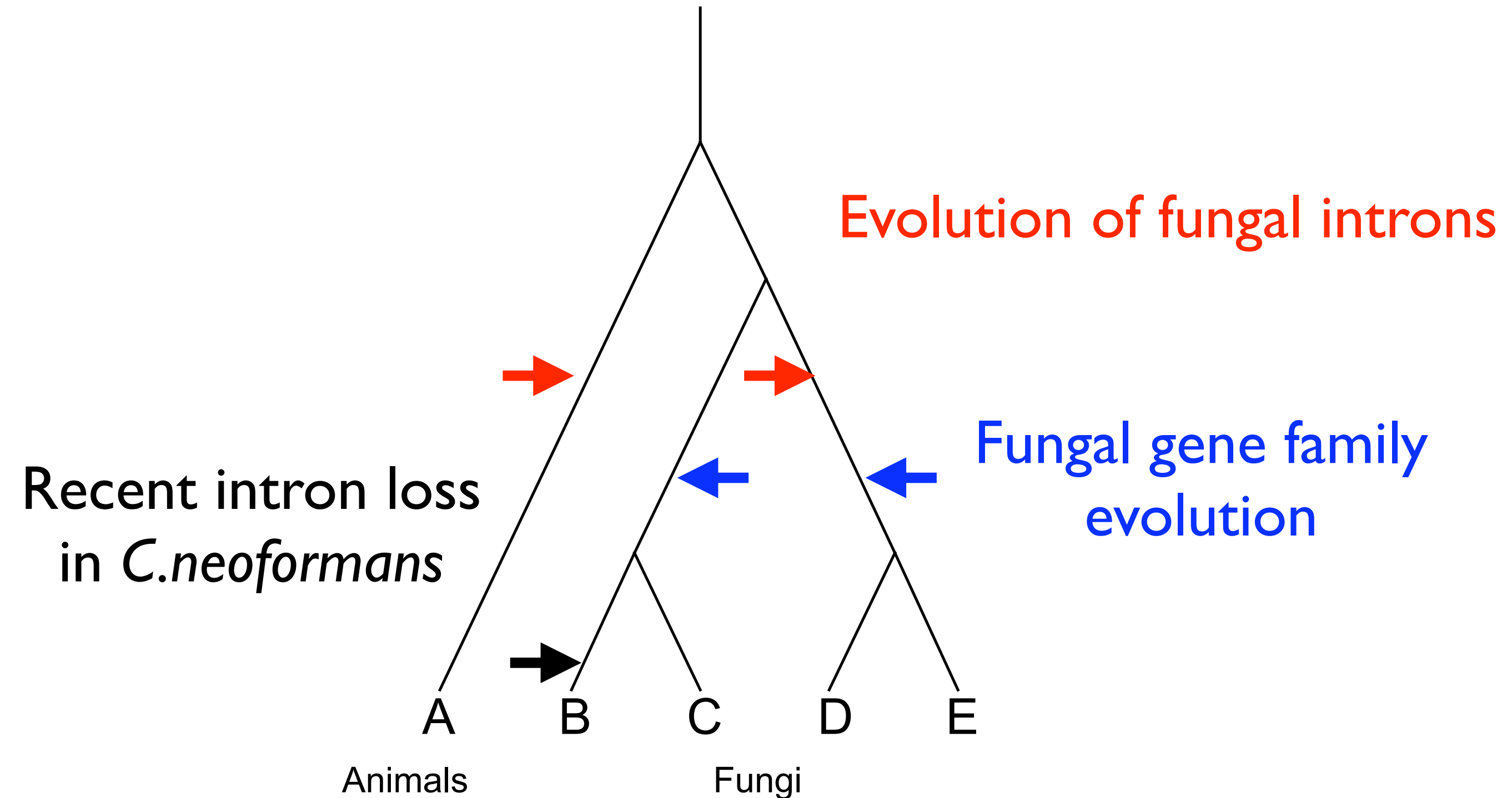
Fungi as genetic models

- Beadle and Tatum (1941) - one gene, one enzyme hypothesis in *Neurospora crassa*
- Cell cycle, cell model - *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe*
- Molecular biology tools to investigate phenotype-genotype
- Evolutionary models

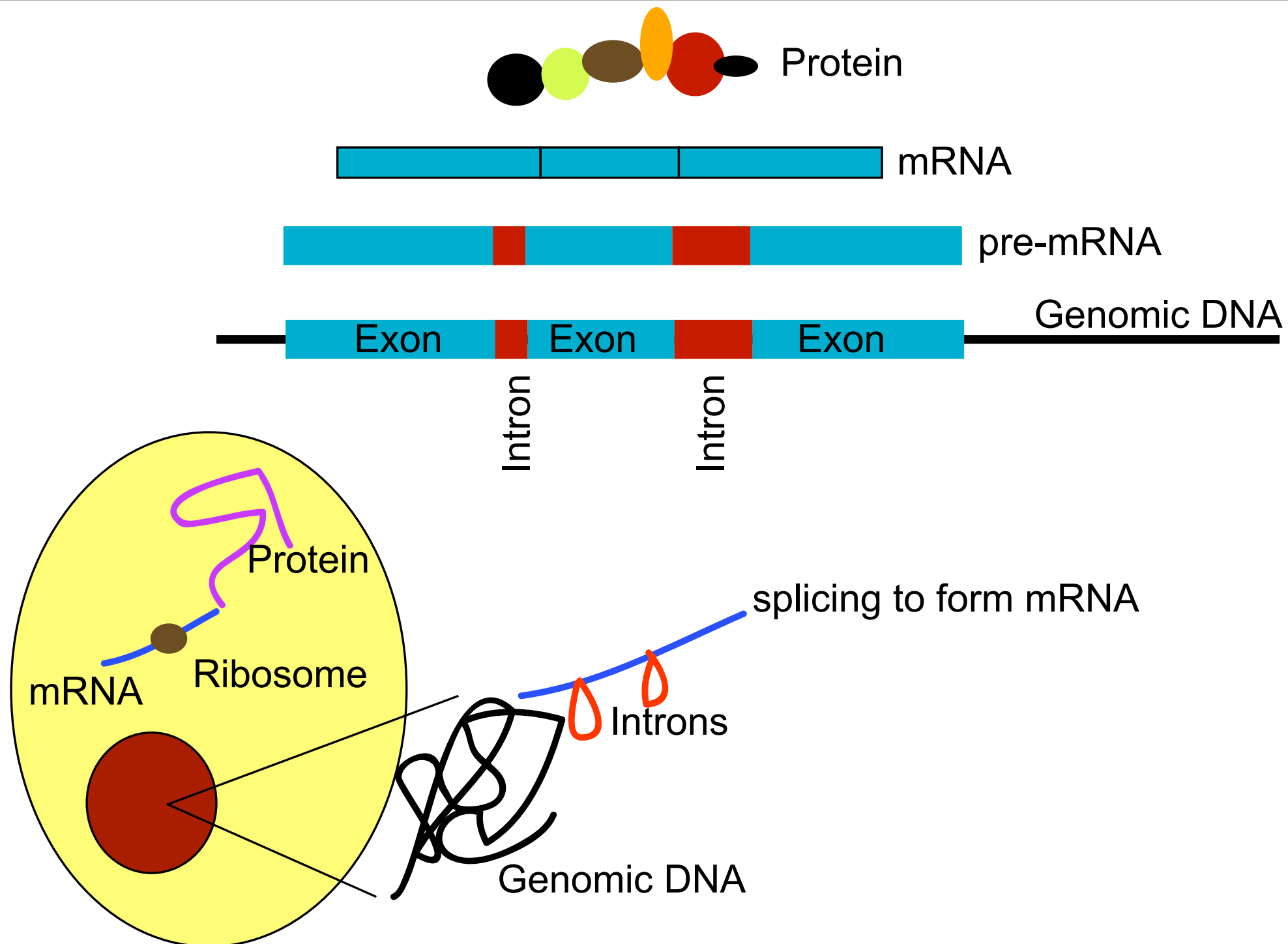
Fungal genomes

- Smaller than most vertebrate and plant genomes
 - *A. gossypii* 8.5 Mb; *S. cerevisiae* 12 Mb
 - *N. crassa* 40 Mb
 - Animals: 100 Mb worm; 3000 Mb Human
- Vary in protein coding gene content
 - 4,700 in *A. gossypii*; 5,800 in *S. cerevisiae*
 - 16,000 in *R. oryzae* or *S. nodorum*
 - 19,000 in Fruitfly; 25,000 in worm

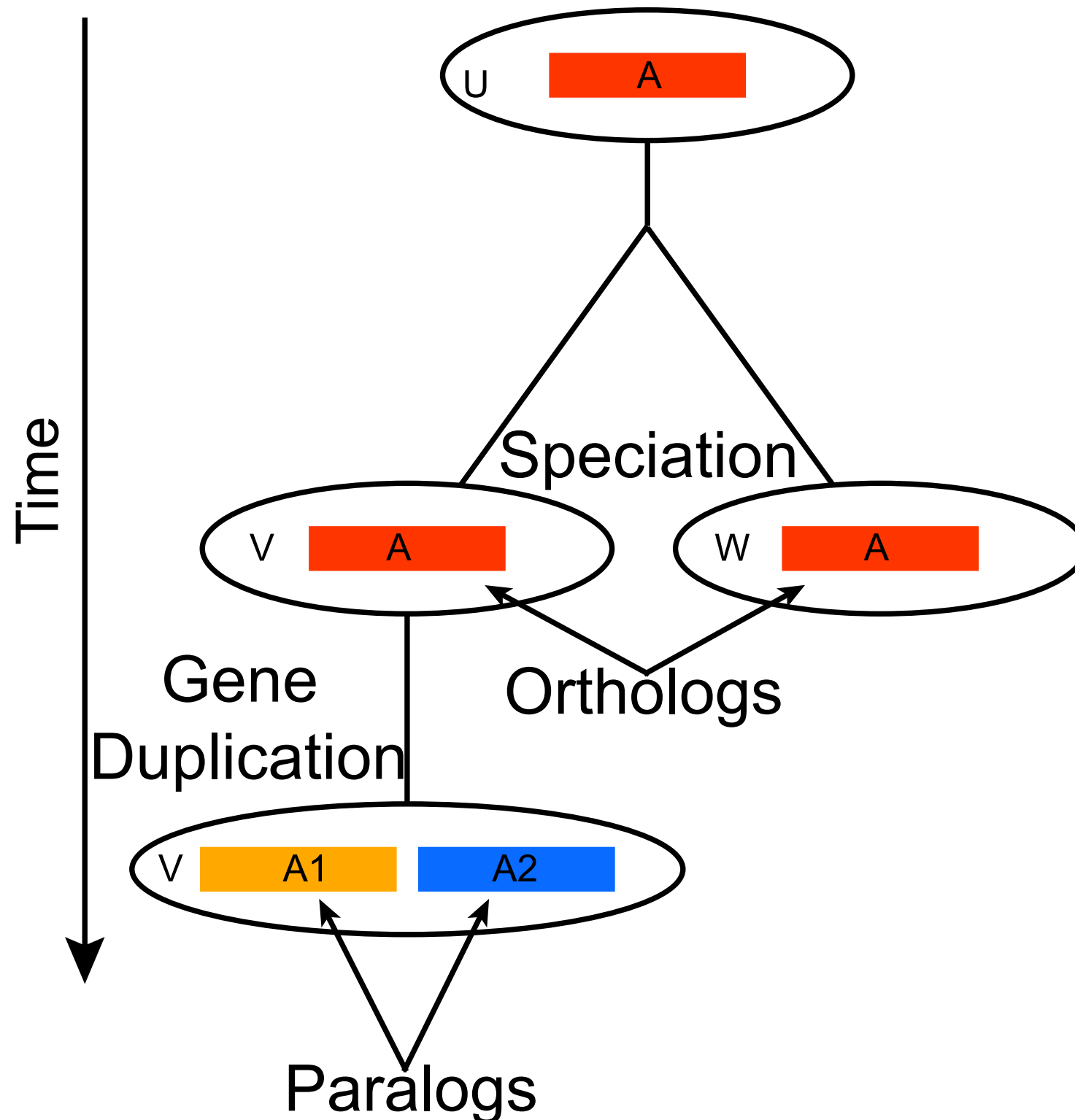
Fungal comparative genomics



Central dogma of eukaryotic biology

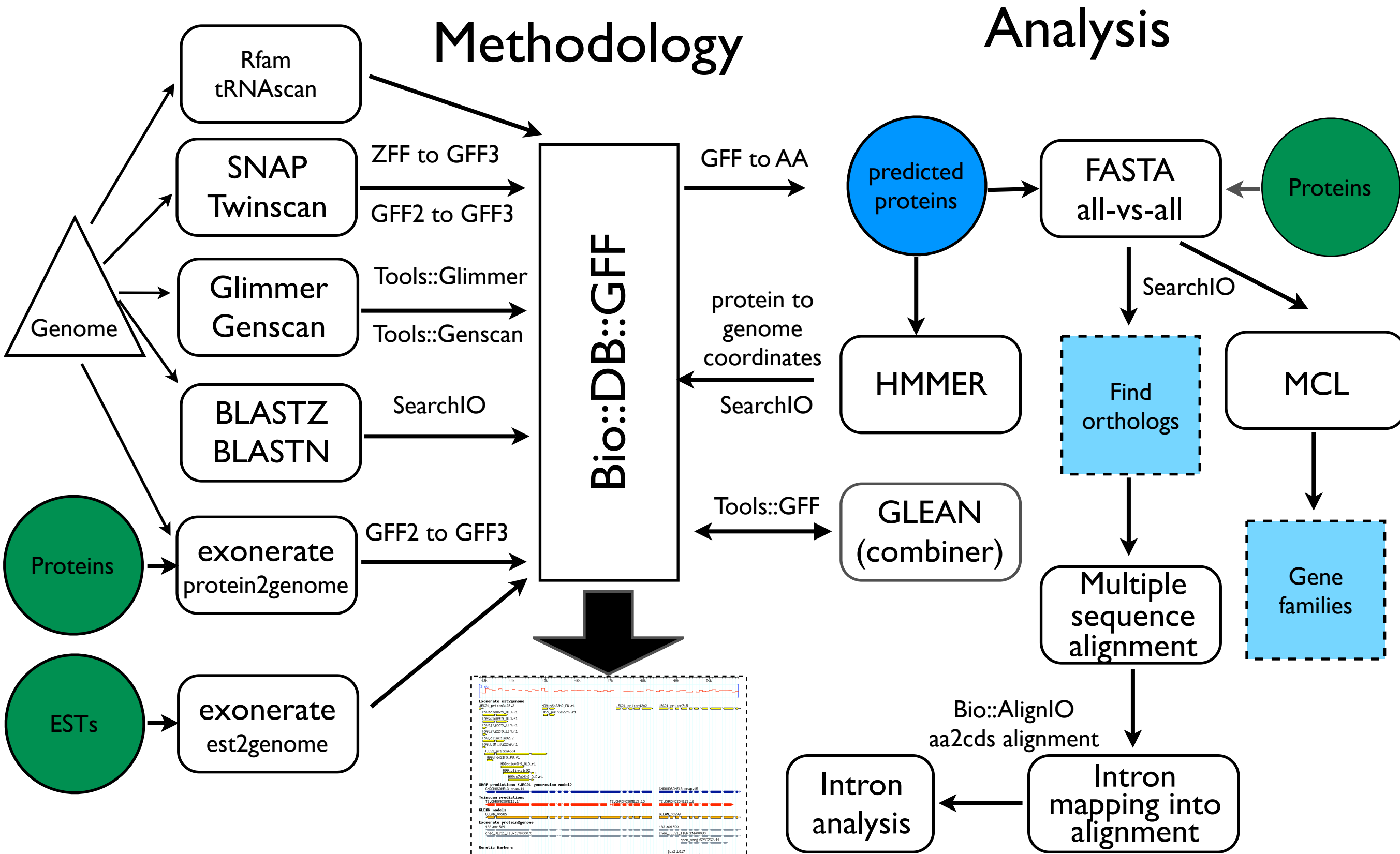


Orthologs and Paralog



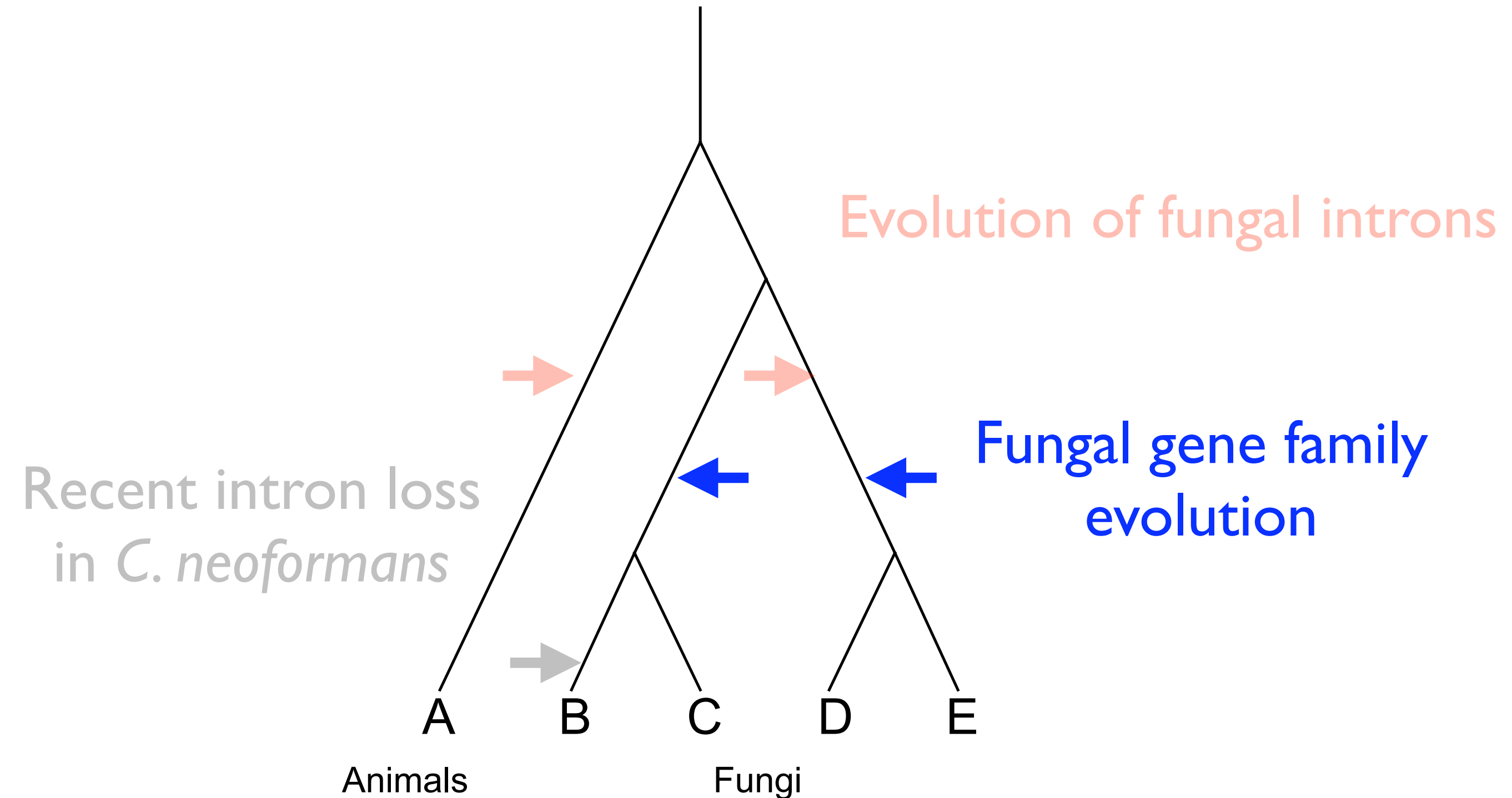
Genome annotation

- Available fungal genomes were only assembled genomic sequence.
- Need systematic and consistent gene predictions for genome comparisons
- Automated annotation pipeline for gene prediction.



<http://fungal.genome.duke.edu>

Fungal comparative genomics



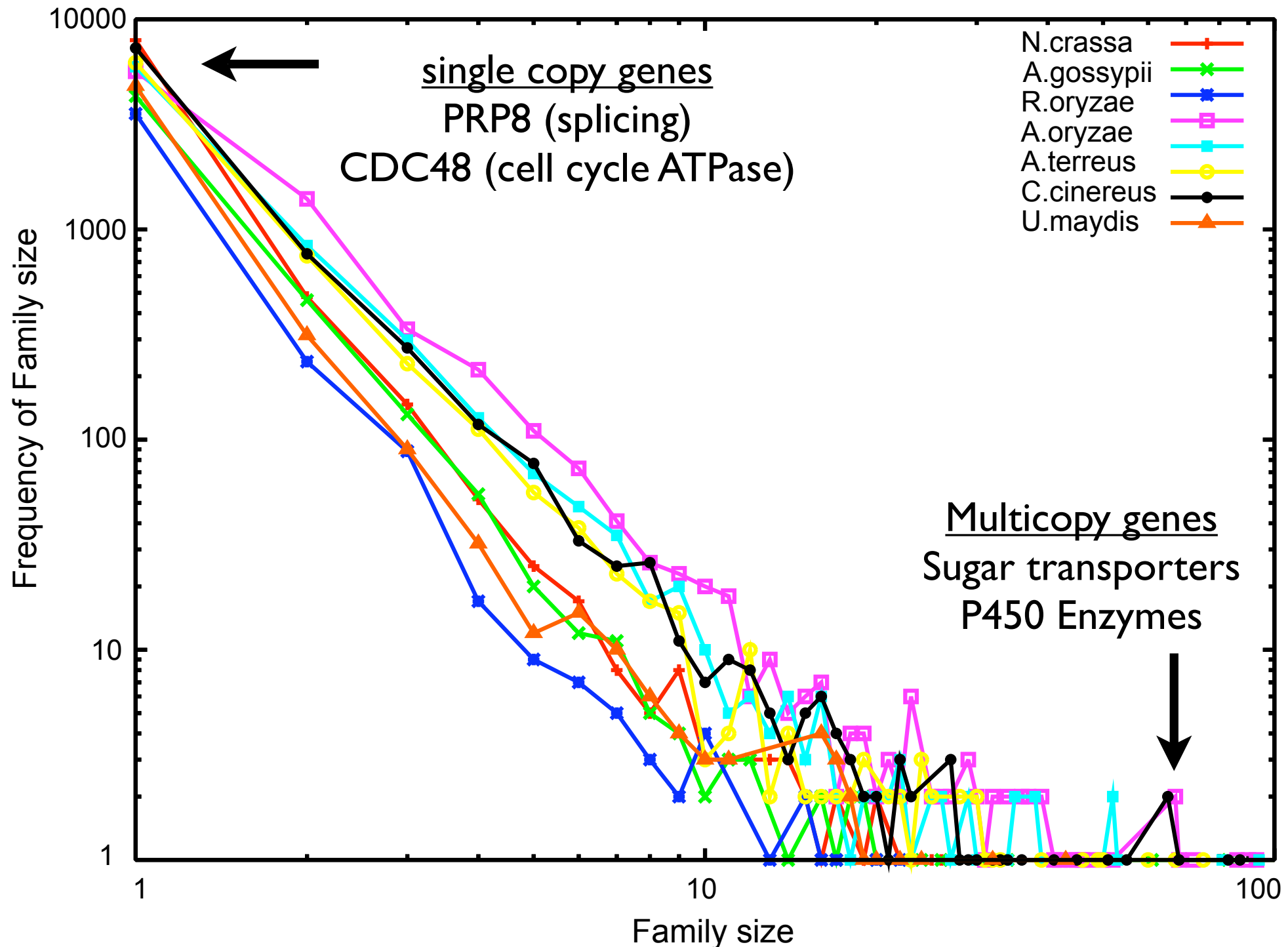
Gene family evolution

- Gene families are the crucible of new genes and thus new functions
- Signature of adaptive evolution often confounded in multi-gene families
- Can we identify families that have unexpectedly large changes in size across a phylogeny?
- Follow up these families with more focused studies

Identifying gene family expansions

- Previous work only considered pairwise
- *Ad hoc* comparison of gene family sizes
 - *C.elegans-C.briggsae* - GPCR family expansions (Stein et al, *PLOS Biology* 2004)
 - *A. gambiae-D. melanogaster* - Mosquito specific family expansions related to symbiotic bacteria (Holt et al, *Science* 2002).

Gene family sizes follow power law distribution



Phylogenetic evaluation of gene family size change

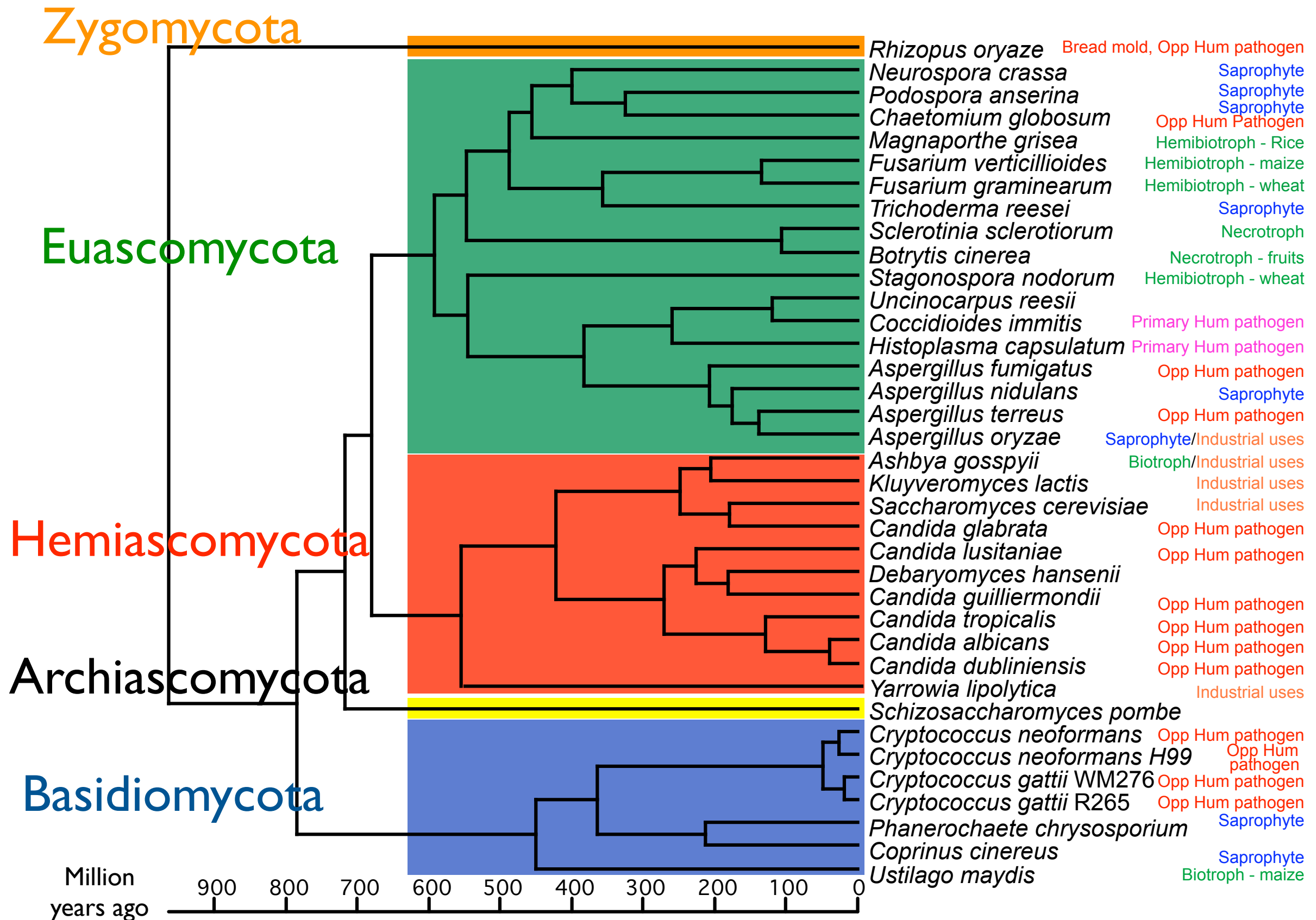
- Previous methods only used *ad hoc* statistics
- Explicit model for gene family size change according to a Birth-Death (BD) models
- Apply BD to family size along phylogeny using probabilistic graph models
- CAFE - Computational Analysis of gene Family Evolution

Hahn et al, *Genome Res* 2005

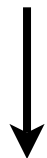
De Bie, et al *Bioinformatics* 2006

Demuth et al, *submitted*

Fully sequenced fungal genomes



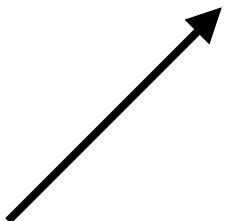
FASTA
all-vs-all



MCL



Gene
families

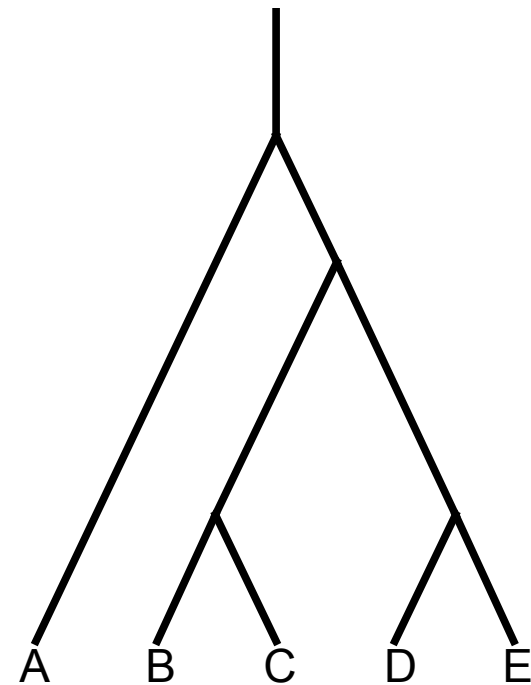


37 fungal species

Family count

	10	1	2
	14	18	2
	7	1	1
	6	1	12
	6	1	8
	3	1	1

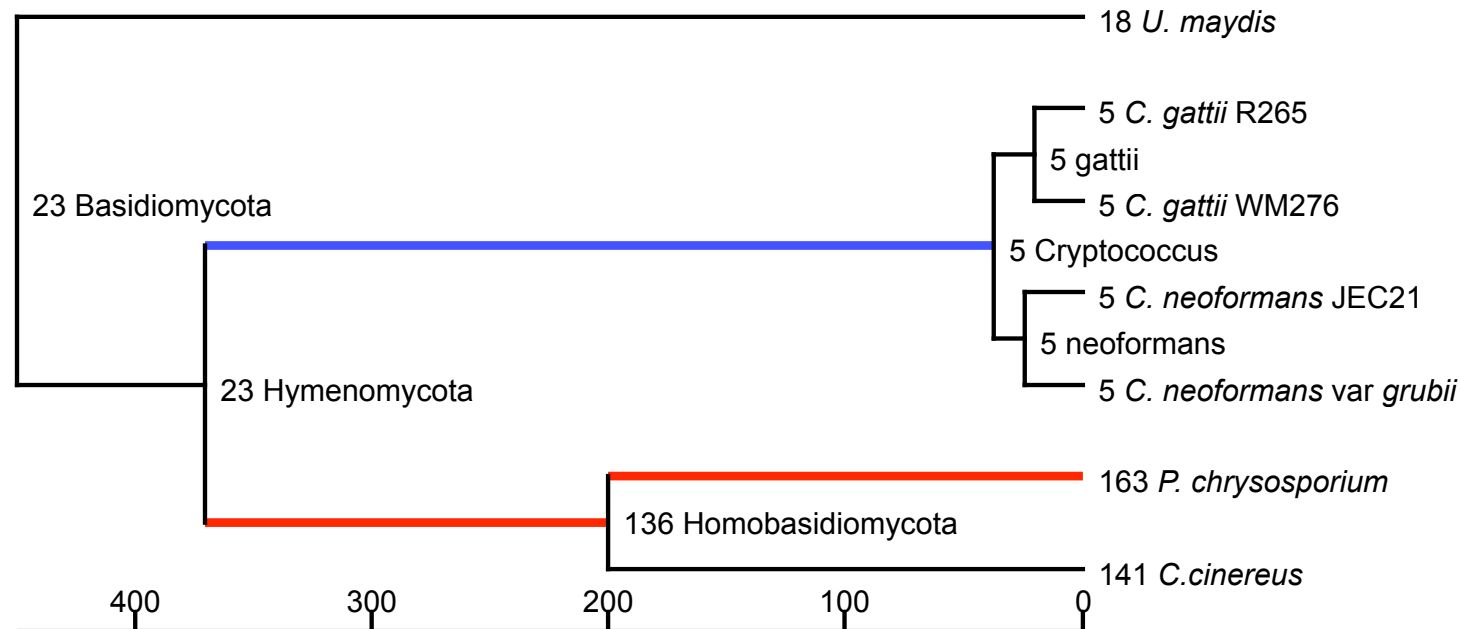
+



CAFE



Family 1	$P < 0.001$	Branch A
Family 2	$P < 0.001$	Branch B
Family 3	$P=0.02$	Branch C,E
Family 4	$P=0.03$	Branch D



Families with significant expansions

Transporters
Kinases
P450
Oxidation

Vitamin & Cofactor transport

Lactose & sugar transport

Amine transport

Myo-inositol, quinate, and
glucose transport

Oligopeptide transport

ABC transporter

MFS, drug pump, & sugar
transport

Transport

Monocarboxylate & sugar
transport

ABC transport

Amino acid permease

Methyltransferase

Cytochrome P450: CYP64

Cytochrome P450: CYP53,57A

Cytochrome P450

Kinase

Subtilase family

NADH flavin oxidoreductase

Aldehyde dehydrogenase

Aldo/keto reductase

Multicopper oxidase

AMP-binding enzyme

Transporters

- Of 45 significant families, 22 were related to transport
- Vitamin and amino acid transport
- Sugar and sugar-like transporters
- Multidrug and efflux pumps
- ABC transporters (ATP Binding Cassette)

Branches with transporter expansions

- Sugar related, Drug pump, and Major Facilitator Superfamily

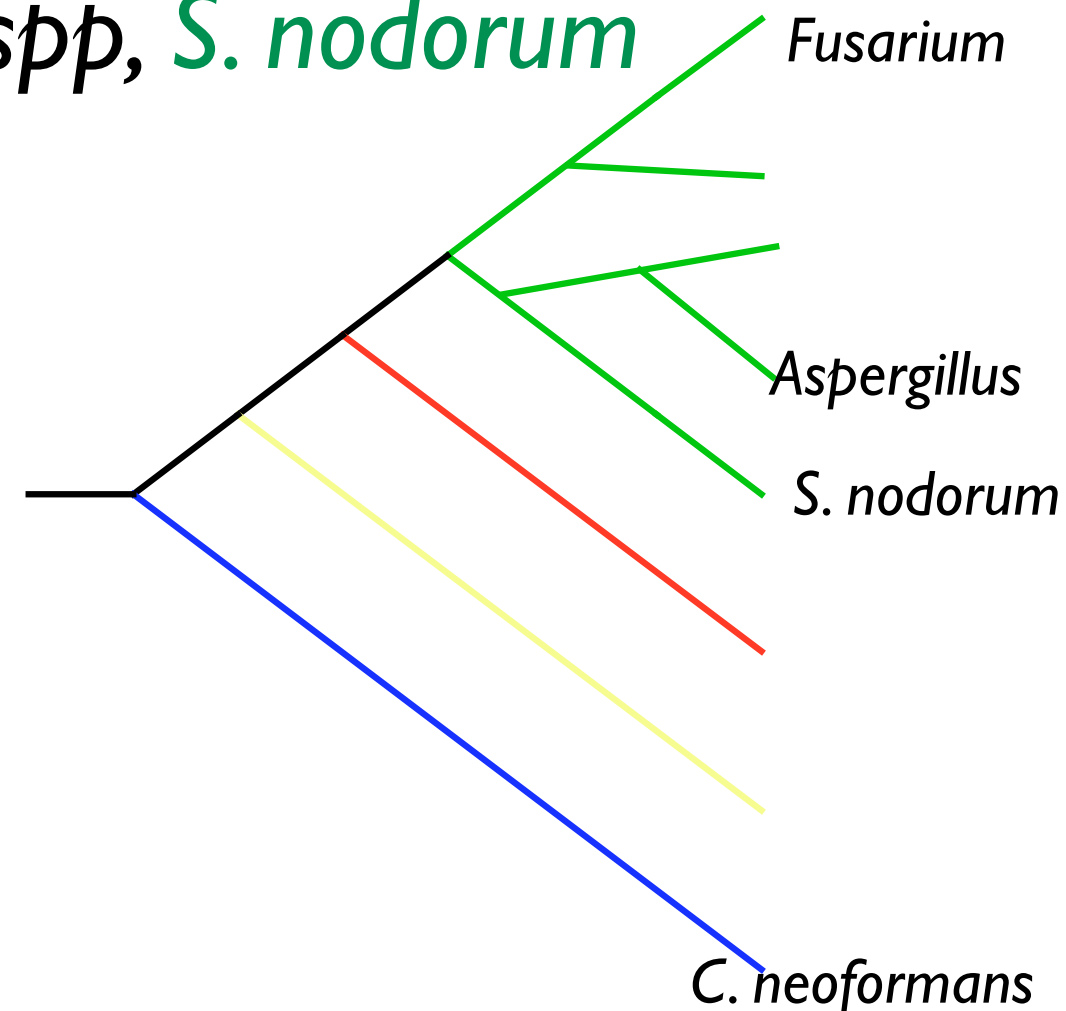
- *Aspergillus* spp, *Fusarium* spp, *S. nodorum*

- *Euscomycota*

- Vitamin transport

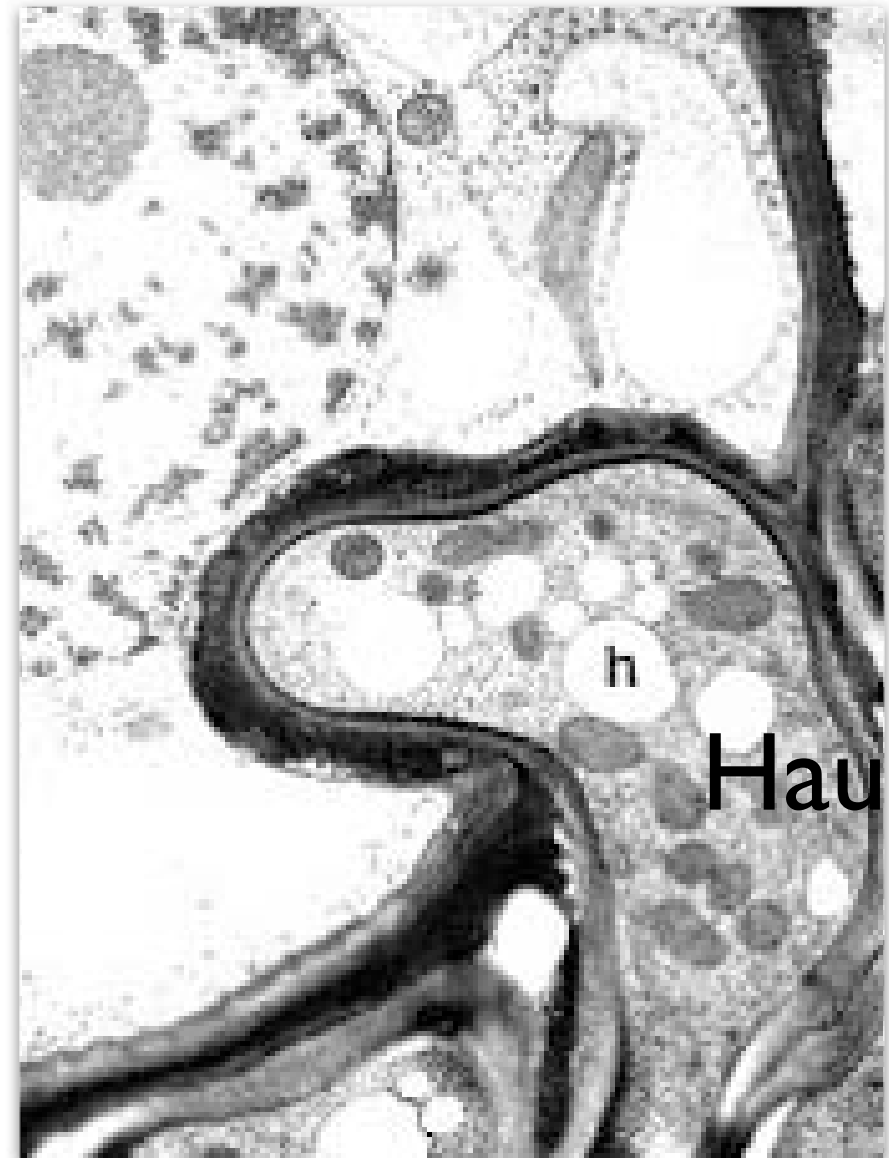
- *C. neoformans*, *Fusarium*

- *A. nidulans* (Biotin)



Sugar transporter use in phytopathogens

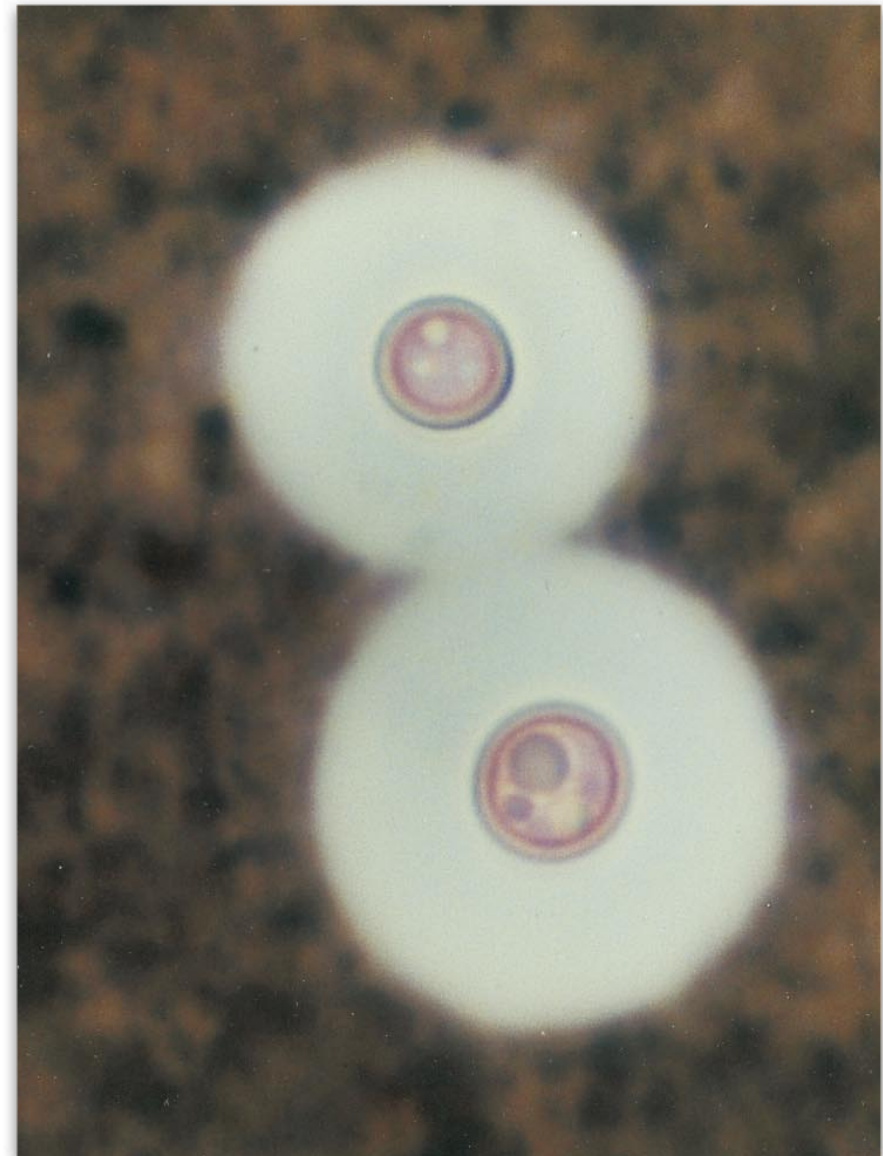
- Sugar transporters are used to extract nutrients from host
- Haustorium: specialized structure for plant parasitism
- Many sugar transporters highly and specifically expressed in haustoria



Robert Bauer <http://tolweb.org/>

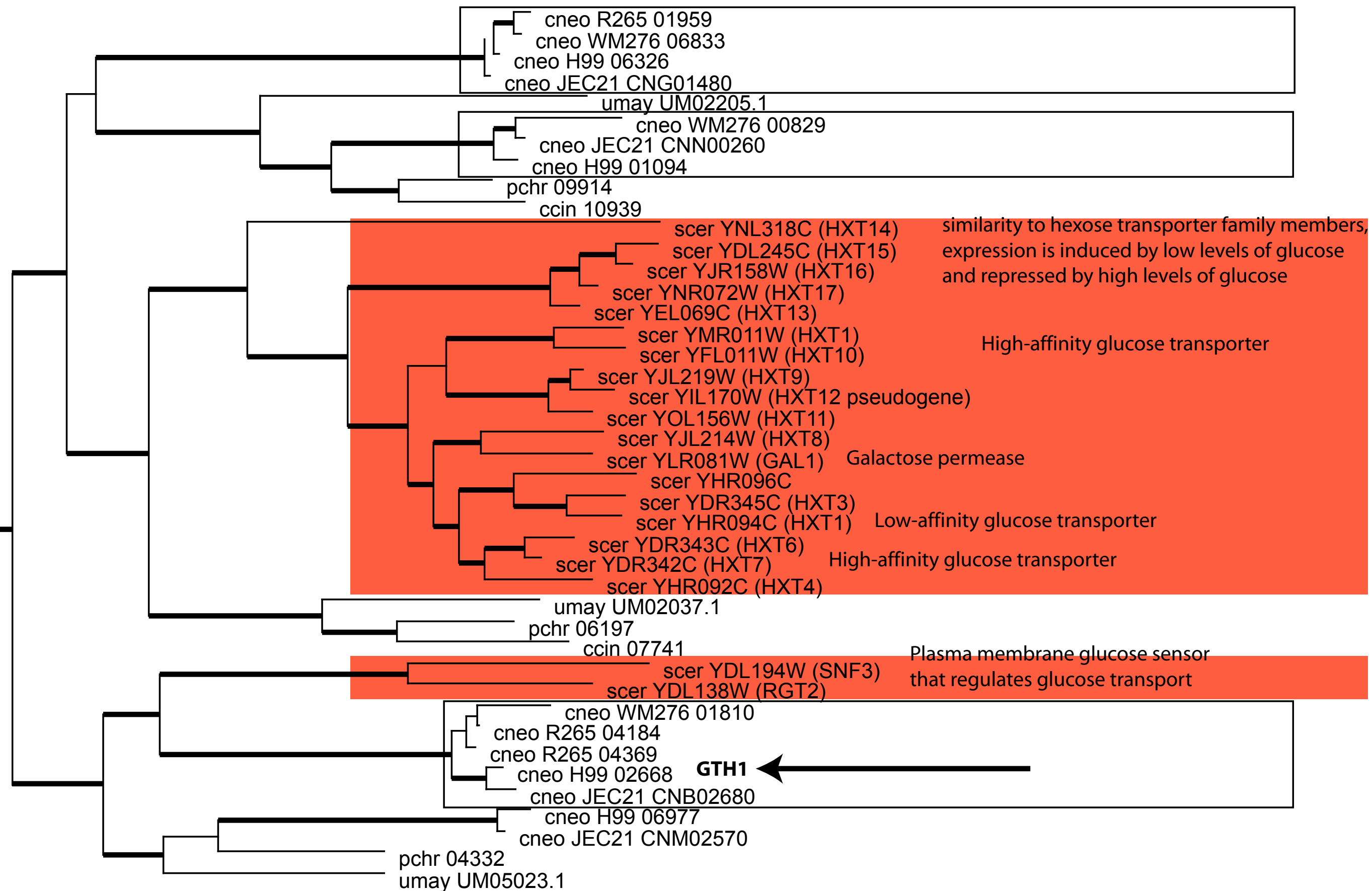
Cryptococcus sugar transporters

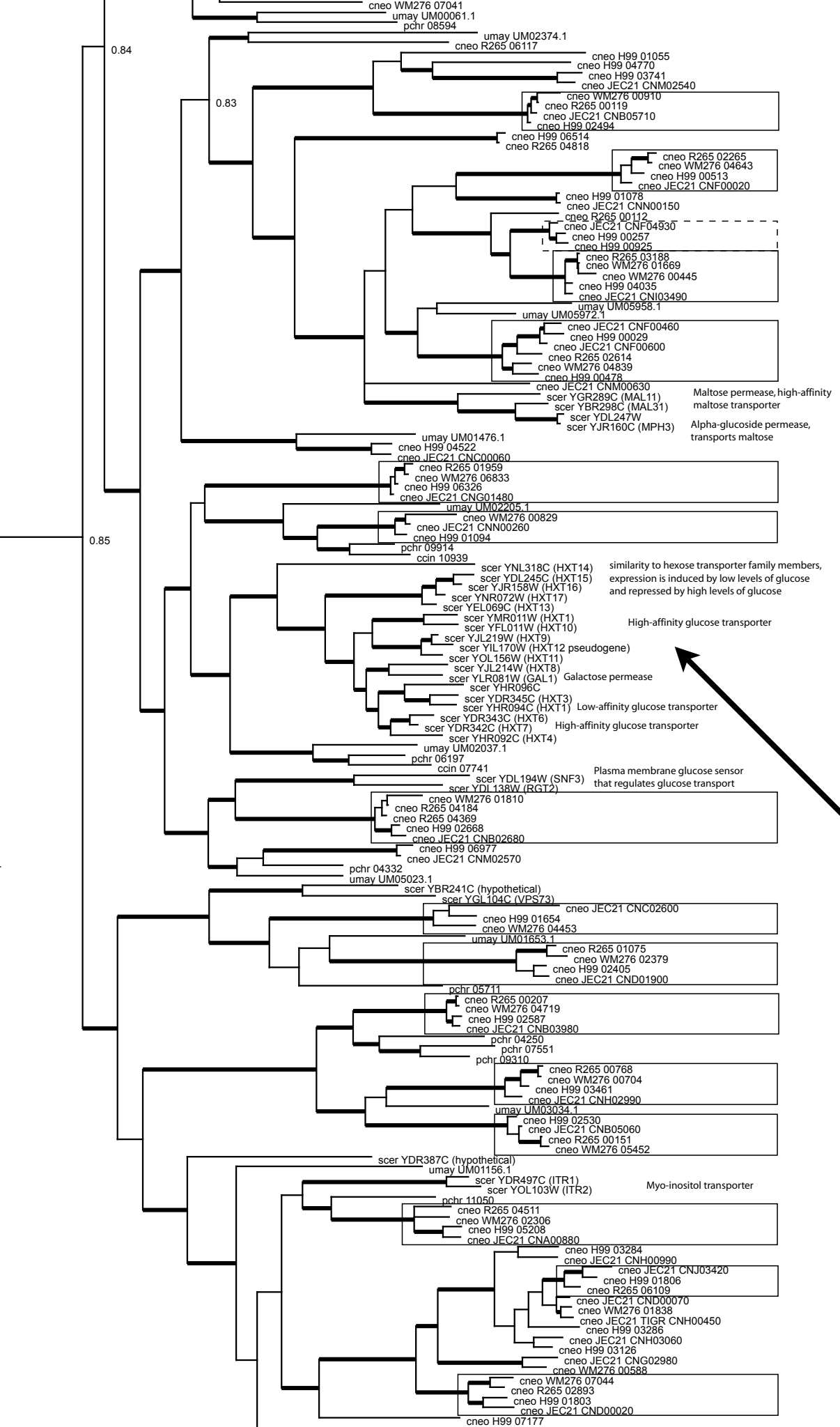
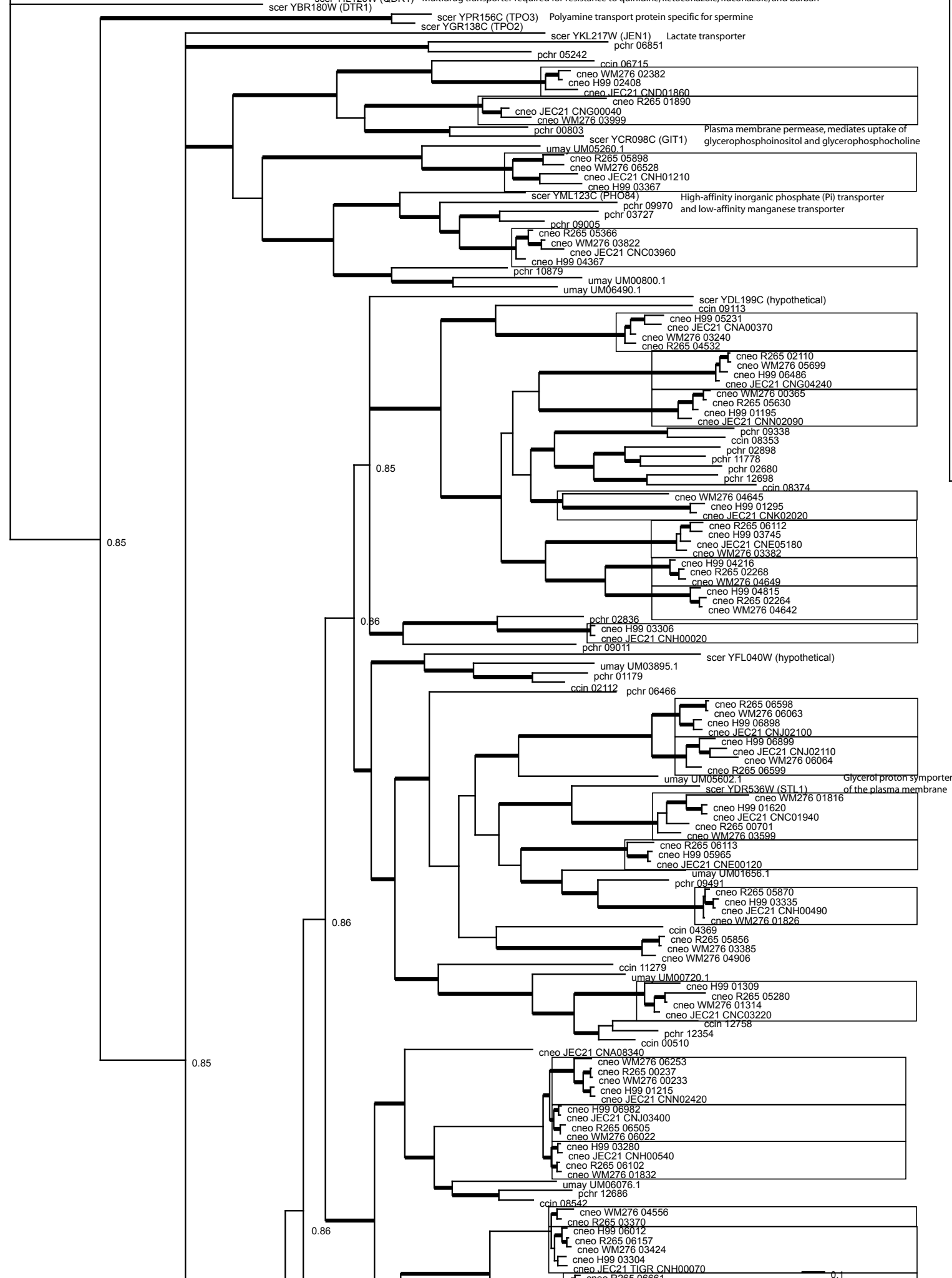
- 3x as many sugar transporters in *C. neoformans* (~50) than other basidiomycetes
- “sugar coated killer”
- Capsule is a mixture of glucose, xylose, and mannose.
- Transporters could be important in capsule synthesis



Zerpa et al, 1996

Analysis of sugar transporter sub-family





GTH I is a Glucose Sensor

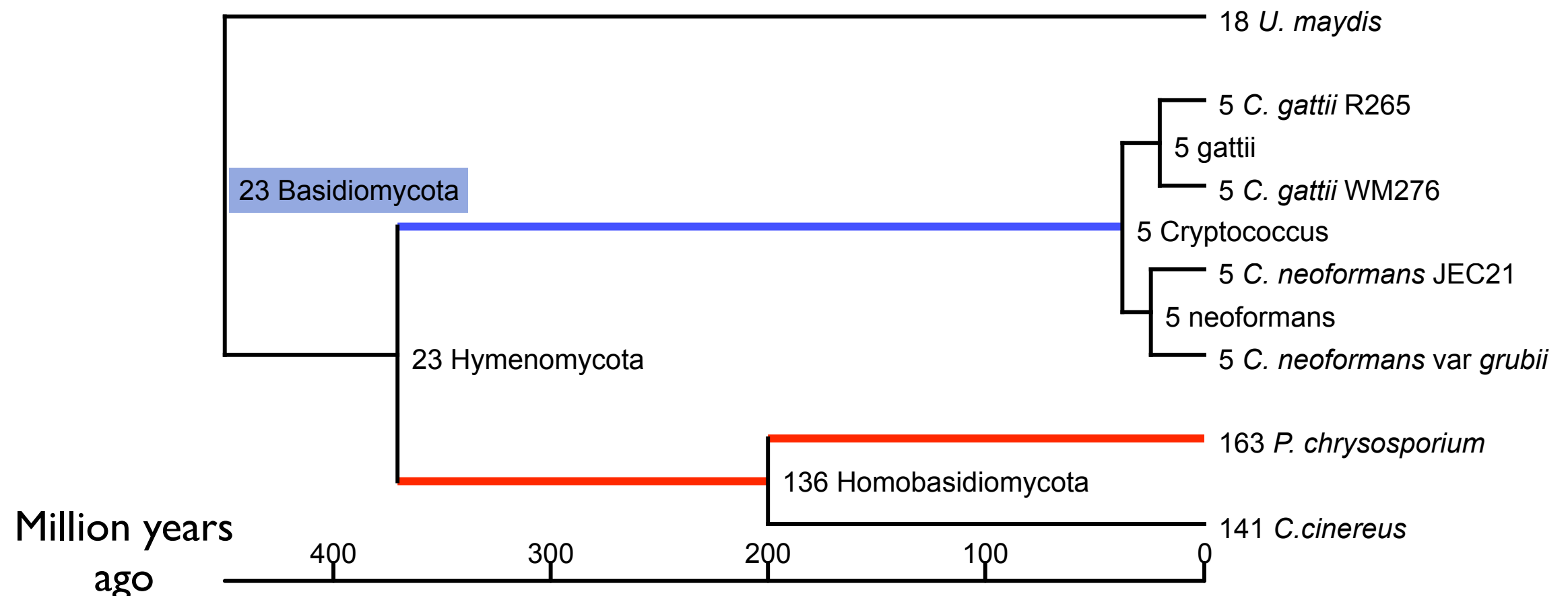
- Melanin production pathway is starvation induced
- GTH I KO is hypermelanized (always glucose starved)
- GTH I overexpression is hypomelanized

Alsbaugh et al 1999

P450 CYP64

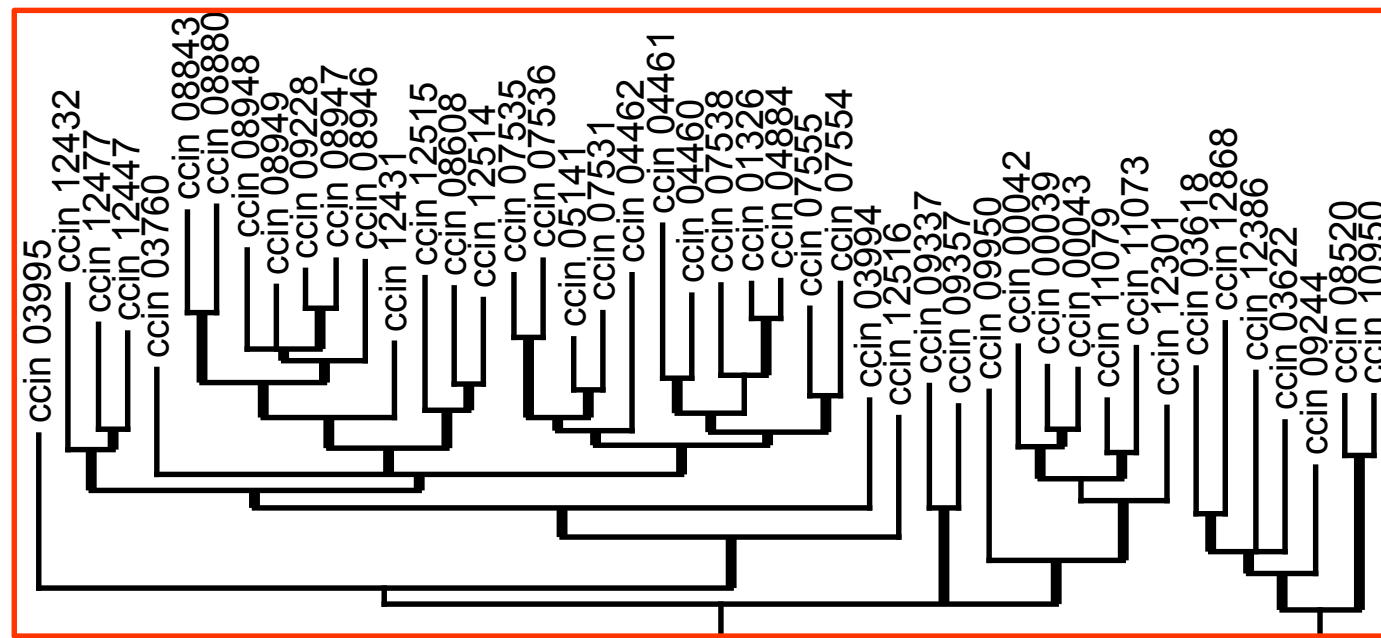
P450 enzymes involved in synthesis and cleavage of chemical bonds. Drug metabolism in animals.

CYP64: Step in *Aspergillus* spp aflatoxin pathway
P. chrysosporium implicated in lignin and hydrocarbon degradation.

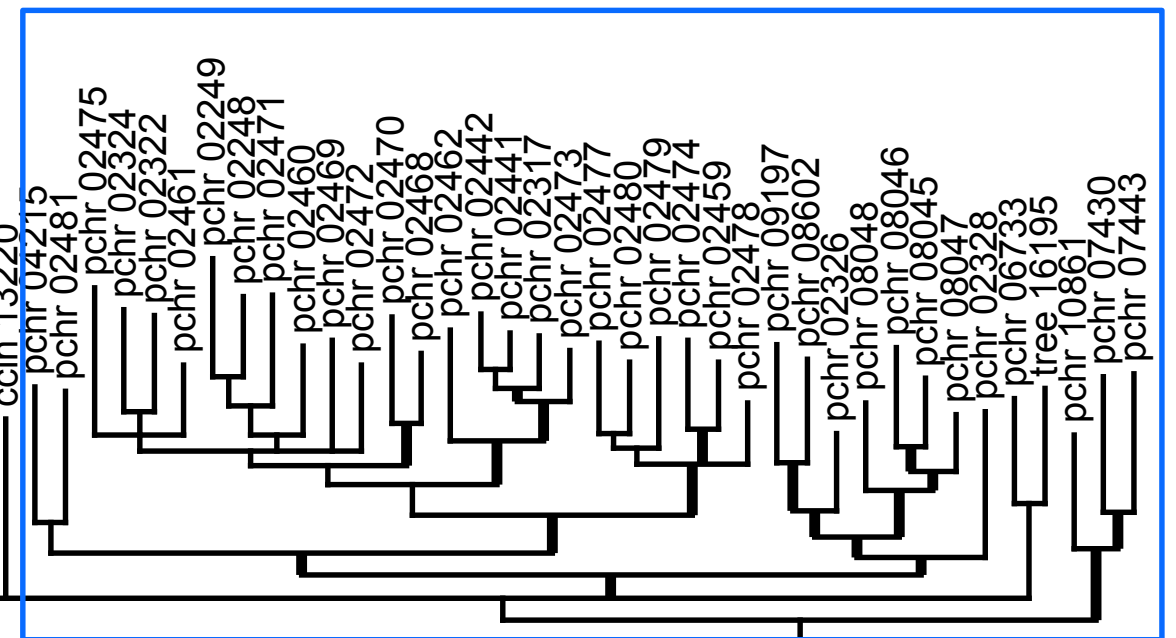


CYP64 was from independent duplication

C. cinereus expansion



P. chrysosporium expansion

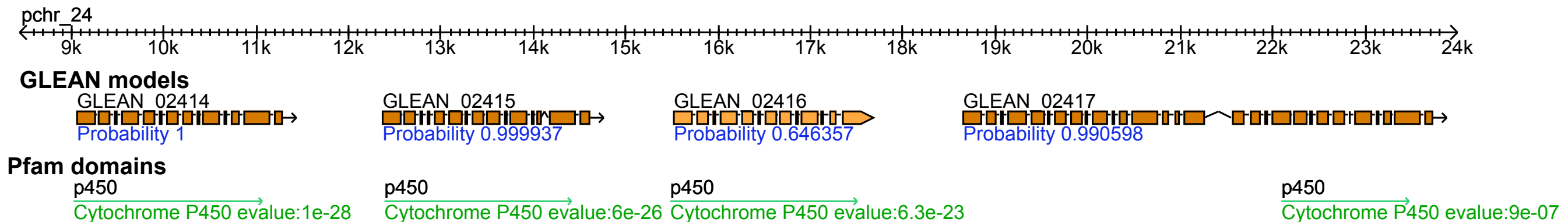


Mario Cervini



Tom Volk

Local duplications created CYP64 expansion



Family size contractions

- *Histoplasma, Coccidioides* many families
- *Hemiascomycetes* - P450
- *C. neoformans* - P450
- *U. maydis* - Lactose transport

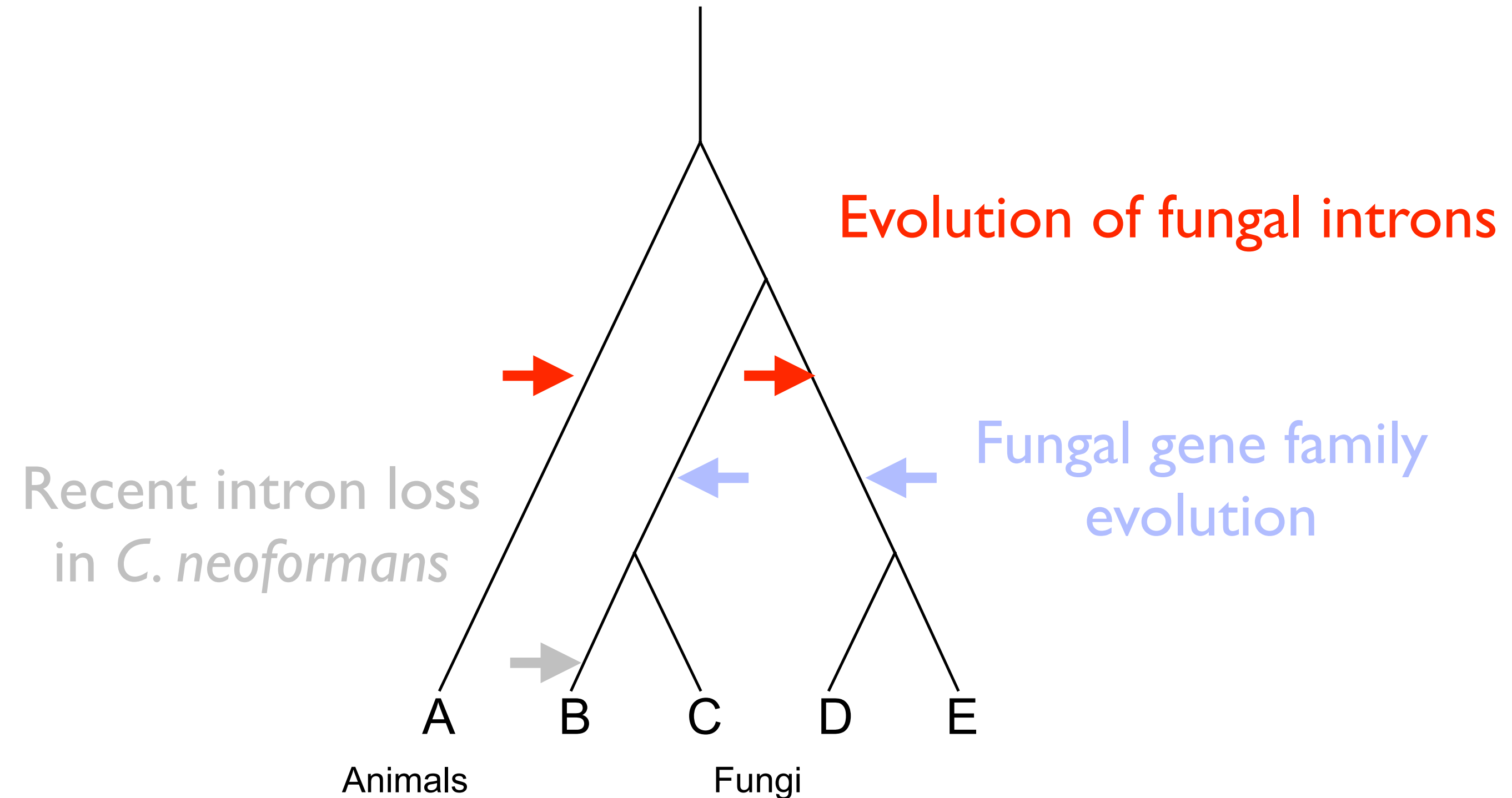
Basidiomycetes

fam	U.maydis	C.cinereus	P.chrysosporium	H99	JEC21	R265	WM276	Cryptococcus & GO Notes	Protein domains
0	1	420	68	140	4	14	70	transposon	Integrase core domain(57);
1	74	108	119	78	81	84	82	kinases. homobasidio	Protein kinase domain(306);
2	44	67	87	48	50	49	52	G-protein coupled receptor protein	WD domain, G-beta repeat(193)
3	42	113	53	27	4	23	106	transposon	Reverse transcriptase
4	18	141	163	5	5	5	5	metal ion binding; cation binding	Cytochrome P450 (39)
5	32	26	76	44	47	46	45	drug transport	Major Facilitator Superfamily(157);
6	17	36	44	47	55	40	39	"hexose transporter	dicarboxylate symporter
7	19	20	23	59	57	47	50	sugar transporter related	Sugar (and other) transporter(167);
8	30	33	64	32	33	31	35	"siderophore biosynthesis?	- dihydroxybenzoat
9	35	39	35	33	40	34	36	transmembrane receptor activity;	Ras family(130); ADP-ribosylation
10	32	44	45	29	27	31	30	localization hydrolase	ABC transporter (112); RecF/RecN/
11	5	125	61	11	11	11	13	hydrolyzing O-glycosyl	WSC domain (24); Glycosyl
18	4	2	5	46	43	22	30	oxidoreductase activity;RNA	Oxidoreductase family, NAD-
20	0	109	36	0	0	0	0	kinase activity; phosphotransferas	

Conclusions

- Sugar transporters are highly expanded in independent lineages
 - Saprophytic and phytopathogenic lifestyles
- P450 CYP64 independent expansions in *Homobasidiomycetes*
 - Lignin degradation and saprophytic lifestyles
- Family size contractions among lineages containing primary pathogens
 - Genome streamlining?

Fungal comparative genomics

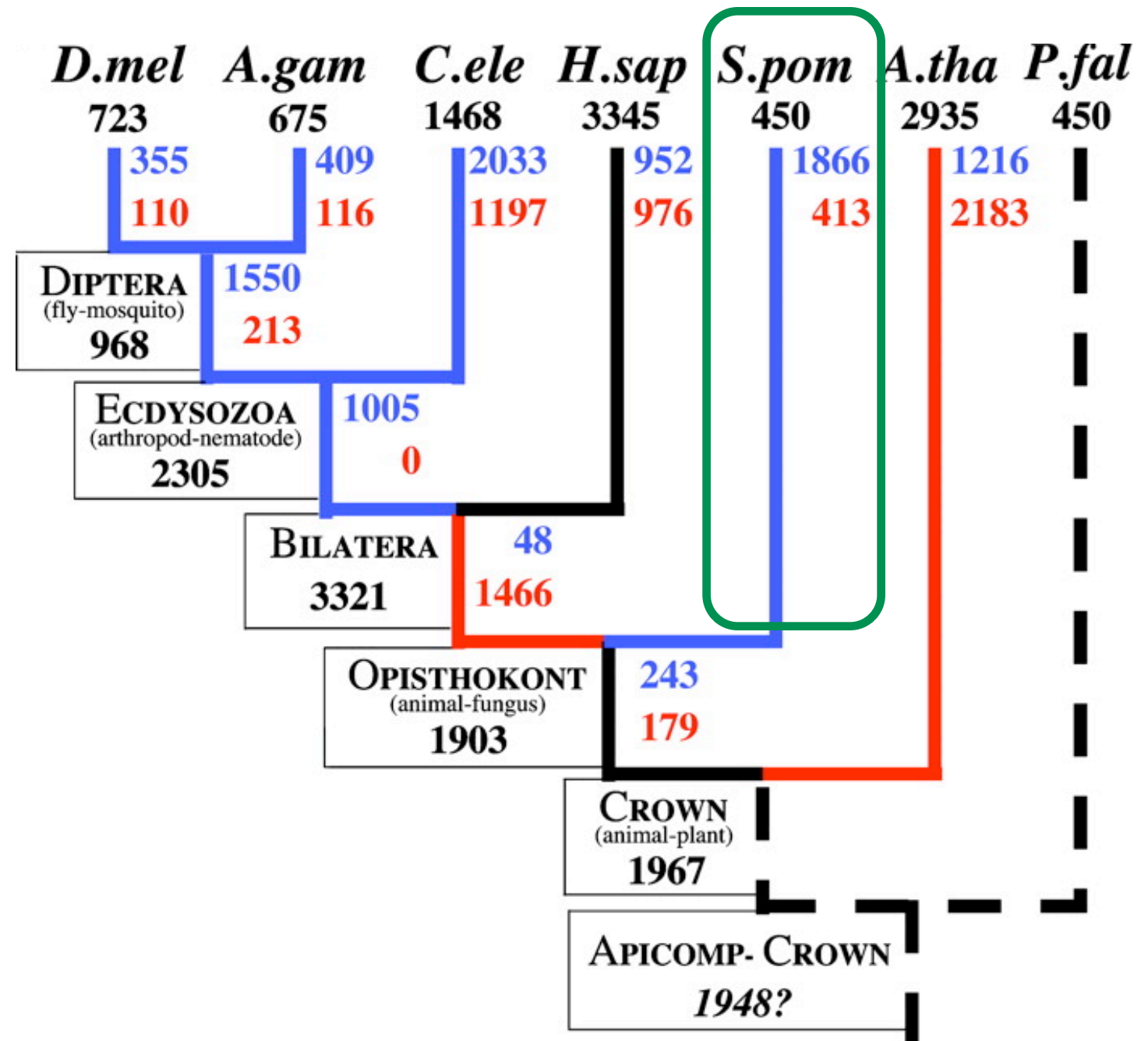


Evolution of gene structure

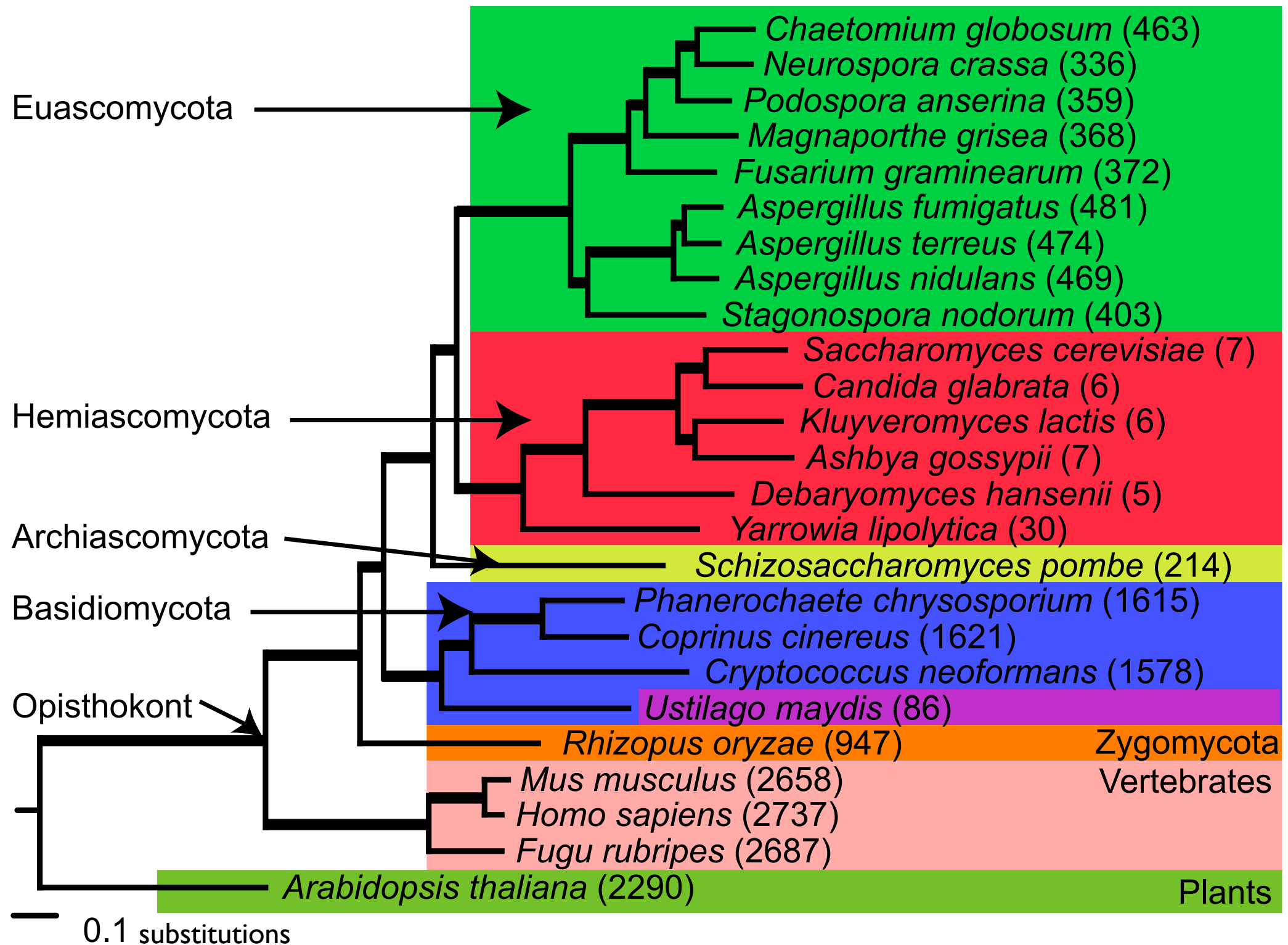
- Present day introns
 - Recent insertions?
 - Introns late hypothesis
- Present in eukaryotic ancestor?
 - Introns early hypothesis / exon theory of genes
- Mixture of two?

Previous work on intron evolution

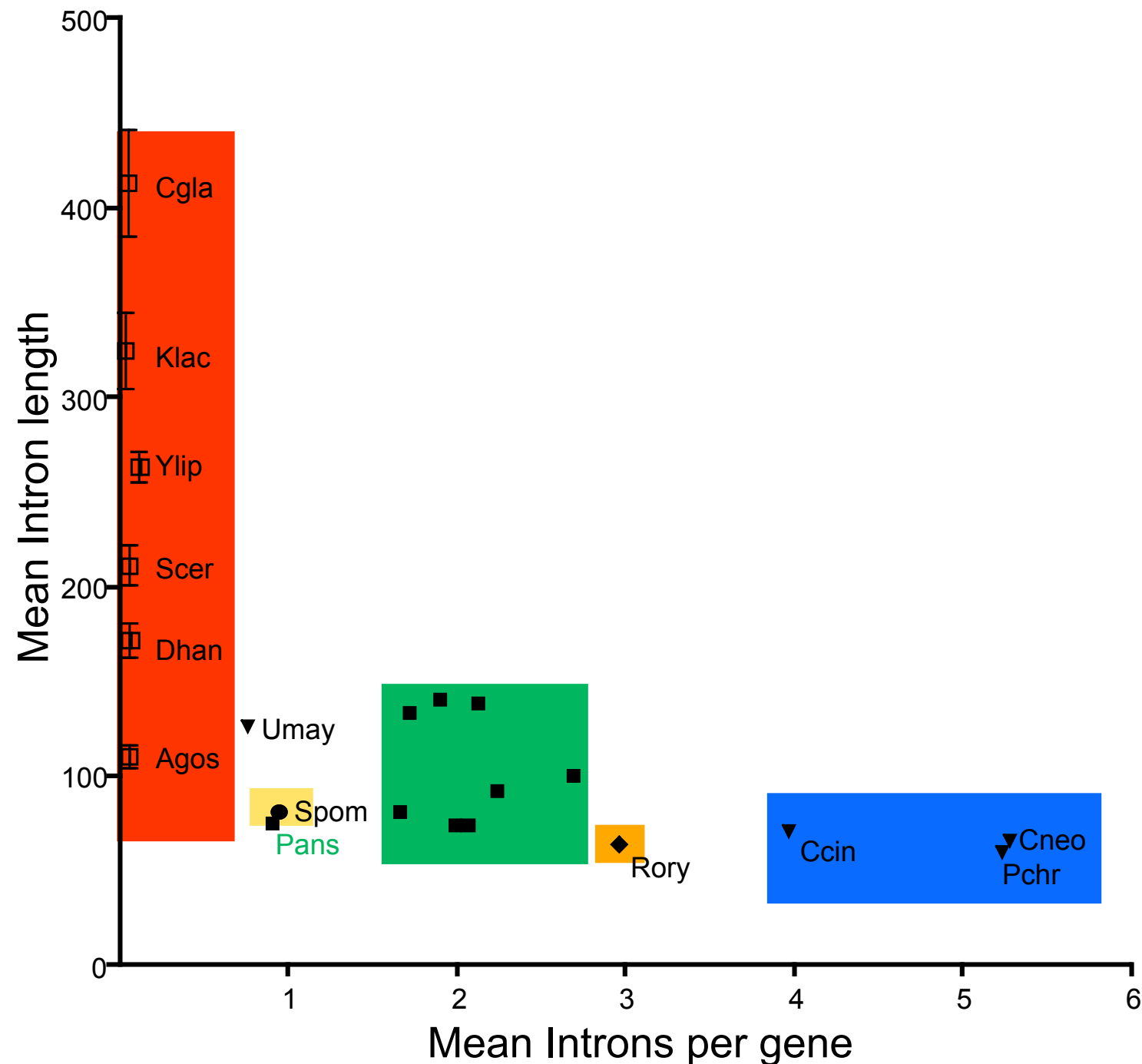
- Rogozin et al. 2003
 - 7 genomes
 - 684 genes, 7236 positions
 - Parsimony analysis
- Analysis methods
 - Roy and Gilbert. 2005
 - Csűrös. 2005
 - Nguyen et al. 2006



Calculating intron densities across a phylogeny



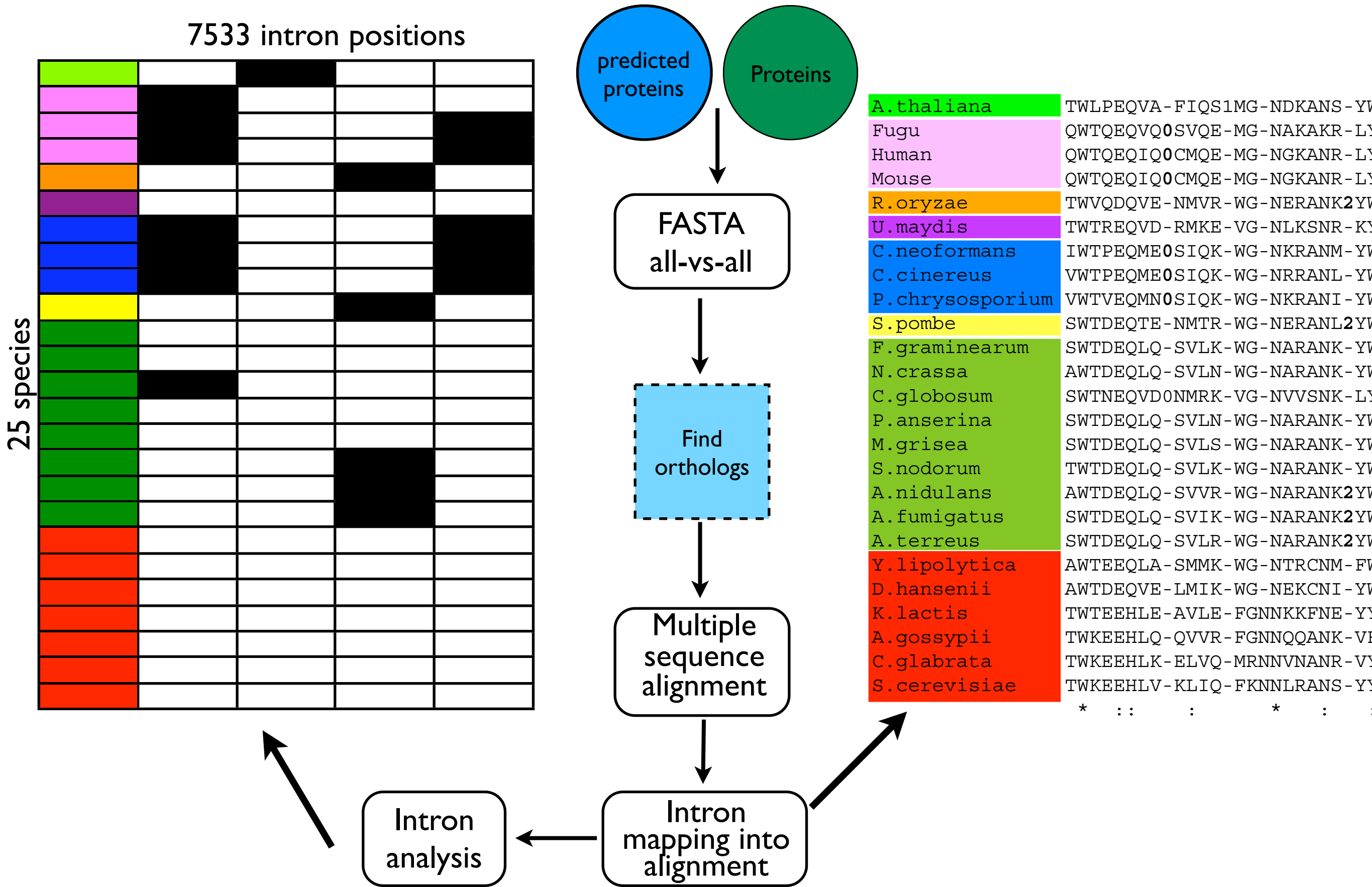
Intron frequency varies among the fungi



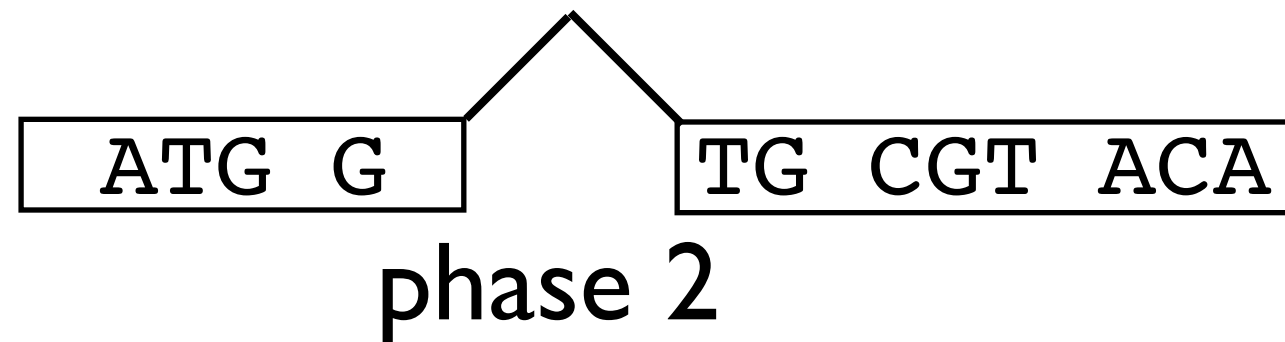
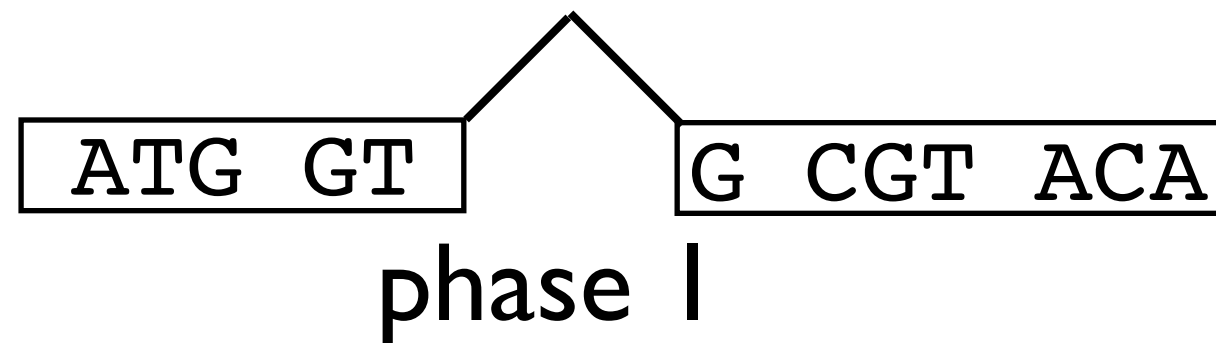
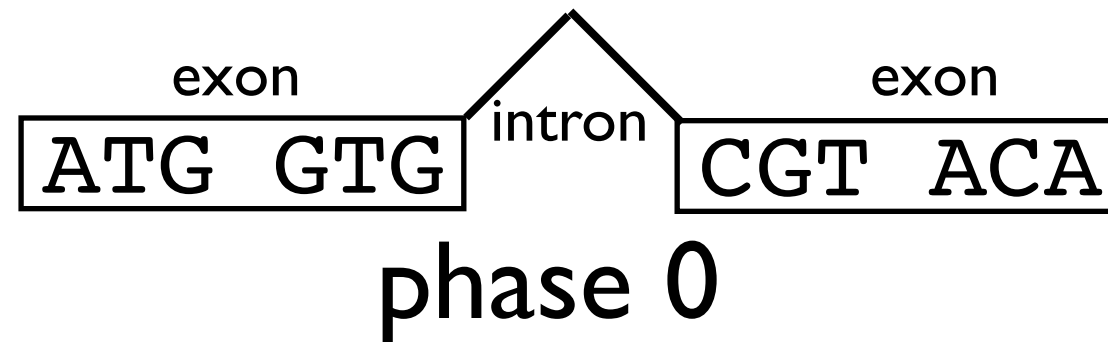
Analysis of whole genomes

- 25 entire genomes
 - 21 fungi, 3 vertebrates, 1 plant
- Largest dataset ever assembled for intron analysis
- 1160 orthologous genes
- 7533 intron positions
- 4.15 Mb conserved coding sequence per genome

Analysis methods



Intron phase



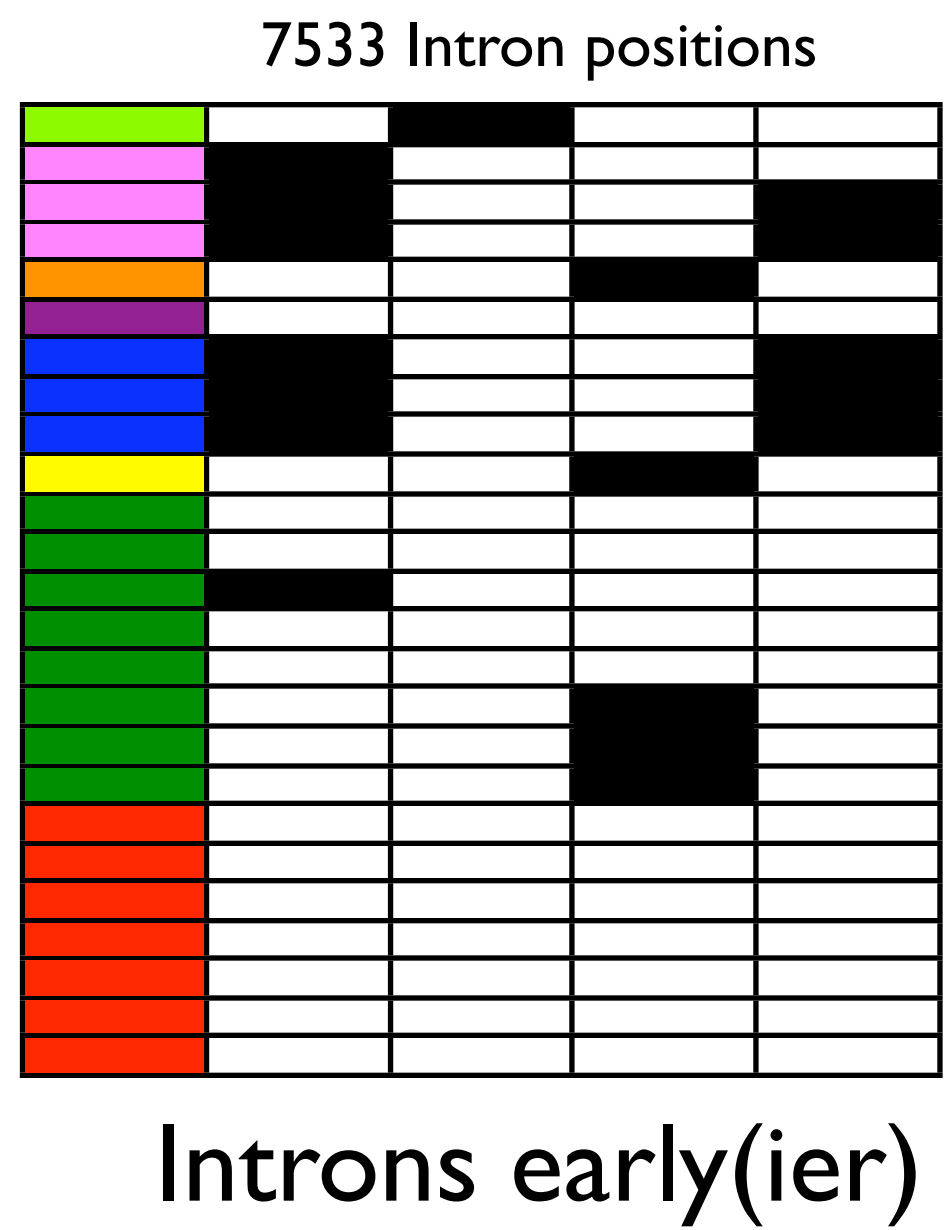
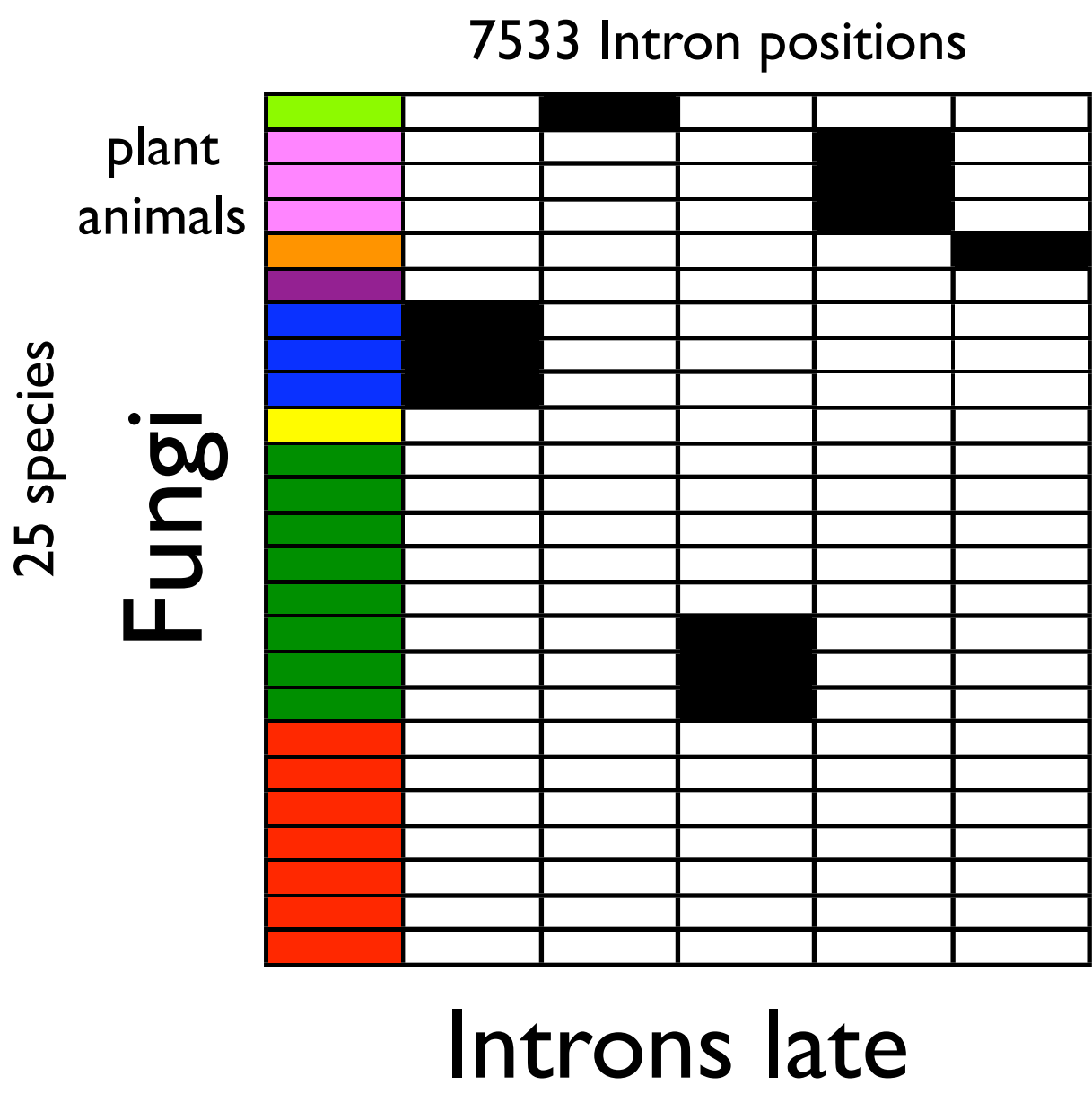
Conserved intron positions



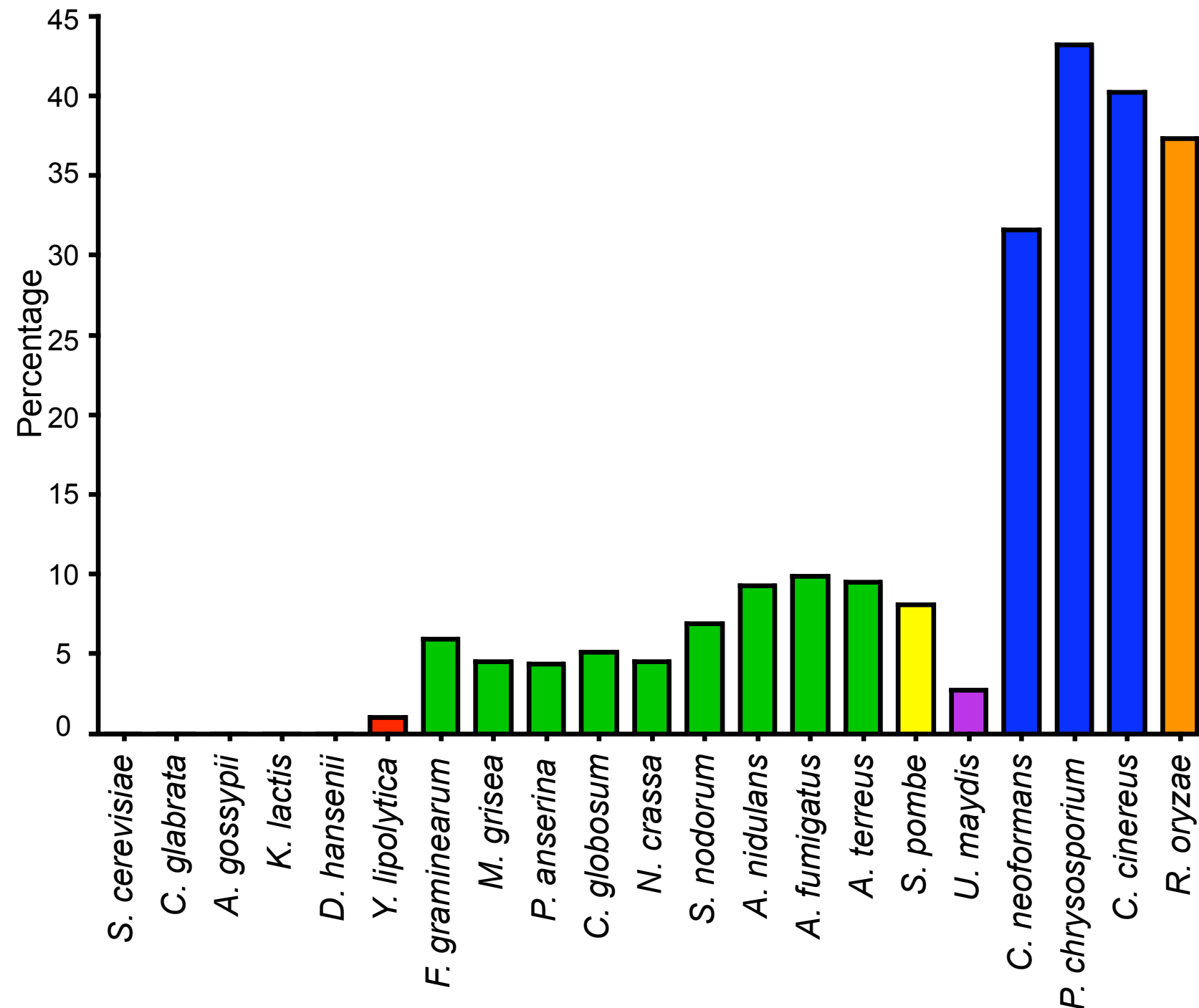
A.thaliana	TWLPEQVA-FIQS1MG-NDKANS-YWEA-----ELPP-----NYD-----RV-GIENFIRAK2Y-----EEKRWV--
Fugu	QWTQEQVQ0SVQE-MG-NAKAKR-LYEA-----FLPK-----CFQRPETDQ-SAEIFIRDK-Y-----DKKKYMDK
Human	QWTQEQIQ0CMQE-MG-NGKANR-LYEA-----YLPE-----TFRRPQIDP2AVEGFIRDK-Y-----EKKKYMDR
Mouse	QWTQEQIQ0CMQE-MG-NGKANR-LYEA-----YLPE-----TFRRPQIDP2AVEGFIRDK-Y-----EKKKYMDR
R.oryzae	TWVQDQVE-NMVR-WG-NERANK2YWEA-----NL-----GDRKPS-ES-NMEMWIRAK-Y-----EQKRWA--
U.maydis	TWTREQVD-RMKE-VG-NLKSNR-KYNPDEMNRNPPT-----NMEESERDS-ELEKYIRRK-Y-----EFRRFV--
C.neoformans	IWTPEQME0SIQK-WG-NKRANM-YWER-----HLKA-----GHI-PS-DH2KIESFIRSK-Y-----ETRRAWA--
C.cinereus	VWTPEQME0SIQK-WG-NRRANL-YWEA-----HLKP-----GHN-PP-EH2KMESFVRSK-Y-----ESRRWA--
P.chrysosporium	VWTVEQMN0SIQK-WG-NKRANI-YWEA-----HLKA-----GHI-PP-DH2KMESFIRSK-Y-----ESKRWA--
S.pombe	SWTDEQTE-NMTR-WG-NERANL2YWEA-----KLAG-----GHV-PS-DS2KIATFIKTK-Y-----EFKKWV--
F.graminearum	SWTDEQLQ-SVLK-WG-NARANK-YWEA-----KLAA-----GHA-PS-EA-KIENFIRTK-Y-----ELKRWV--
N.crassa	AWTDEQLQ-SVLN-WG-NARANK-YWEA-----KLAQ-----GHV-PS-ES-KIENFIRTK-Y-----ELKRWV--
C.globosum	SWTNEQVD0NMRK-VG-NVVSNK-LYNPDN---KNPPVPIDADEA---DG-AMERFIRQK-YIARTLSIGKRRPGGD
P.anserina	SWTDEQLQ-SVLN-WG-NARANK-YWEA-----KLAP-----GHV-PS-EA-KIENFIRTK-Y-----ELKRWV--
M.grisea	SWTDEQLQ-SVLS-WG-NARANK-YWES-----KLAA-----GHA-PS-EA-KIENFIRTK-Y-----ELKRWV--
S.nodorum	TWTDEQLQ-SVLK-WG-NARANK-YWEA-----KLAP-----GHV-PS-EA-KIENFIRTK-Y-----ESKRWT--
A.nidulans	AWTDEQLQ-SVVR-WG-NARANK2YWEA-----KLAP-----GHV-PP-EA2KIENFIRTK-Y-----ESKRWV--
A.fumigatus	SWTDEQLQ-SVIK-WG-NARANK2YWEA-----KLAP-----GHV-PS-EA2KIENFIRTK-Y-----ESKRWV--
A.terreus	SWTDEQLQ-SVLR-WG-NARANK2YWEA-----KLAP-----GHV-PS-EA2KIENFIRTK-Y-----ESKRWV--
Y.lipolytica	AWTEEQLA-SMMK-WG-NTRCNM-FWEA-----KLPK-----GHV-PD-DN-KIENFIRTK-Y-----DMKKWA--
D.hansenii	AWTDEQVE-LMIK-WG-NEKCNI-YWES-----KLPD-----GYV-PD-QL-KIDNFIRTK-Y-----DLKKWV--
K.lactis	TWTEEHLE-AVLE-FGNNKKFNE-YYEN-----KLGG-----GTYPVD-QS-KIGQFIRTK-Y-----ELKKWV--
A.gossypii	TWKEEHLQ-QVVR-FGNNQQANK-VFEG-----RLGG-----GSYVPD-QS-KMGQFIKTK-Y-----EVRKWY--
C.glabrata	TWKEEHLK-ELVQ-MRNNVNANR-VYEA-----KLPDSSKFNGKSLGNDIN-LLQEFIRQK-Y-----ERKRWM--
S.cerevisiae	TWKEEHLV-KLIQ-FKNNLRANS-YYEATL-ADELKQ-----RKI-TD-TS-SLQNFINKN-Y-----EYKKWI--

* :: : * : :: :: * *

Patterns of conservation



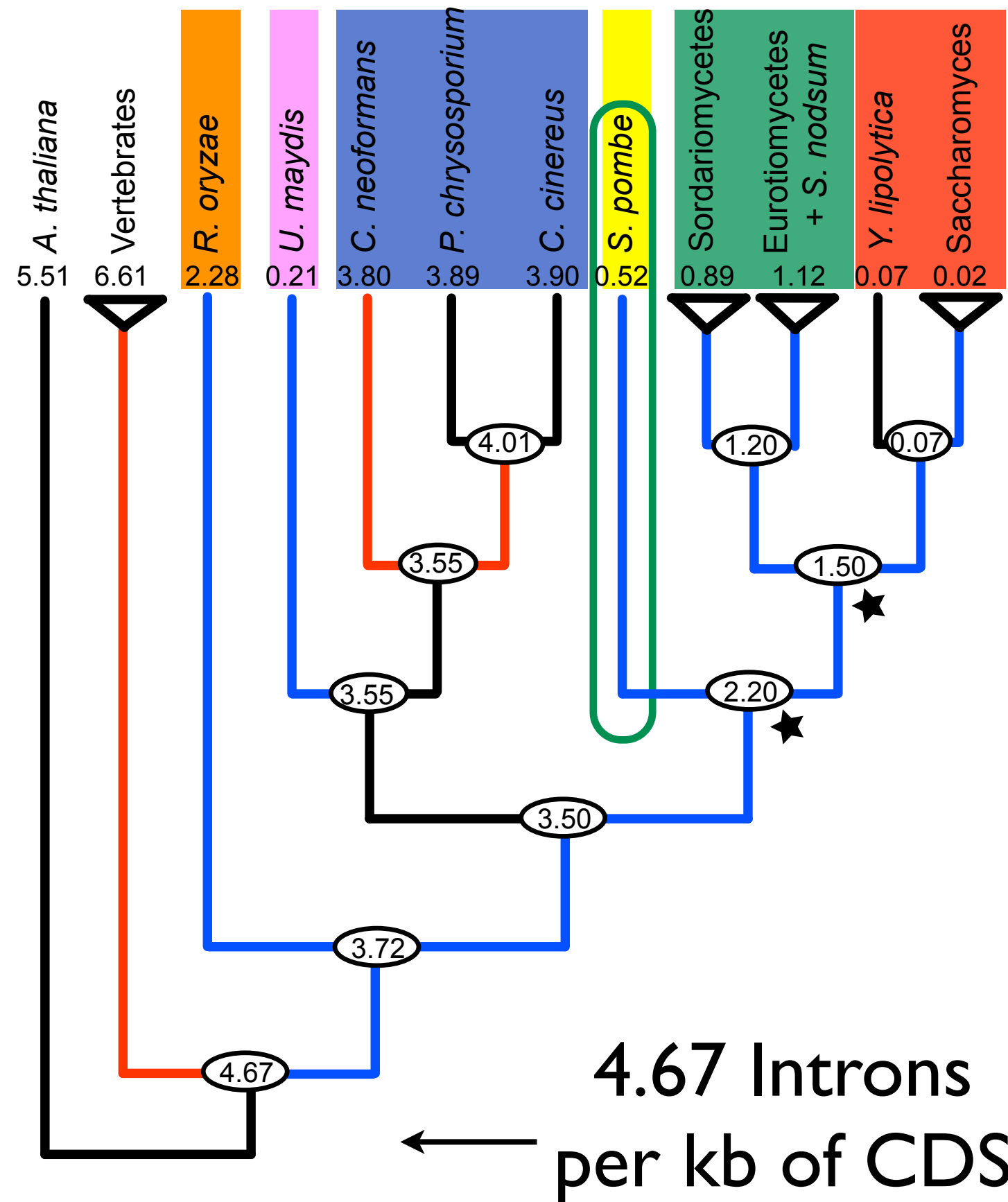
Intron positions shared with animals or plants



Intron position reconstruction

- 3 Methods
 - Roy and Gilbert. 2005
 - Csűrös. 2005
 - Nguyen et al. 2006
- Methods agree for all but 2 nodes in tree

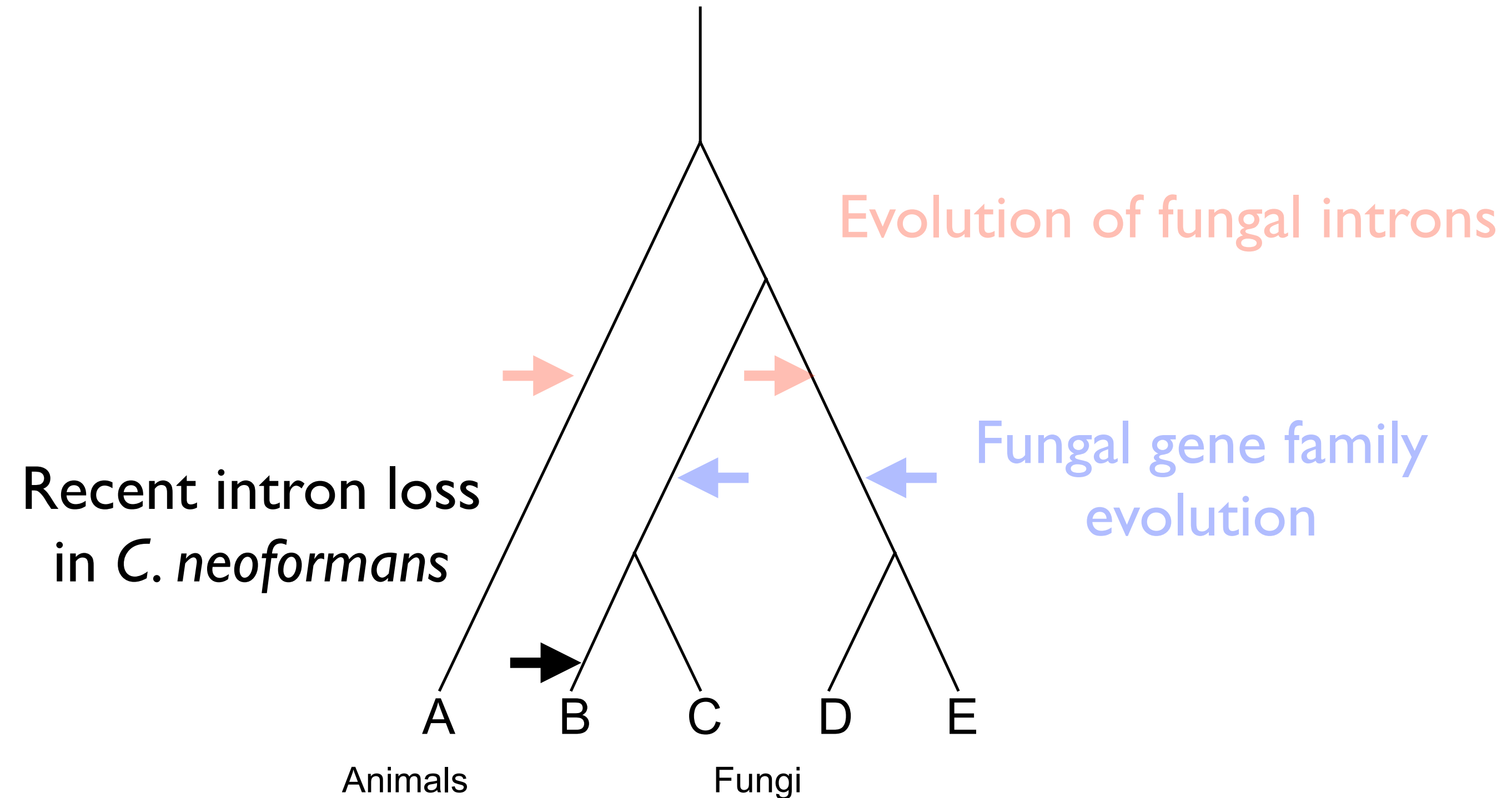
Reconstruction of ancestral intron densities



Conclusions

- Early eukaryotic crown genes were complex!
- Ancestor had 70% of the introns in vertebrates - many more than previously reported
- Intron loss has dominated among the fungi
 - **Hemiascomycota** experienced loss
- No significant evidence for intron sliding or double insertions
- Sampling can bias interpretations - all fungi are not equal.

Fungal comparative genomics



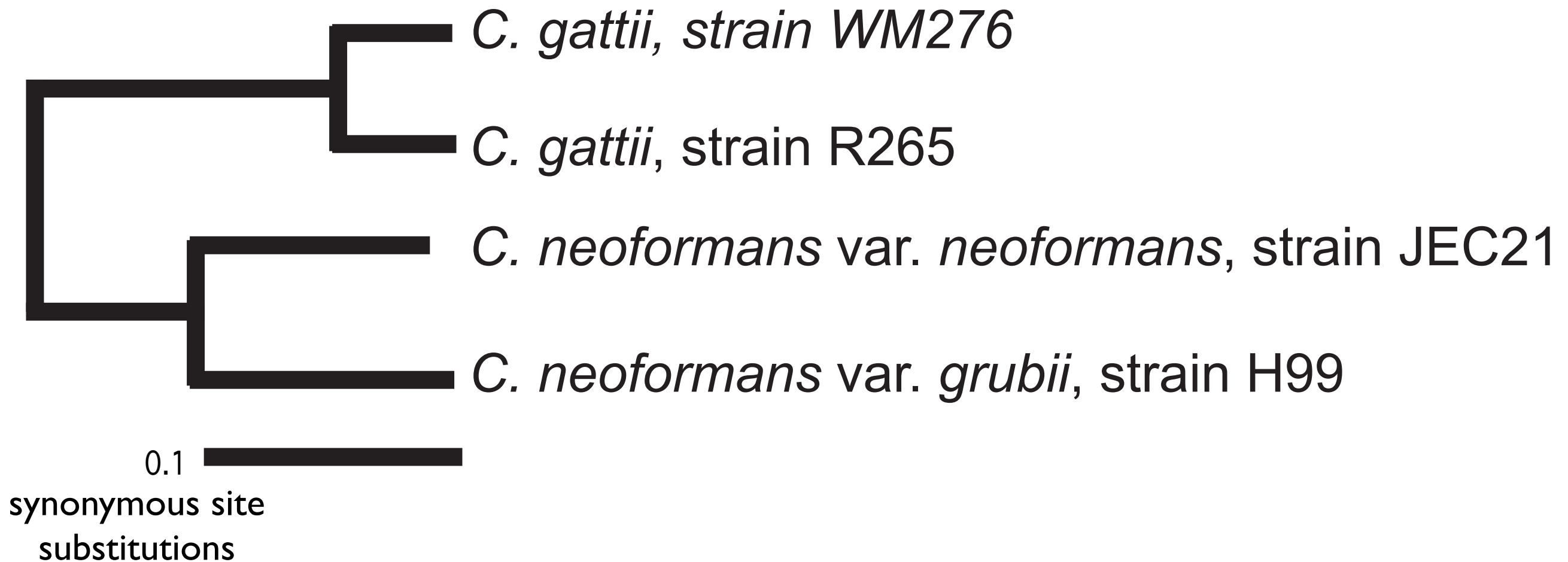
Mechanism of intron loss in fungi

- *S. cerevisiae* and **Hemiascomycota** have undergone intron loss.
- How are introns lost from the genome?
 - Are they lost independently?
 - Are they lost many at a time?
- What is the molecular mechanism of loss?

Models of intron loss

- All introns in *S. cerevisiae* are in 5' end of gene.
- G. Fink proposed transcripts recombine with genome 3' -> 5' explaining 5' retention bias.
- Most lost events in *S. cerevisiae* occurred too long ago to find evidence of mechanism with a comparative approach.

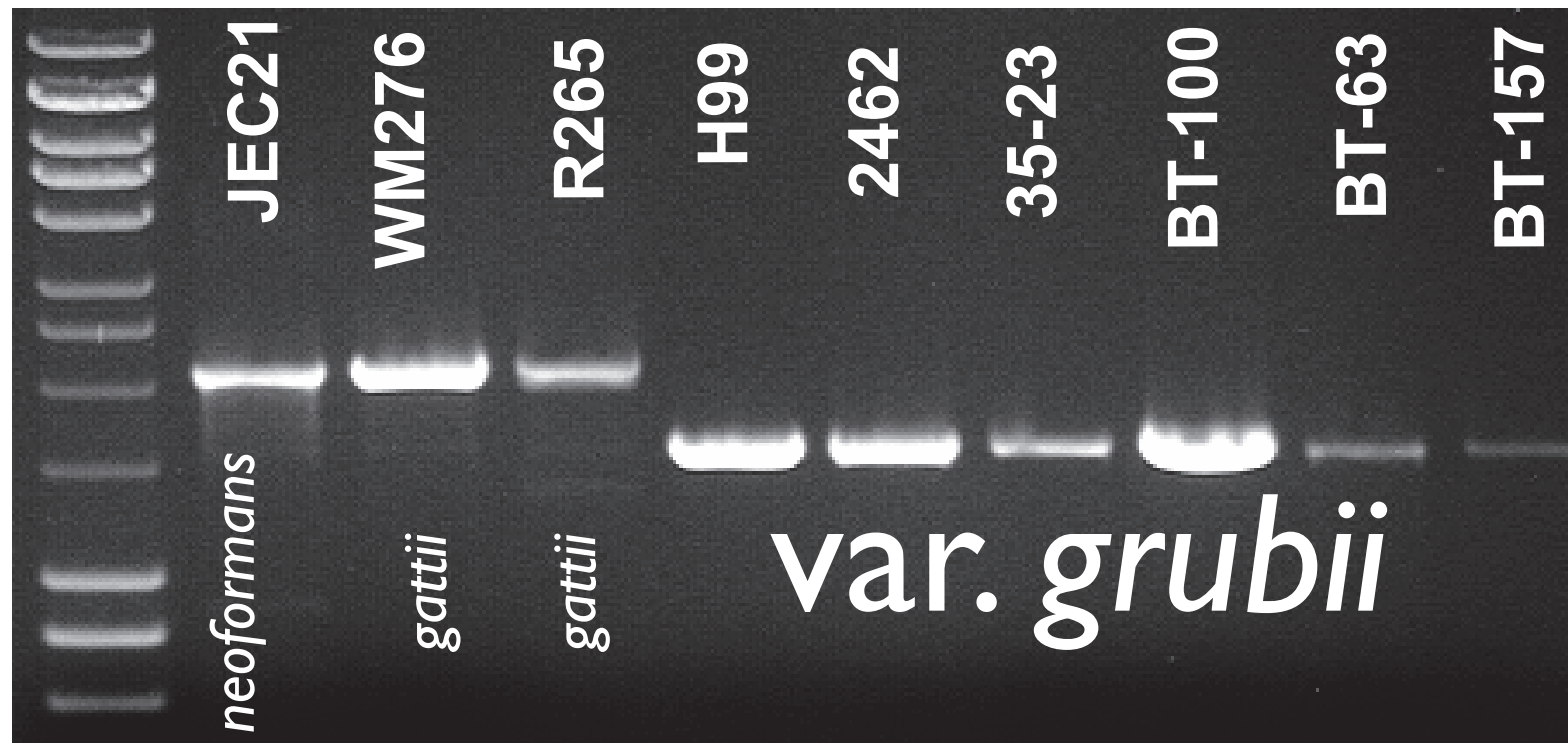
Sequenced *Cryptococcus* genomes



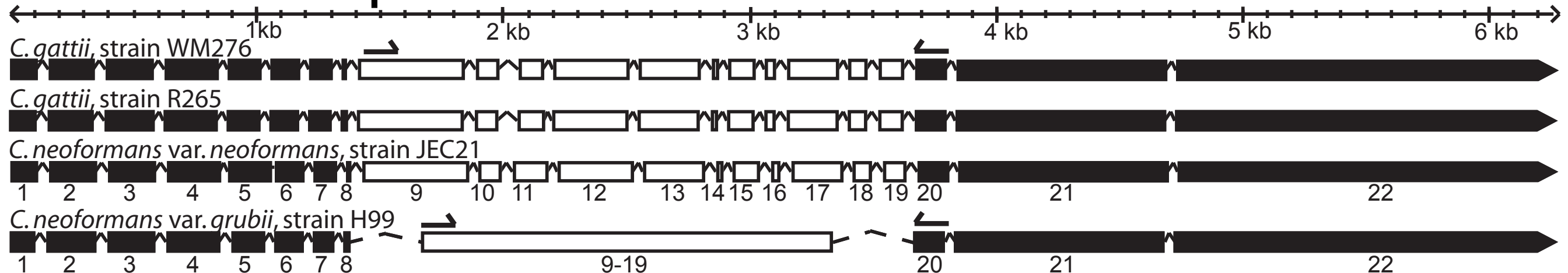
Screen for intron changes

- Annotate 3 *Cryptococcus* genomes (var. *grubii* and 2 var. *gattii* genomes)
- Identify and align 4-way orthologous genes
 - 5298 orthologous genes (out of ~6500)
- Identify intron position changes


Intron loss in var. *grubii*



CNI01550 - putative RNA helicase

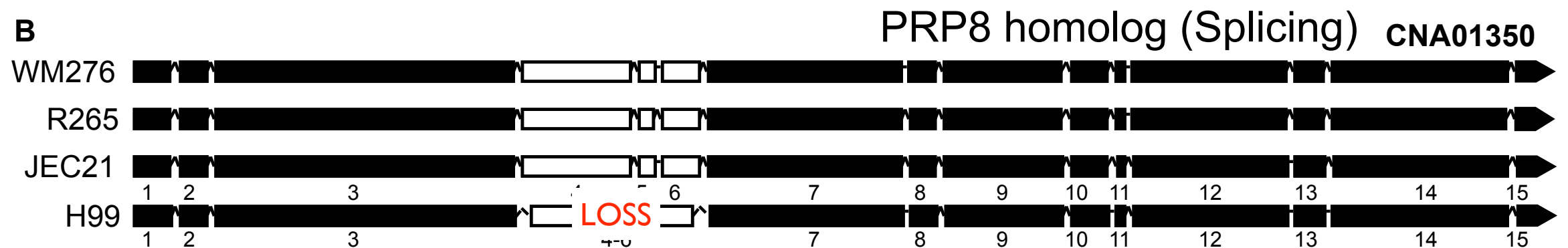
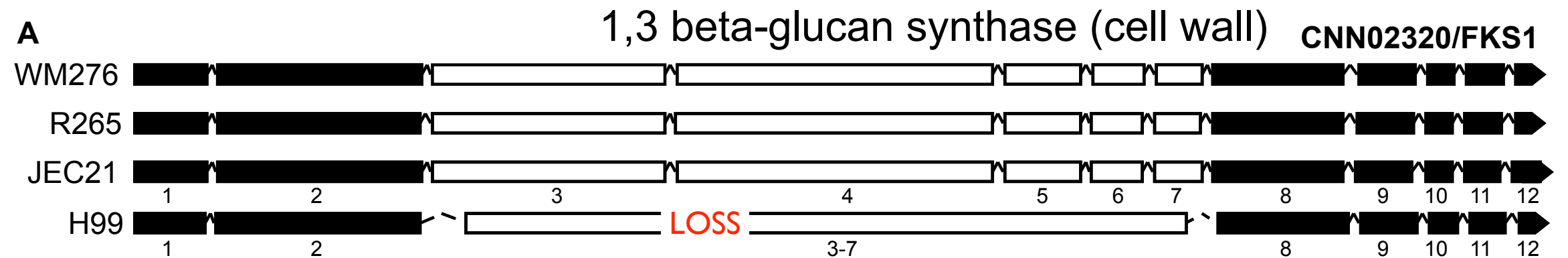


Intron loss was a precise excision

R265	CGACAAGTACATAAACTTTTTTTGTGCCTGGCGCAAAGACTTTCCATTGCTGACAGAAAACAGGTTGAA
WM276	AGACAAGTACATAAACTTTTTTTGTCTCTCCTCCAAACATTTTTCATTGCTGACAGAAAACAGGTTGAA
H99	AGACAA ←  → -GTTGAA
JEC21_CDS	AGACAA-----GTTGAA
JEC21	AGACAAG GT ACATACTAGTCCTTGTG---CTATCCCAAAGACTTT-CATTGCTGACAGAAAAC AG GTTGAA

R265	CGCTGCCGAATTATGTCGATGTTGGAGATTTCTTGAGGTAAGCAACAGACTCGTAACAGCTTGTTTCGGTC
WM276	CGCTGCCGAACACTATGTCGATGTTGGAGATTTCTTGAGGTAAGCAACAGACTCGTAACAGCTTGTTTCGGTC
H99	CCCTGCCGAATTATGTCGACGTTGGAGATTTCTTGAG-----
JEC21_CDS	CCCTGCCGAATTATGTCGATGTTGGAGATTTCTTGAG-----
JEC21	CCCTGCCGAATTATGTCGATGTTGGAGATTTCTTGAG GT ACGTCGCAAACACTCGTAACAGCTTGTTTCGATC
	* *****

Other examples of loss



Conclusions

- Intron loss via homologous recombination with spliced transcript.
- Multiple adjacent introns are lost.
- Precise deletions of introns.
- Loss biased towards the middle of gene not 3'.

Overall conclusions

- Multiple genome sequences have helped resolve several outstanding questions in evolution introns
- Origin of introns still mystery, but early eukaryotic genes were complex
- Suggested intron function in splicing
- Gene family expansions can be important in identifying molecular basis for adaptation