# A Comparative Genomic Investigation of Fungal Genome Evolution

Jason Stajich
Duke University
University Program in Genetics & Genomics

# Evolutionary genomics

## Evolution & Organismal

Phenotype
Population structure
Ecological adaptation
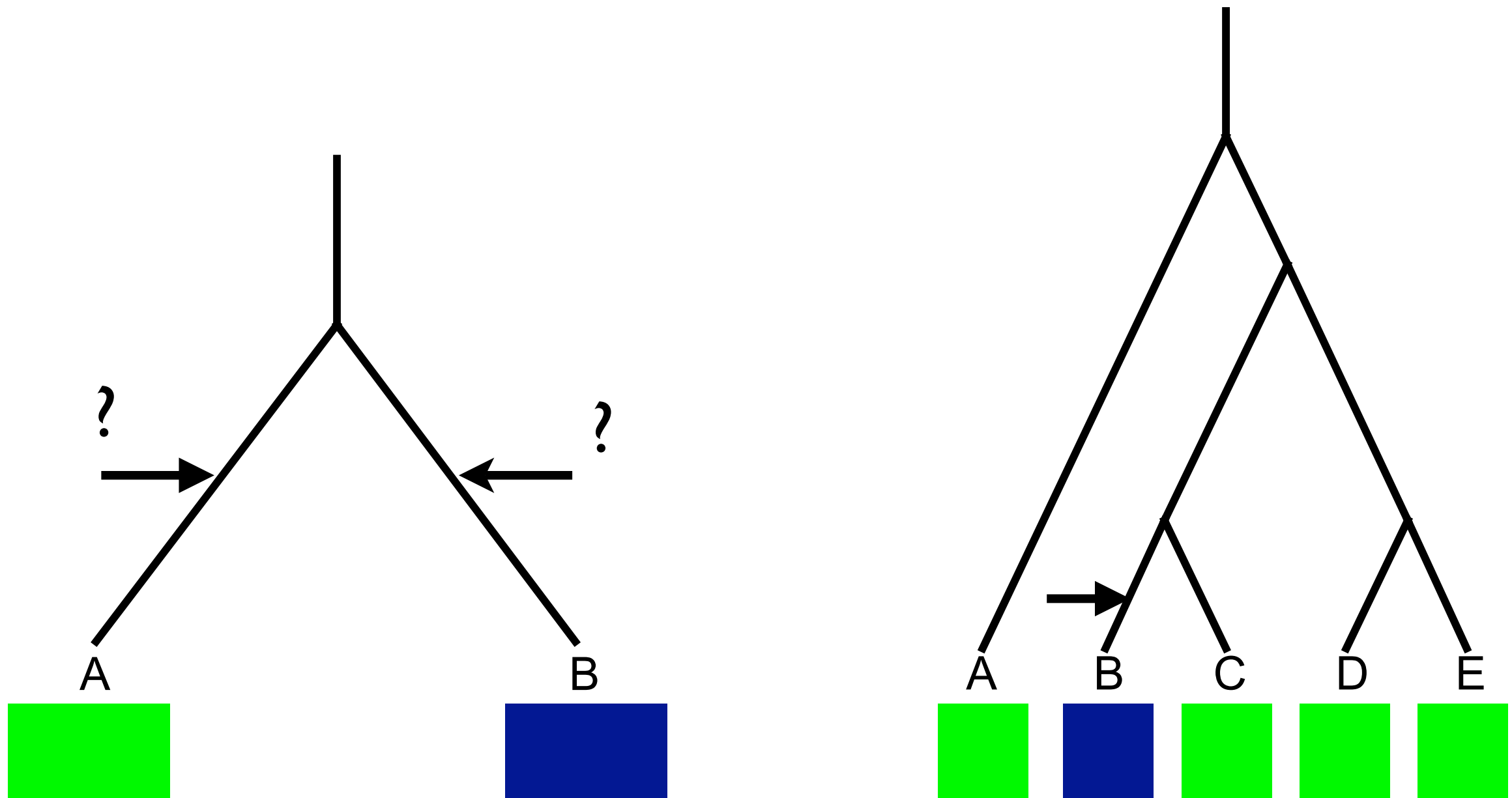Niche changes
Phylogeny

## Comparative Genomics

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

## 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*

# Agricultural impact of fungi



...s of plant disease is caused by

USDA

A.G. Bölker

...posit mycotoxins - e.g. ergot

...al fungi provide nutrient

...and nitrogen fixation

Western Committee on Plant Diseases

www.gov.mb.ca

# Impact of fungi on human health

- Biggest risk for immunocompromised individuals

- Primary pathogens

  - *Histoplasma, Coccidioides, Cryptococcus gattii*

- Opportunistic pathogens

  - *Candida albicans, Aspergillus fumigatus, Cryptococcus neoformans, Rhizopus oryzae*
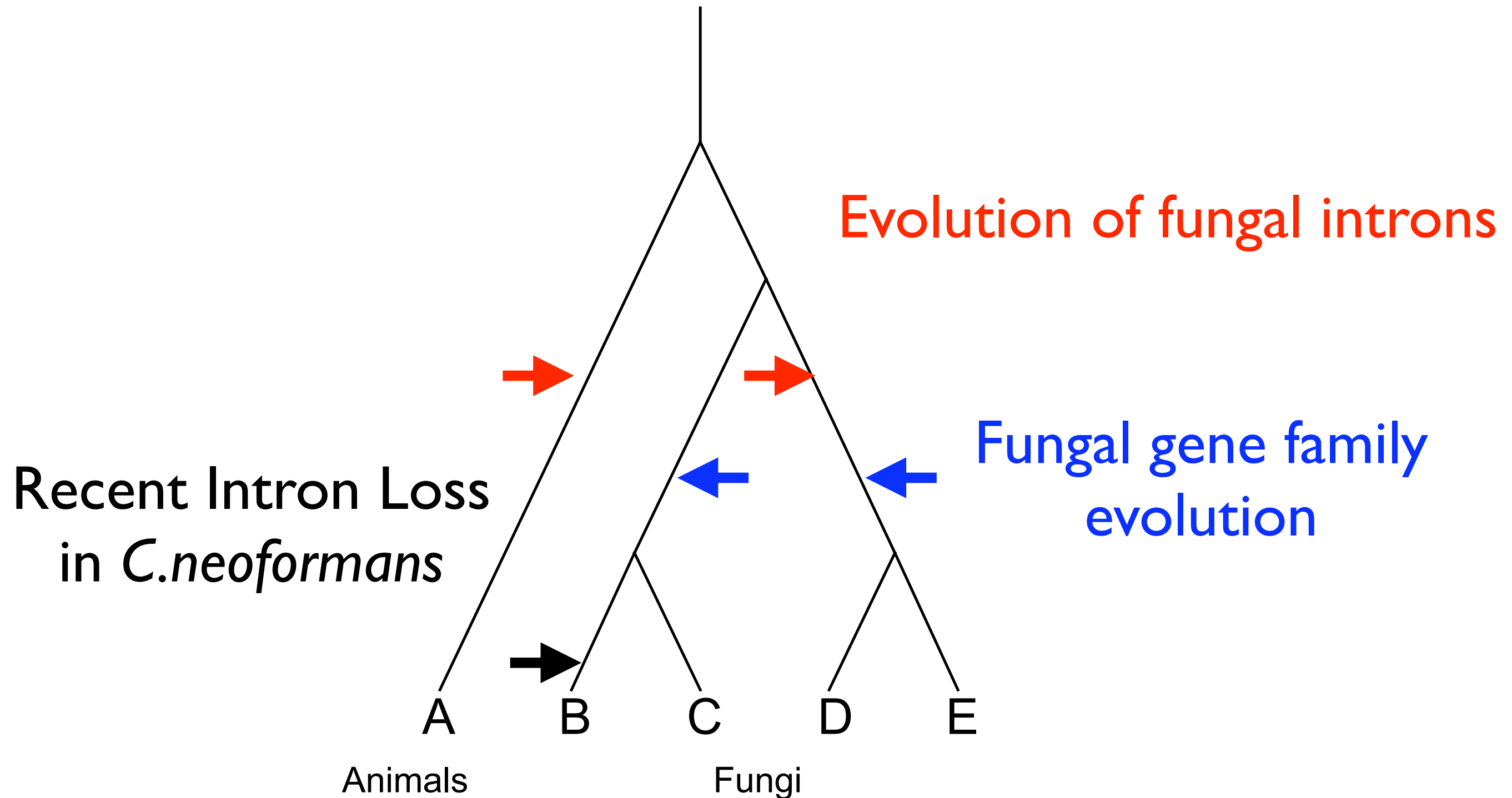
# Fungi as genetic models

- Beadle and Tatum (1941) - one gene, one enzyme hypothesis in *Neurospora crassa*

- Cell cycle, cell model - *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe*

- Straightforward molecular biology tools to investigate phenotype-genotype
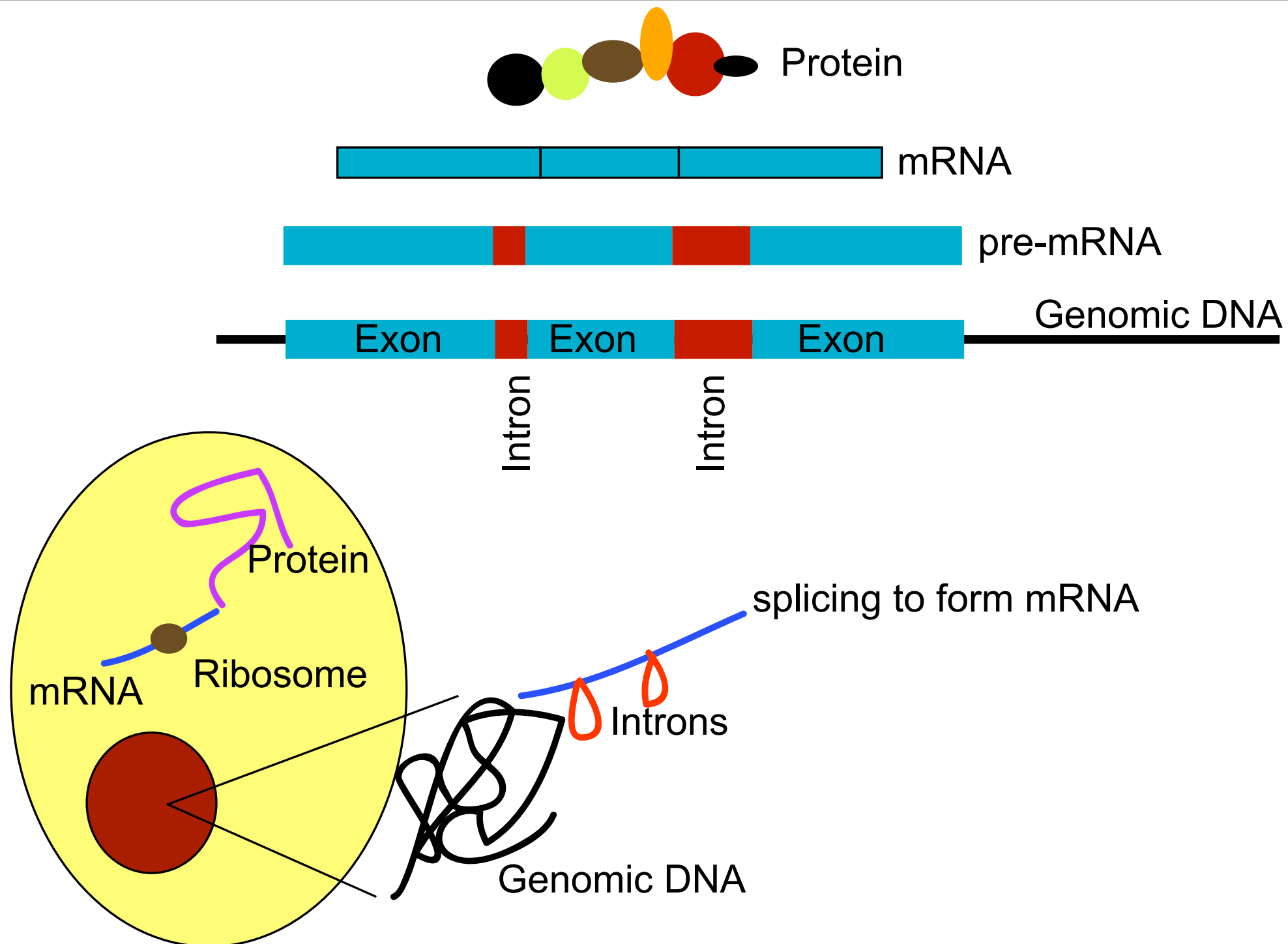
- Evolutionary models

# 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
  - 4700 in *A. gossypii*, 5800 in *S. cerevisiae*
  - 16,000 in *R. oryzae* or *S. nodorum*
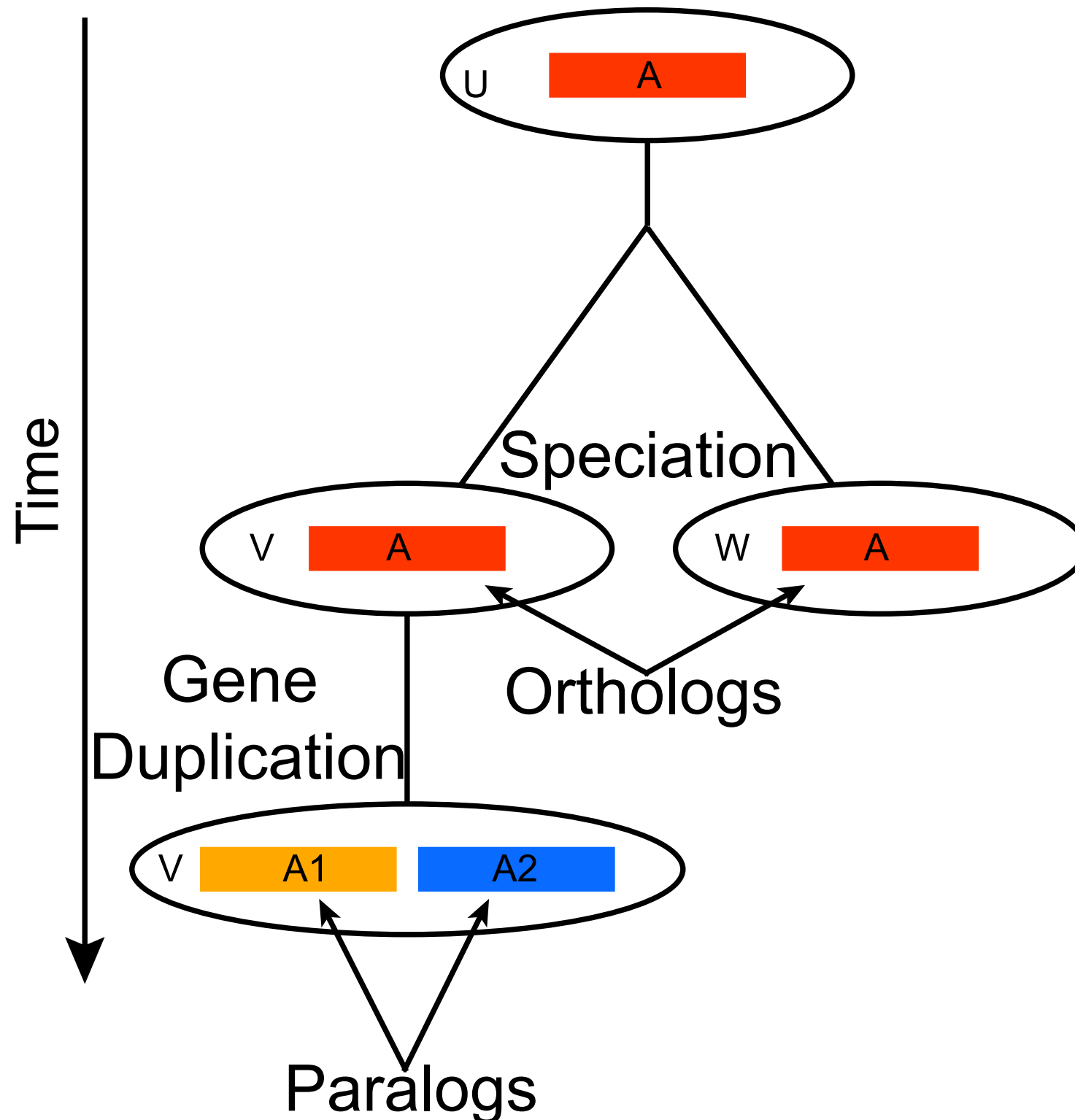  - 19,000 in Fruitfly, 25,000 in worm

# Fungal comparative genomics



Evolution of fungal introns

Fungal gene family evolution

Recent Intron Loss in *C.neoformans*

A B C D E

Animals Fungi
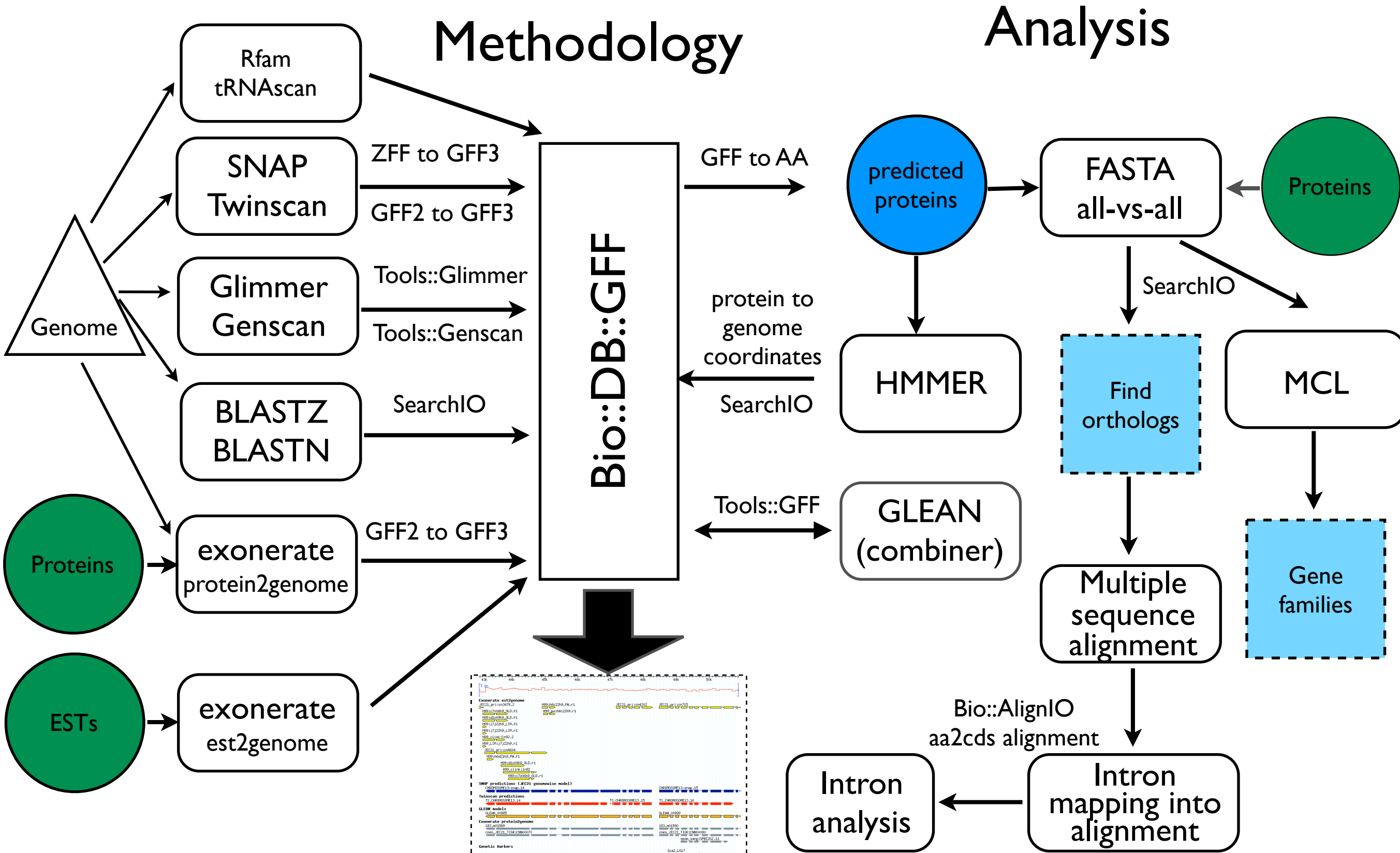
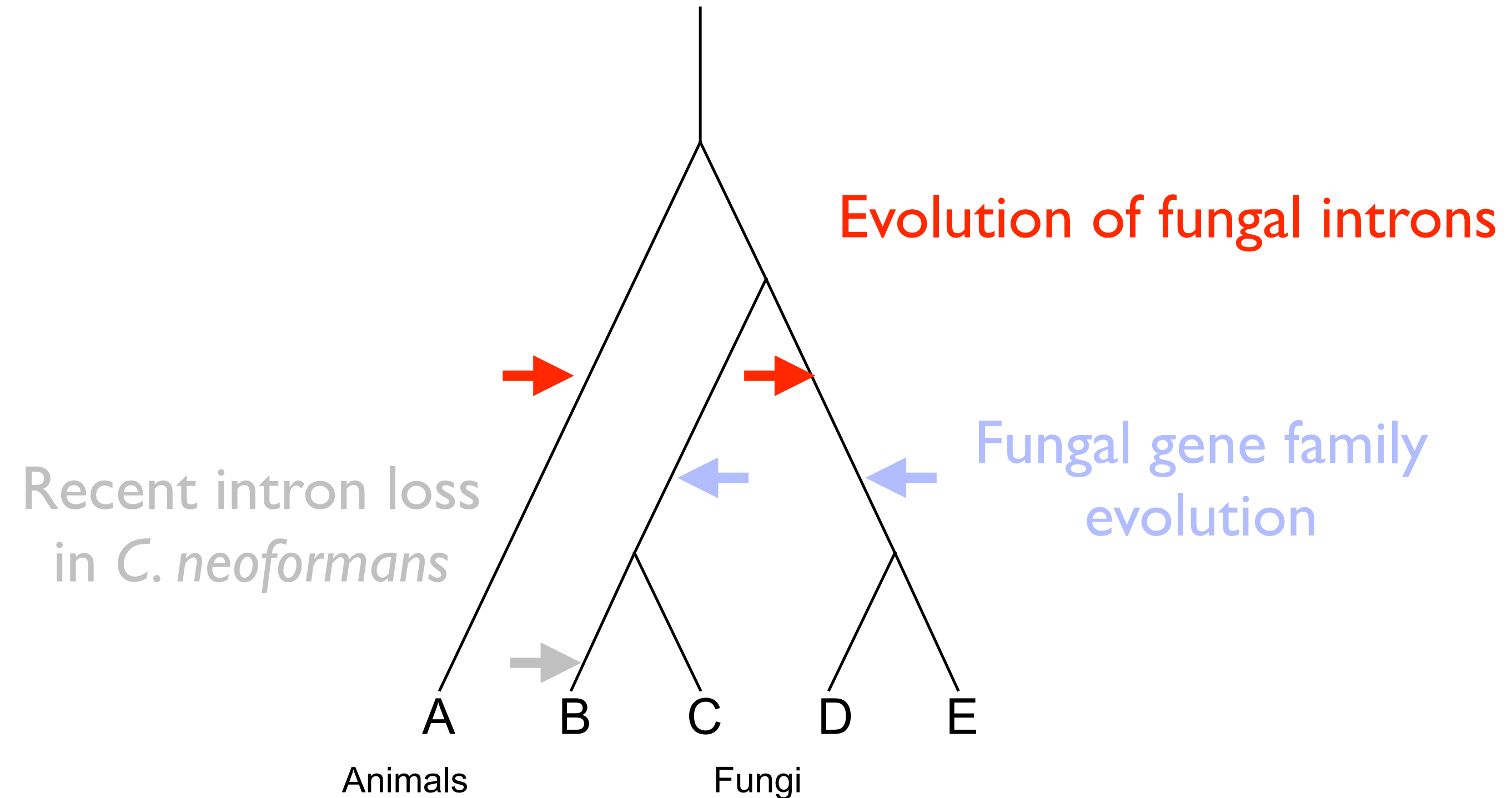# Central dogma of eukaryotic biology

# Orthologs and Paralogs

# Genome annotation

- Many of the fungal genomes were only assembled genomic sequence.

- Automated annotation pipeline was built to generate to get systematic gene prediction.

- Several gene prediction programs were trained and results were combined to produce composite gene calls

Methodology

Analysis

http://fungal.genome.duke.edu

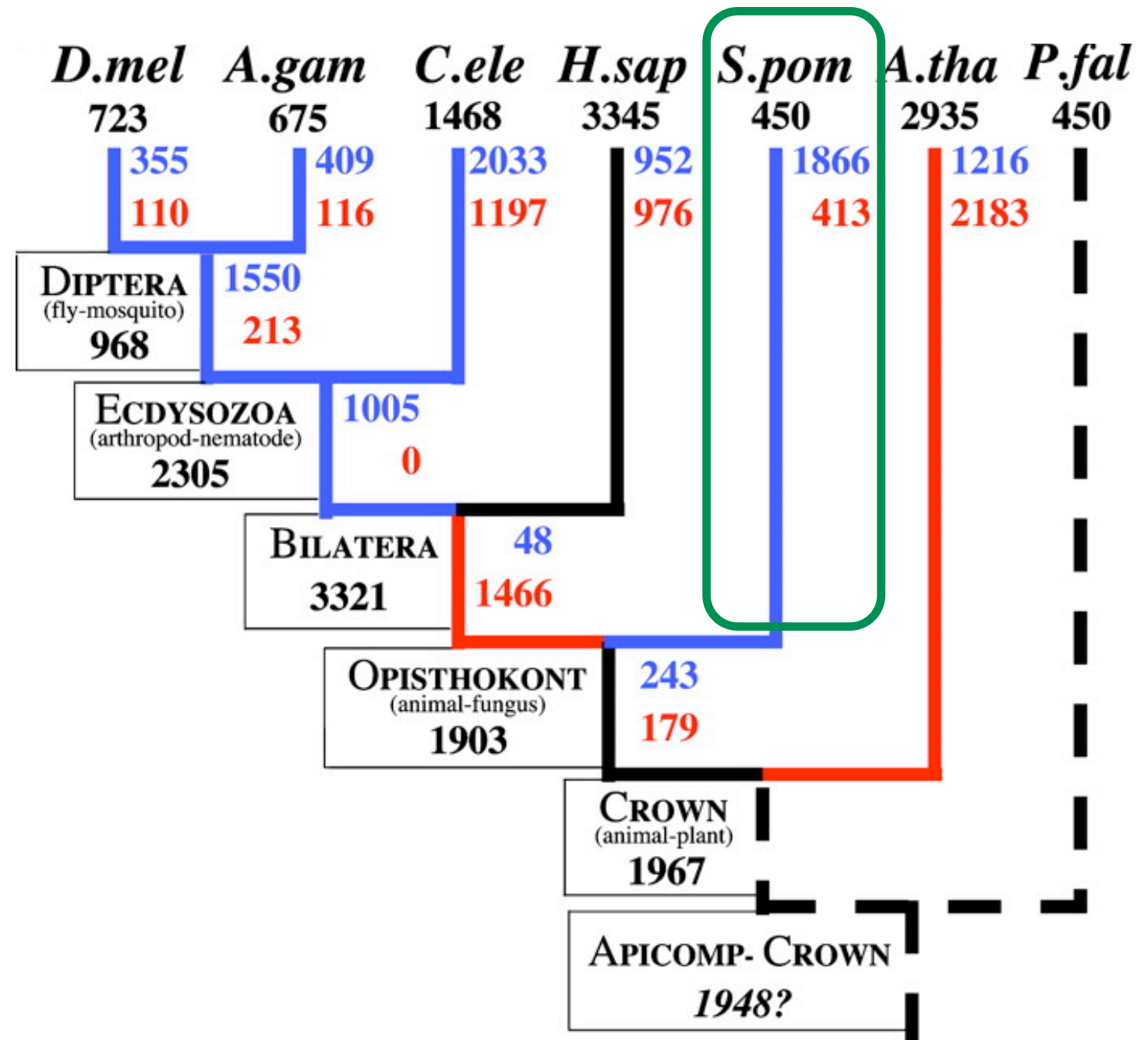# Fungal comparative genomics
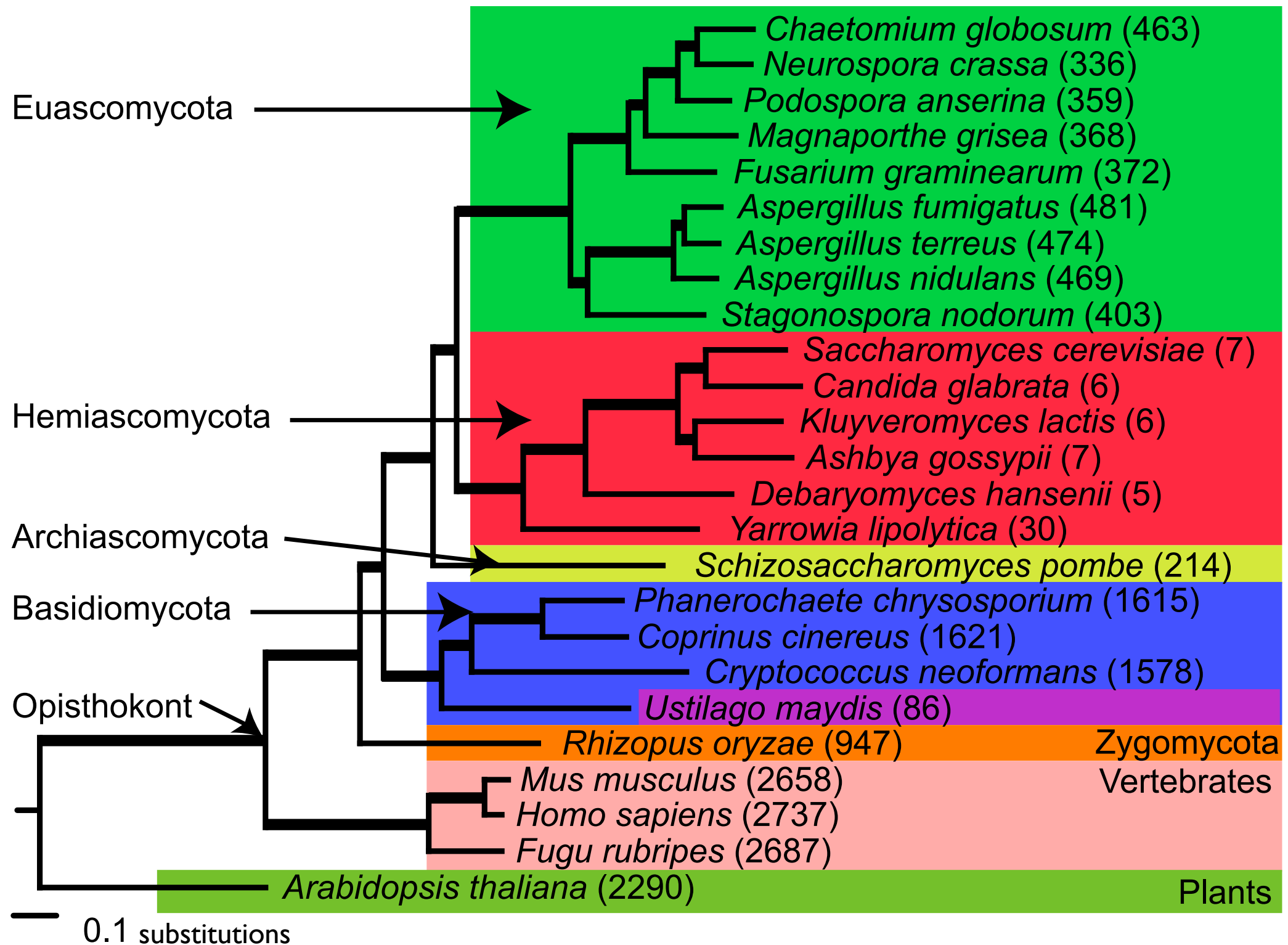
# Evolution of gene structure

- Present day introns

  - Recent insertions?

    - Introns late hypothesis

  - Formed in eukaryotic ancestor?

    - Introns early hypothesis / exon theory of genes

  - Mixture of two?
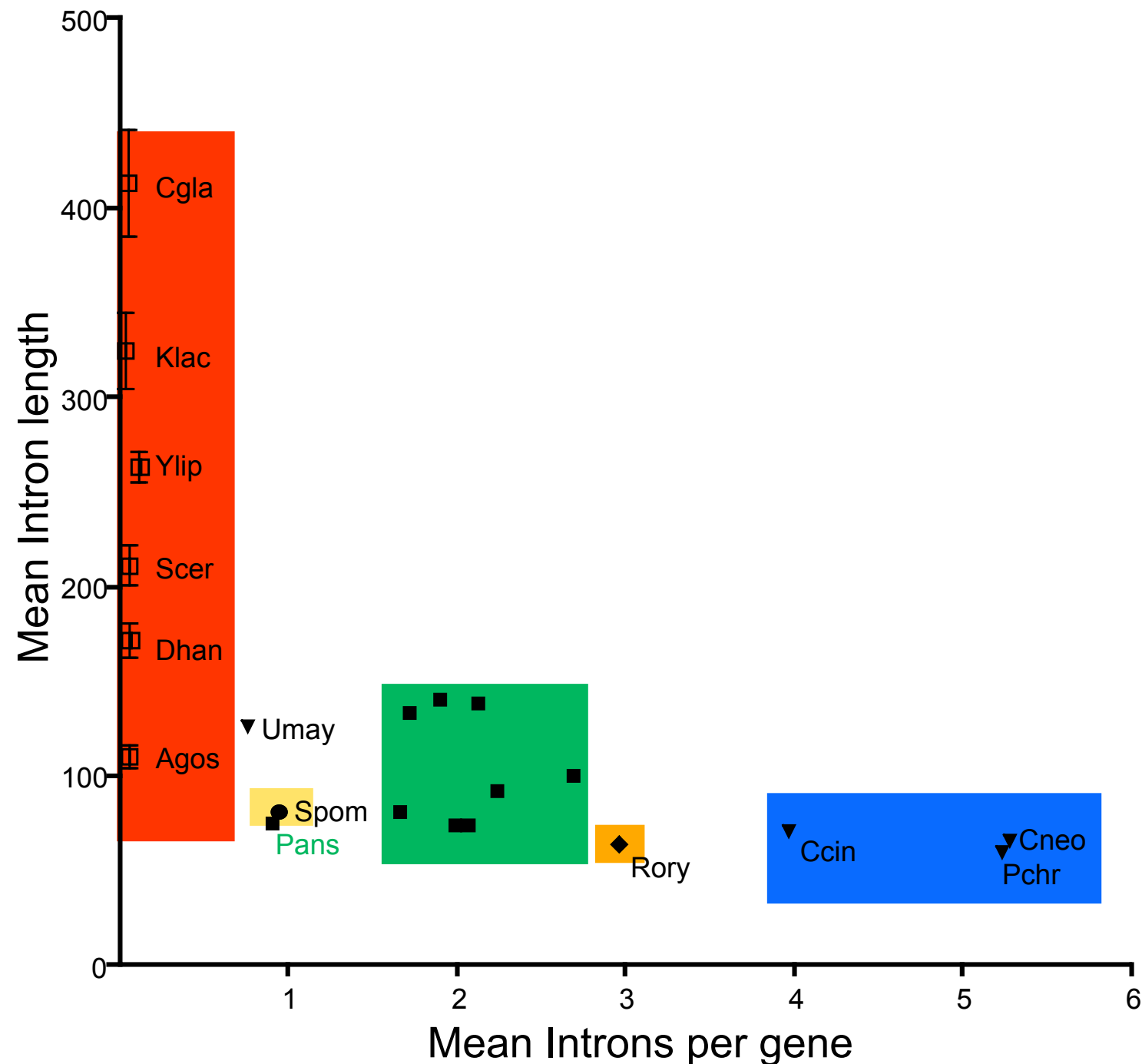
# Previous work on intron evolution

- Rogozin et al. 2003

  - 7 genomes

  - 684 genes, 7236 positions

- Other methods

  - Roy and Gilbert. 2005

  - Csũrös. 2005

  - Nguyen et al. 2006
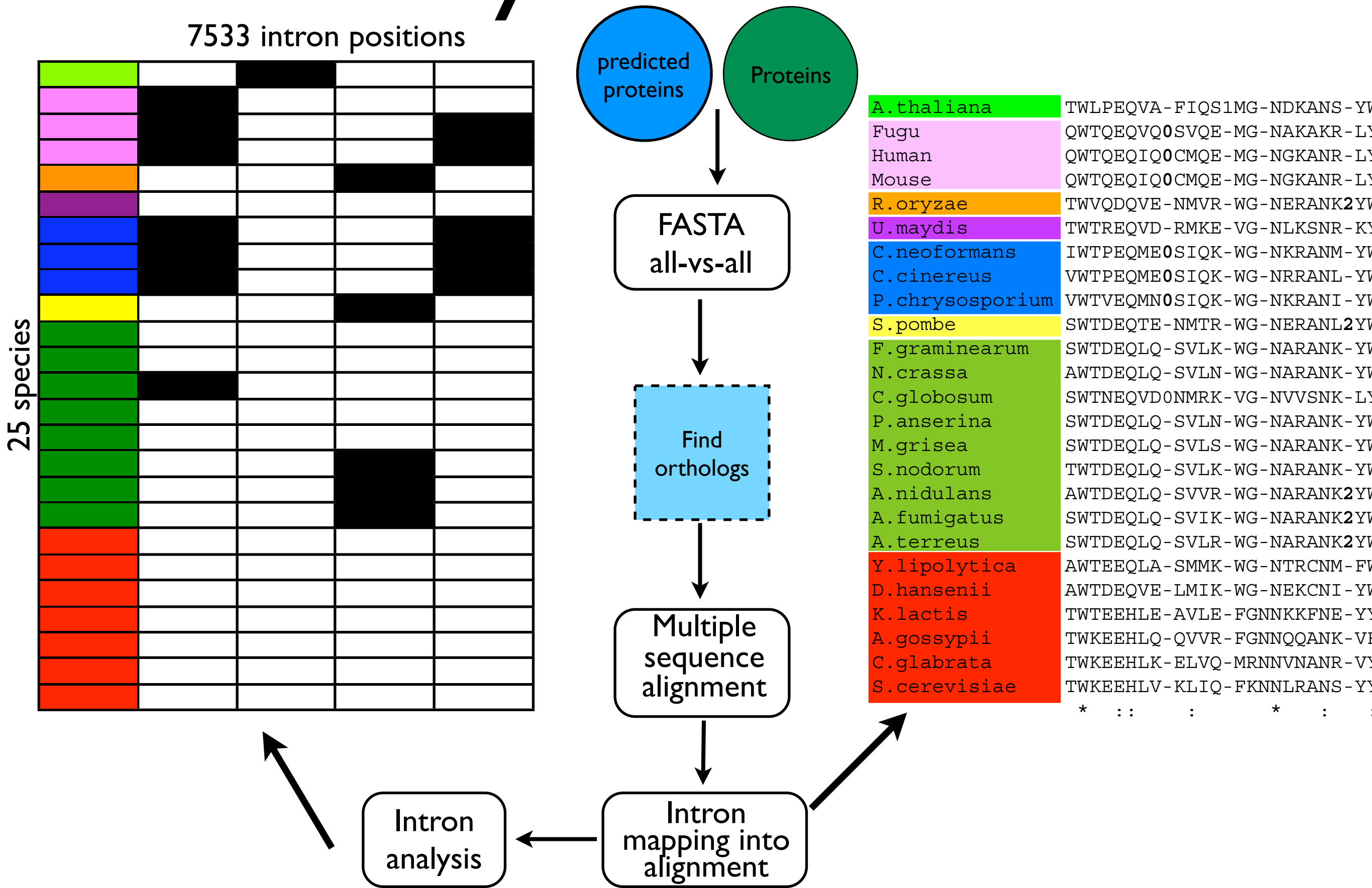
# Calculating intron densities across a phylogeny

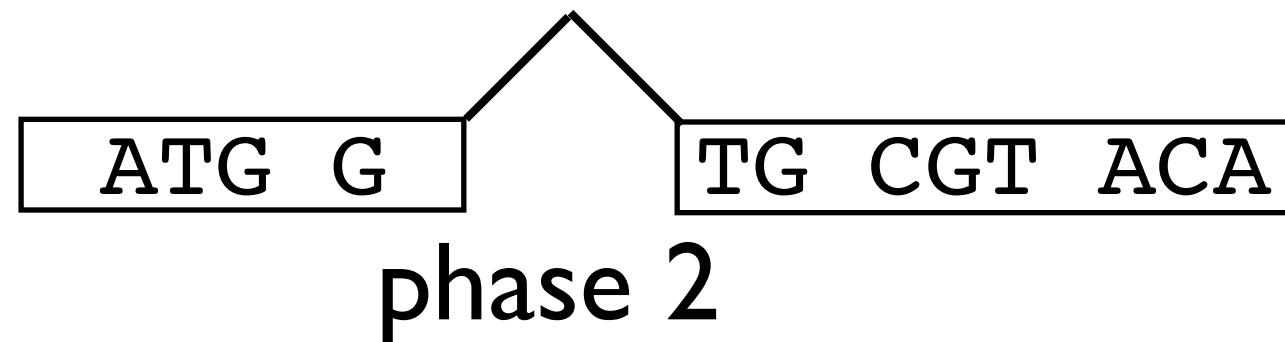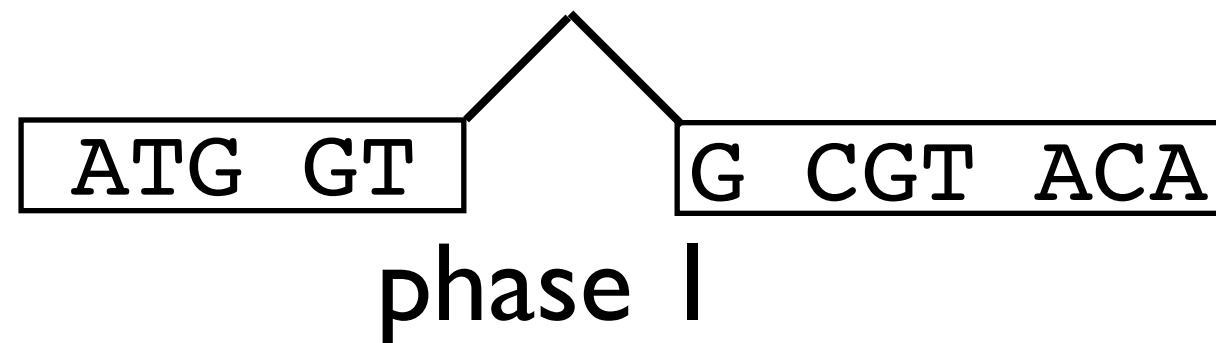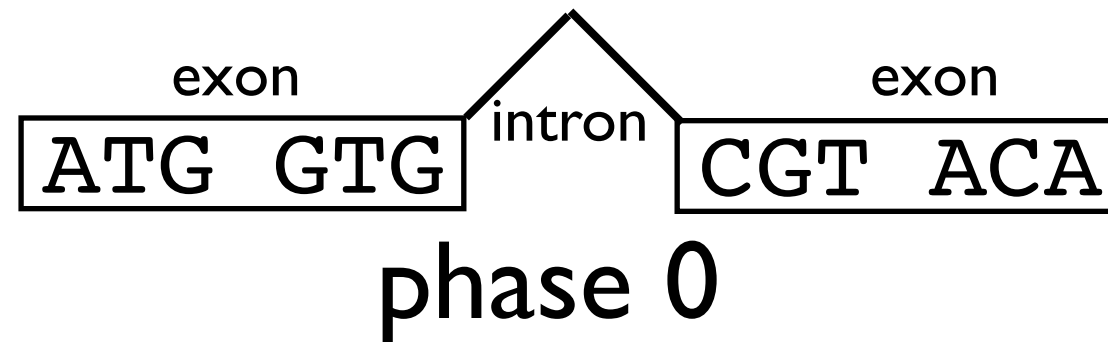# Intron frequency varies among the fungi

# 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 coding sequence (CDS) per genome
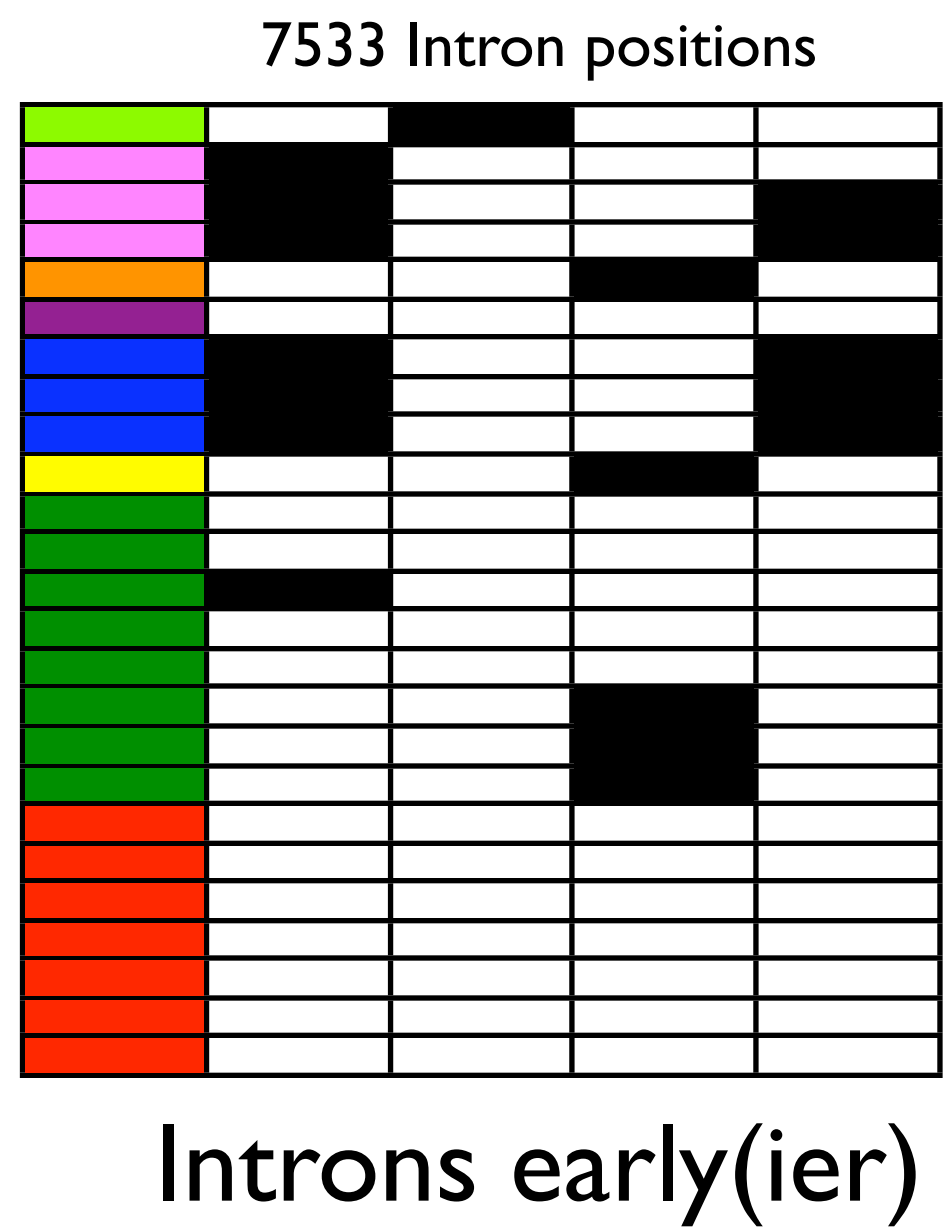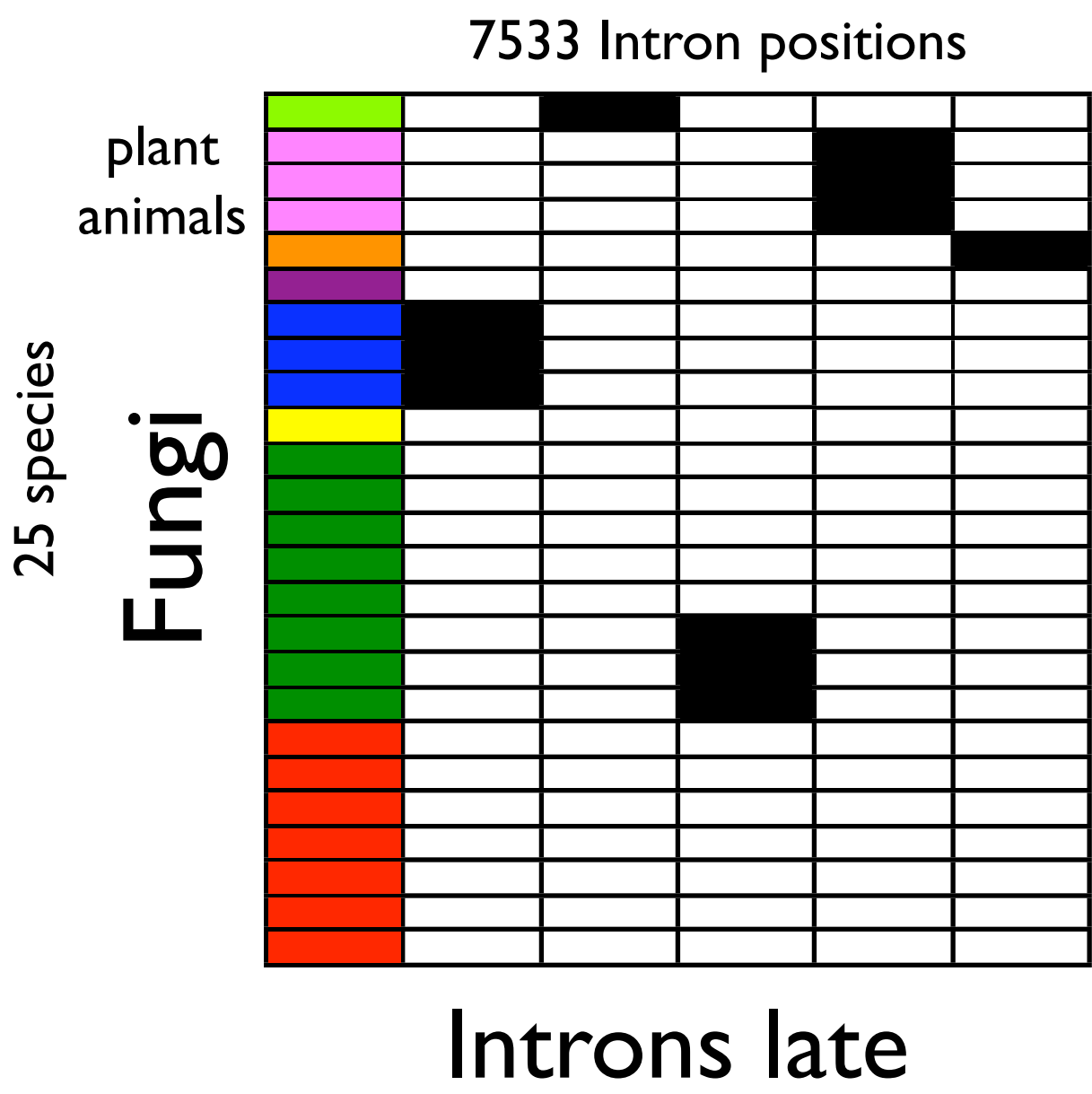
# Analysis Methods

7533 intron positions

25 species
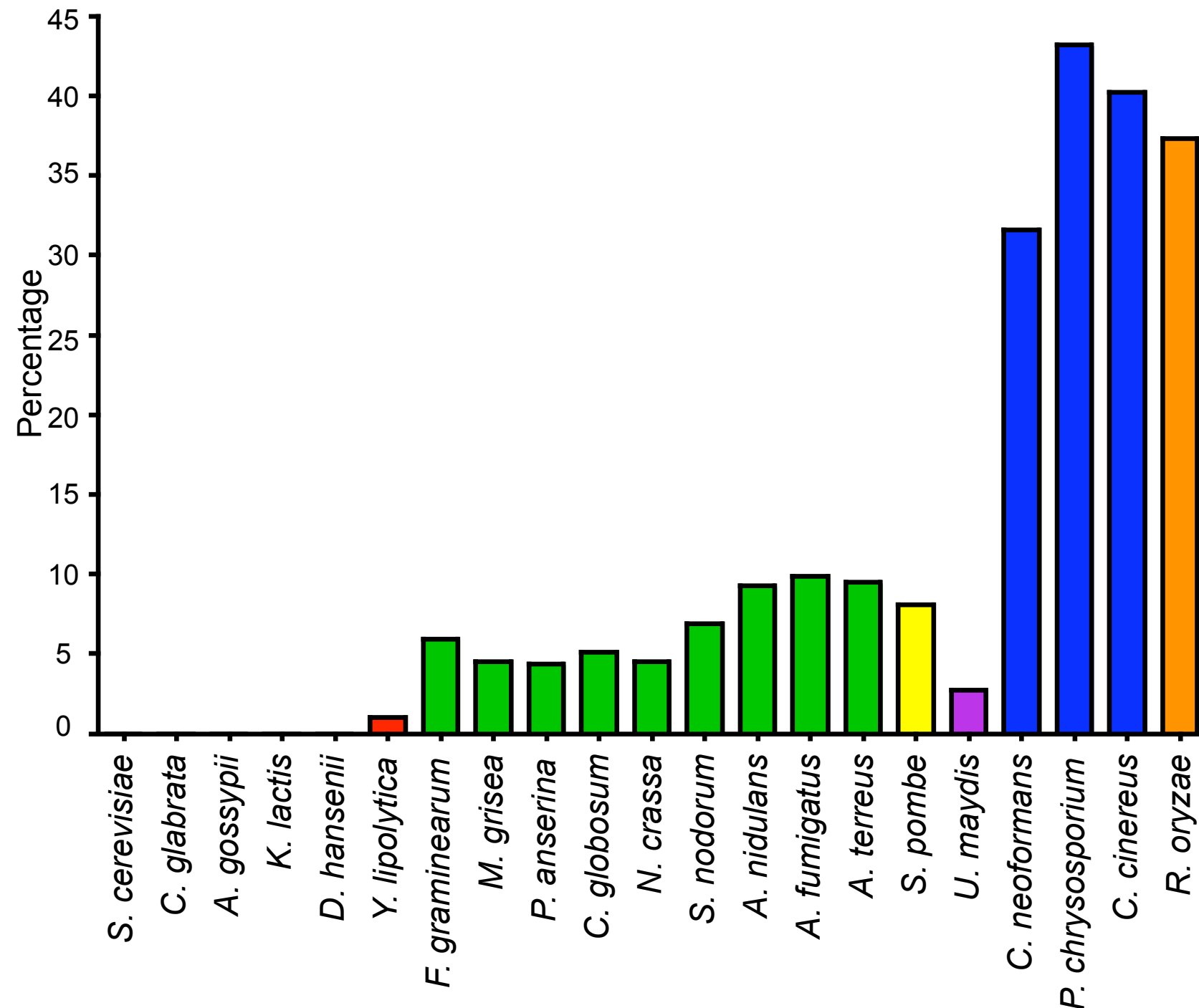
predicted proteins

Proteins

FASTA all-vs-all

Find orthologs

Multiple sequence alignment

Intron mapping into alignment

Intron analysis

| | |
|---|---|
| A.thaliana | TWLPEQVA-FIQS1MG-NDKANS-YW |
| Fugu | QWTQEQVQ0SVQE-MG-NAKAKR-LY |
| Human | QWTQEQIQ0CMQE-MG-NGKANR-LY |
| Mouse | QWTQEQIQ0CMQE-MG-NGKANR-LY |
| R.oryzae | TWVQDQVE-NMVR-WG-NERANK2YW |
| U.maydis | TWTREQVD-RMKE-VG-NLKSNR-KY |
| C.neoformans | IWTPEQME0SIQK-WG-NKRANM-YW |
| C.cinereus | VWTPEQME0SIQK-WG-NRRANL-YW |
| P.chrysosporium | VWTVEQMN0SIQK-WG-NKRANI-YW |
| S.pombe | SWTDEQTE-NMTR-WG-NERANL2YW |
| F.graminearum | SWTDEQLQ-SVLK-WG-NARANK-YW |
| N.crassa | AWTDEQLQ-SVLN-WG-NARANK-YW |
| C.globosum | SWTNEQVD0NMRK-VG-NVVSNK-LY |
| P.anserina | SWTDEQLQ-SVLN-WG-NARANK-YW |
| M.grisea | SWTDEQLQ-SVLS-WG-NARANK-YW |
| S.nodorum | TWTDEQLQ-SVLK-WG-NARANK-YW |
| A.nidulans | AWTDEQLQ-SVVR-WG-NARANK2YW |
| A.fumigatus | SWTDEQLQ-SVIK-WG-NARANK2YW |
| A.terreus | SWTDEQLQ-SVLR-WG-NARANK2YW |
| Y.lipolytica | AWTEEQLA-SMMK-WG-NTRCNM-FW |
| D.hansenii | AWTDEQVE-LMIK-WG-NEKCNI-YW |
| K.lactis | TWTEEHLE-AVLE-FGNNKKFNE-YY |
| A.gossypii | TWKEEHLQ-QVVR-FGNNQQANK-VI |
| C.glabrata | TWKEEHLK-ELVQ-MRNNVNANR-VY |
| S.cerevisiae | TWKEEHLV-KLIQ-FKNNLRANS-YY |

```
 *  ::   :      *   :
```

# Intron phase

exon                intron                exon
| ATG  GTG |                | CGT  ACA |
                    phase 0

| ATG  GT |                | G  CGT  ACA |
                    phase 1

| ATG  G |                | TG  CGT  ACA |
                    phase 2

# Conserved intron positions

# Patterns of conservation

# Intron positions shared with animals or plants

# Phylogenetic signal in intron positions

Species Tree
30 proteins

Parsimony

reconstruc-
tion

| | |
|---|---|
| afum | afum |
| ater | ater |
| anid | anid |
| snod | cglo |
| cglo | pans |
| ncra | ncra |
| pans | fgra |
| mgri | mgri |
| fgra | snod |
| agos | agos |
| klac | cgla |
| cgla | klac |
| scer | dhan |
| dhan | scer |
| ylip | ylip |
| spom | umay |
| ccin | spom |
| pchr | ccin |
| cneo | pchr |
| umay | cneo |
| rory | rory |
| frub | frub |
| hsap | hsap |
| mmus | mmus |
| atha | atha |

# Intron position reconstruction

- 3 Methods
  - Roy and Gilbert. 2005
  - Csũrös. 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
  - More introns than previously reported
- Intron loss has dominated among the fungi
  - Hemiascomycota experienced loss
- Sampling can bias interpretations - all fungi are not equal.

# Fungal comparative genomics

Evolution of fungal introns

Fungal gene family evolution

Recent intron loss in *C. neoformans*

A    B    C    D    E

Animals          Fungi

# Mechanism of intron loss

- *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?

- 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.

- In *S. cerevisiae* most intron loss events occurred too long ago so little evidence supporting any mechanism

# Sequenced *Cryptococcus* genomes



C. gattii, strain WM276

C. gattii, strain R265

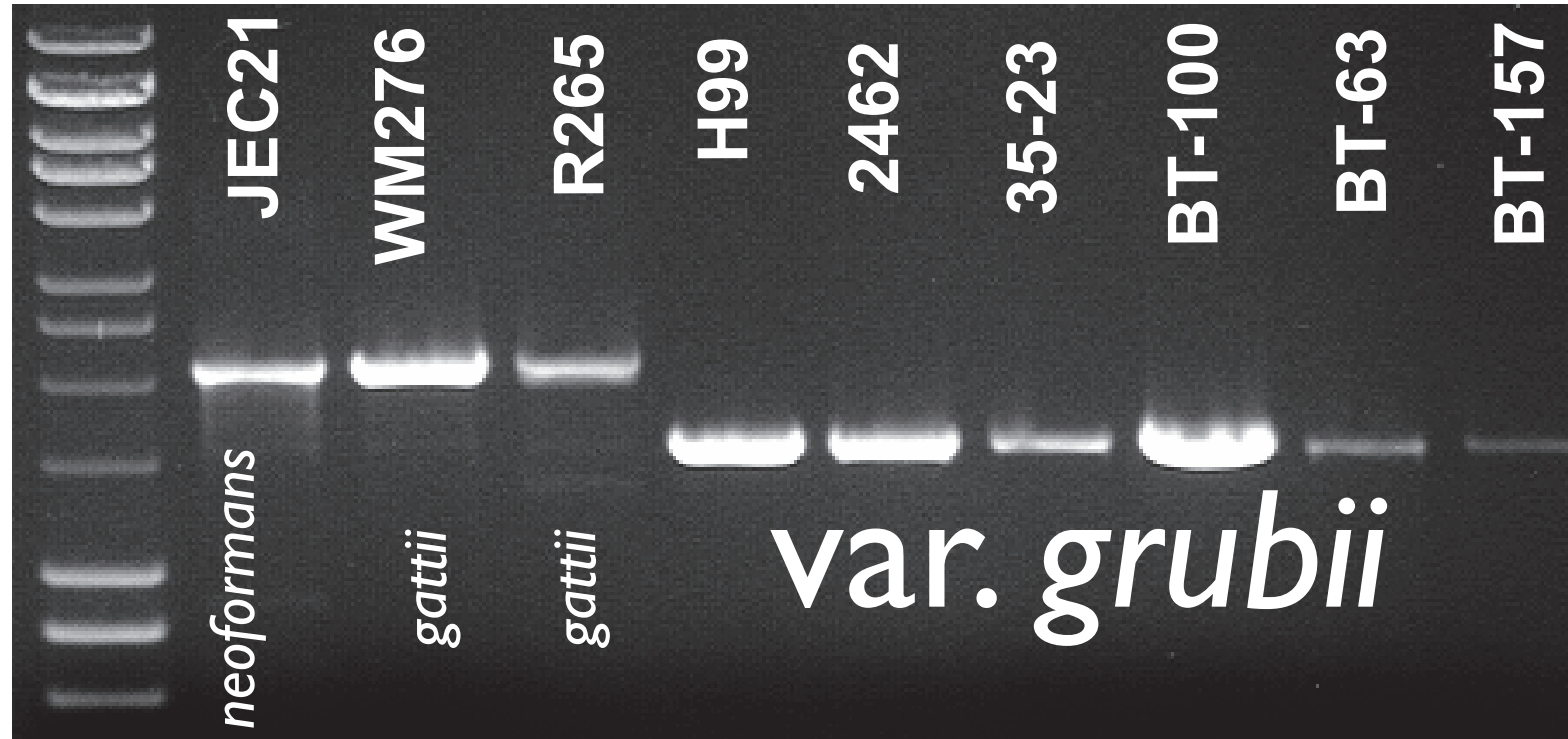C. neoformans var. neoformans, strain JEC21

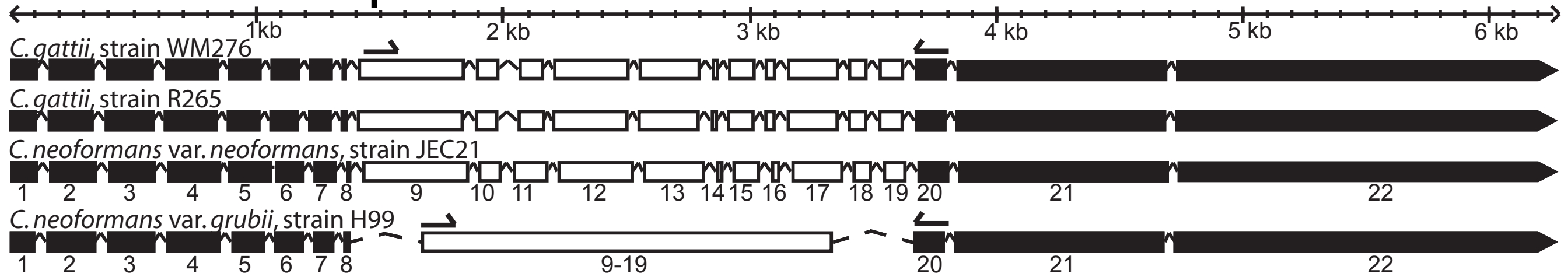C. neoformans var. grubii, strain H99

0.1

substitutions

# 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

Stajich and Dietrich. Euk Cell *In press.*

# Intron loss in var. *grubii*
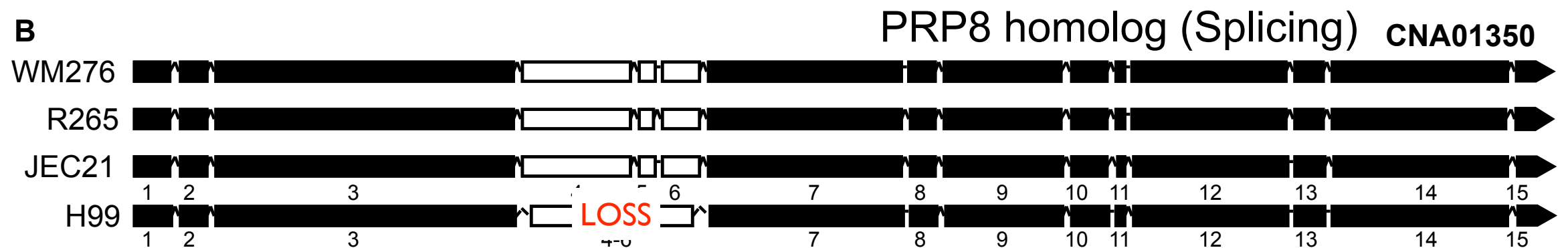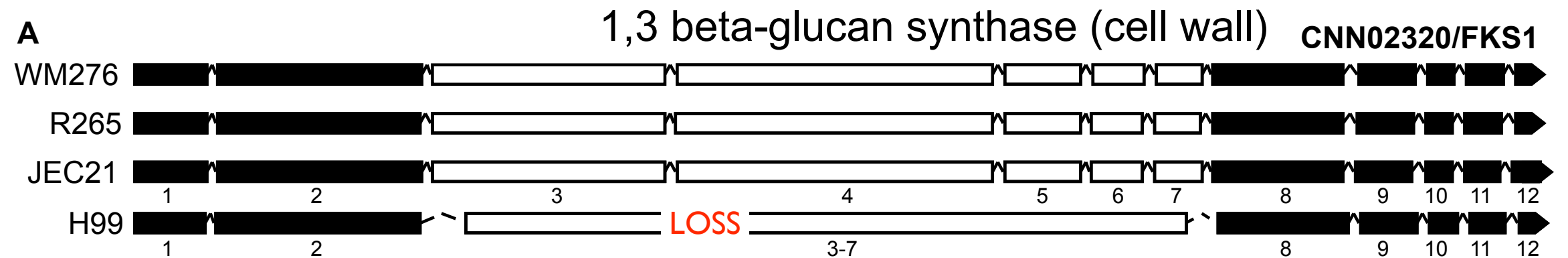


CNI01550 - putative RNA helicase

# Intron loss was a precise excision

```
R265        CGACAAGTACATAAAACTTTTTTTTGTGCCTGGCGCAAAGACTTTCCATTGCTGACAGAAAACAGGTTGAA
WM276       AGACAAGTACATAAAACTTTTTTTTGTGTCTCGTCCAAAGATTTTTCATTGCTGACAGAAAACAGGTTGAA
H99         AGACAA←─────────────────── Intron Missing ───────────────────→-GTTGAA
JEC21_CDS   AGACAA-----------------------------------------------------------GTTGAA
JEC21       AGACAAGTACATACTAGTCCTTGTG---CTATCCCAAAGACTTT-CATTGCTGACAGAAAACAGGTTGAA
                  *****                                              ******
```

```
R265        CGCTGCCGAATTATGTCGATGTTGGAGATTTCTTGAGGTAAGCAACAGACTCGTAACAGCTTGTTCGGTC
WM276       CGCTGCCGAACTATGTCGATGTTGGAGATTTCTTGAGGTAAGCAACAGACTCGTAACAGCTTGTTCGGTC
H99         CCCTGCCGAATTATGTCGACGTTGGAGATTTCTTGAG---------------------------------
JEC21_CDS   CCCTGCCGAATTATGTCGATGTTGGAGATTTCTTGAG---------------------------------
JEC21       CCCTGCCGAATTATGTCGATGTTGGAGATTTCTTGAGGTACGTCGCAAACTCGTAACAGCTTGTTCGATC
            *  ******** ******* ****************
```

Stajich and Dietrich. Euk Cell *In press.*

# Other examples of loss



**A**             1,3 beta-glucan synthase (cell wall)   **CNN02320/FKS1**

WM276

R265

JEC21

H99            LOSS

**B**             PRP8 homolog (Splicing)   **CNA01350**

WM276

R265

JEC21

H99            LOSS

**C**             Ubiquitin protein ligase (protein degradation)   **CNG04610**

WM276

R265

JEC21            GAIN

H99

Stajich and Dietrich. Euk Cell *In press.*
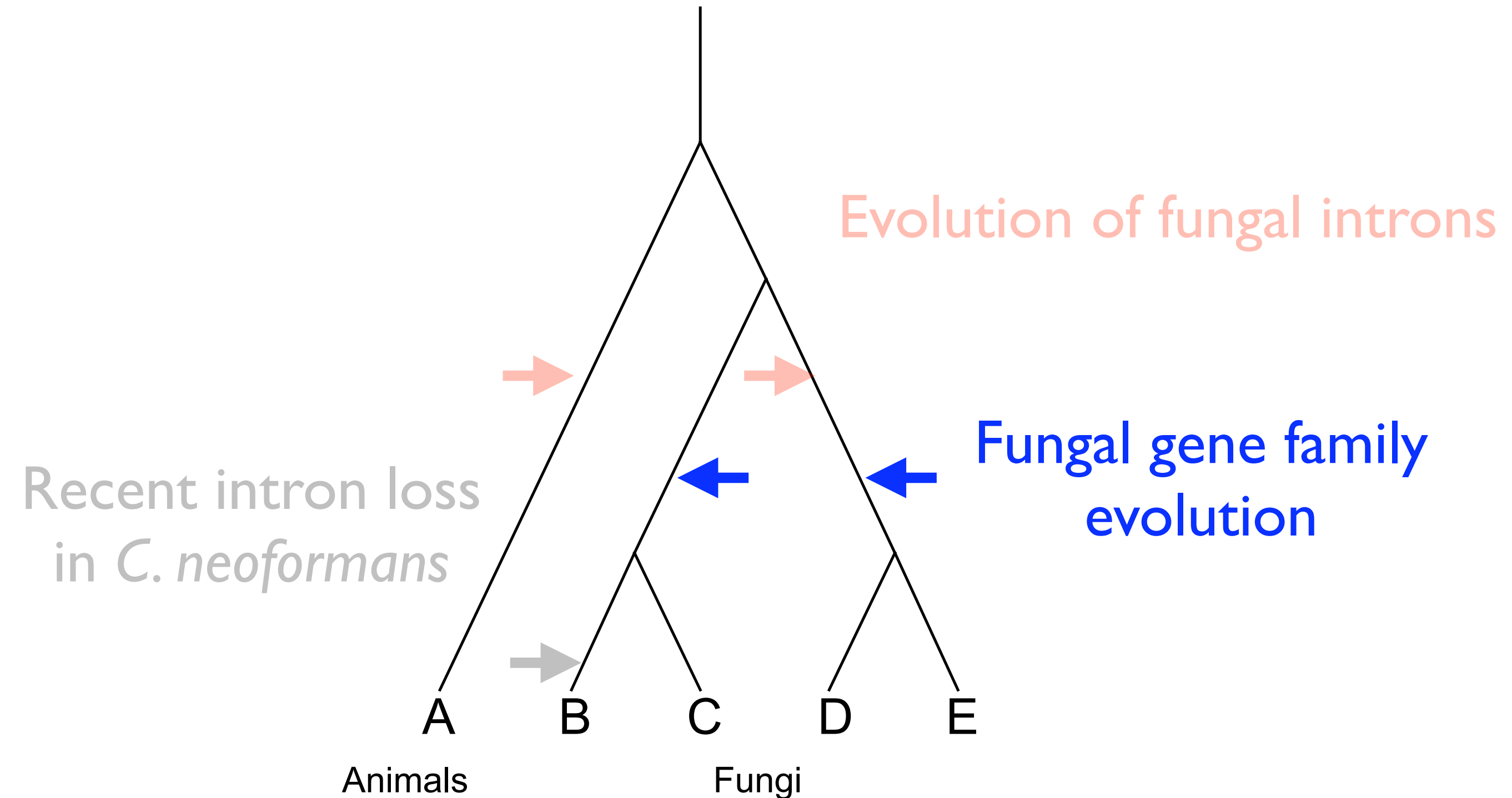
# Conclusions

- Intron loss via homologous recombination with spliced transcript

  - Large losses are all adjacent introns

  - Precise deletion

- Loss biased towards the middle of gene not 3'

# Fungal comparative genomics

Evolution of fungal introns

Fungal gene family evolution

Recent intron loss in *C. neoformans*

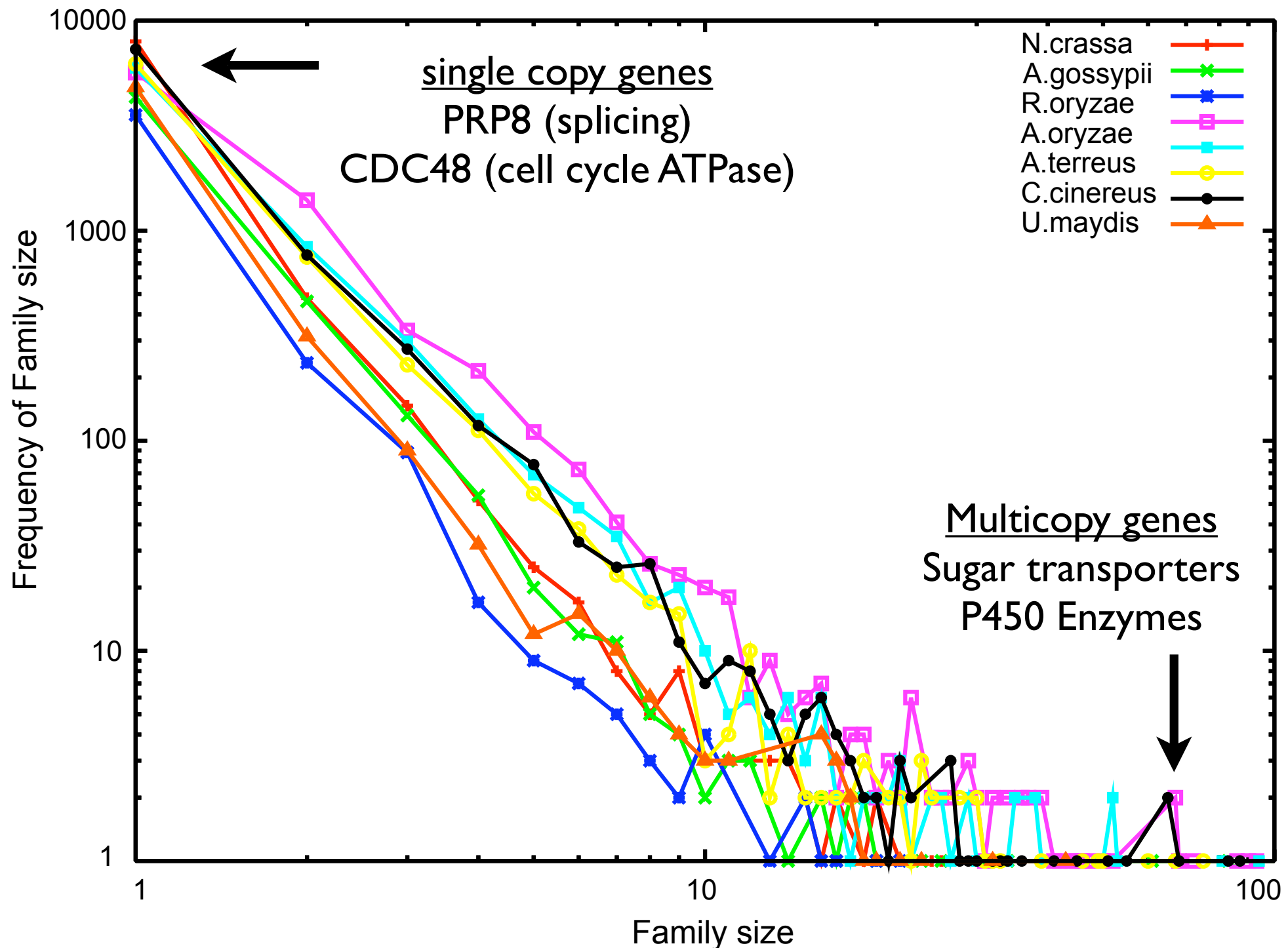A          B     C     D          E

Animals              Fungi

# 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 are have unexpectedly large changes in size across a phylogeny?

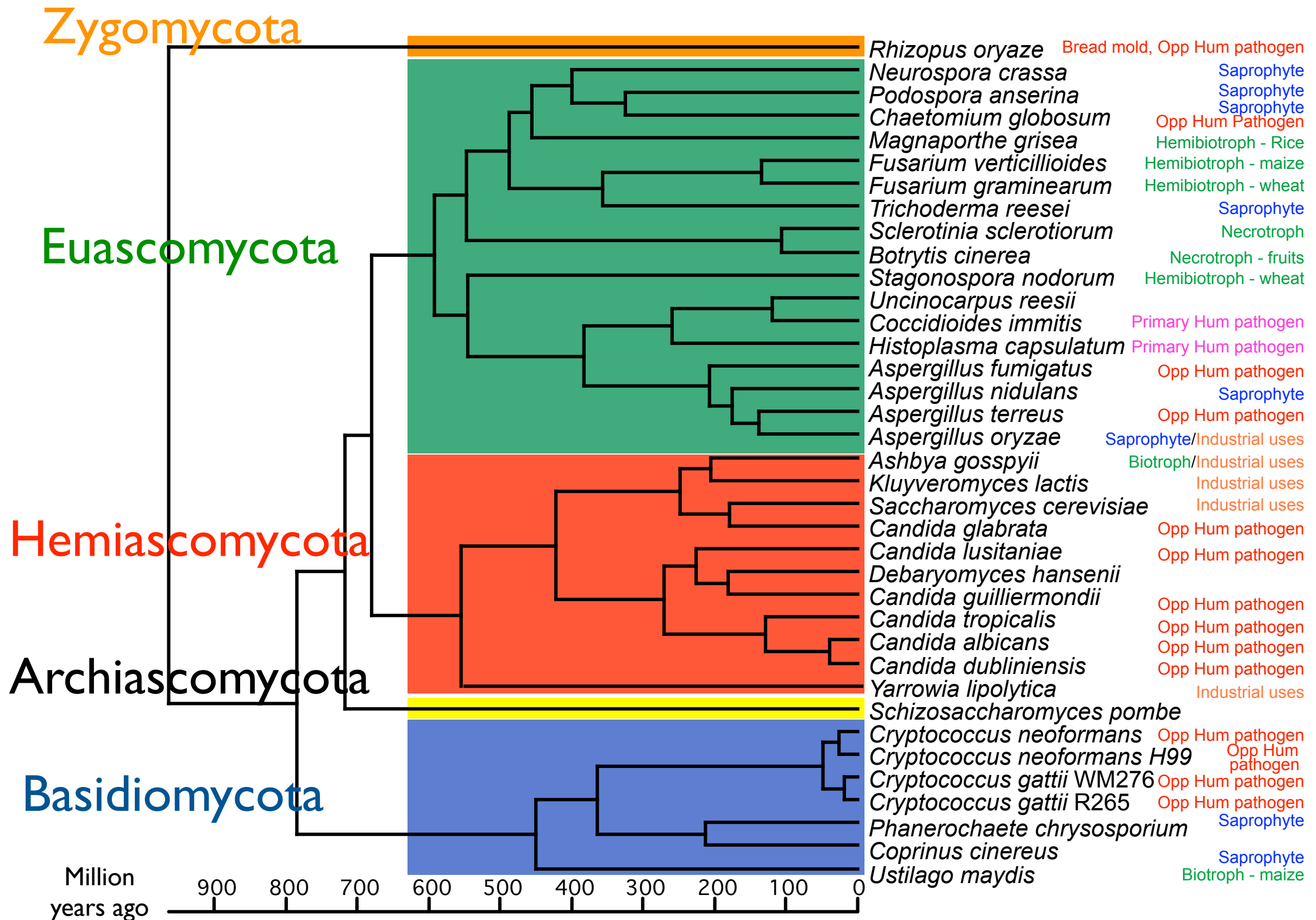  - Follow up these families with more focused studies

# Identifying 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 el, *Science* 2002).

# Gene family sizes follow power law distribution

Fully sequenced fungal genomes

# 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 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
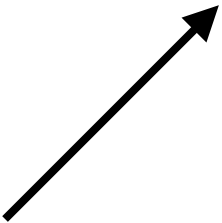De Bie, et al *Bioinformatics* 2006
Demuth et al, *submitted*

# Families with significant expansions
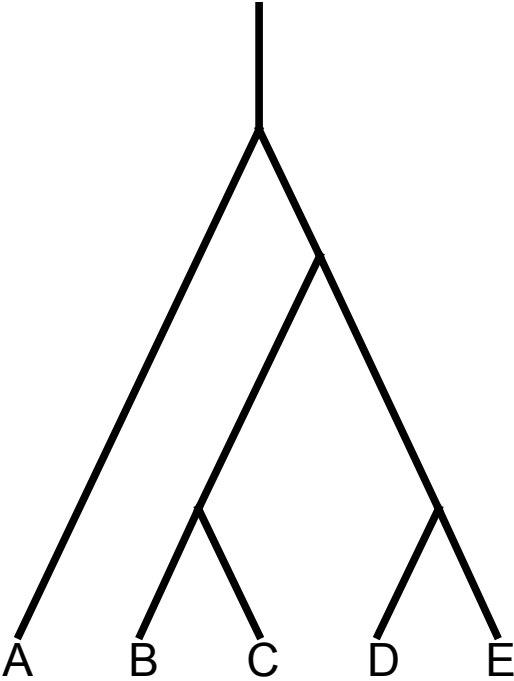
Transporters
Kinases
P450
Oxidation

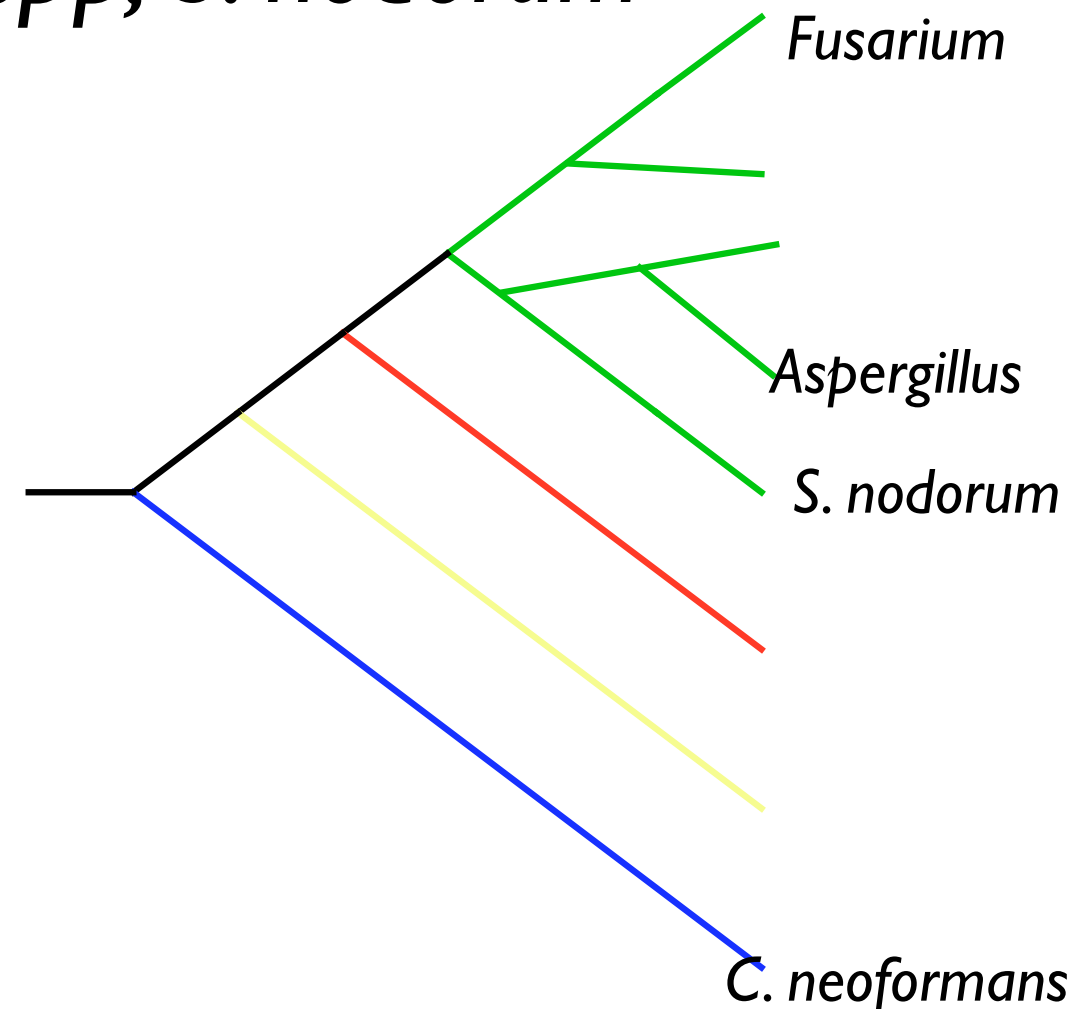| | |
|---|---|
| Vitamin & Cofactor transport | Methytransferase |
| Lactose & sugar transport | Cytochrome P450: CYP64 |
| Amine transport | Cytochrome P450: CYP53,57A |
| Myo-instol, quinate, and glucose transport | Cytochrome P450 |
| Oligopeptide transport | Kinase |
| ABC transporter | Subtilase family |
| MFS, drug pump, & sugar transport | NADH flavin oxidoreductase |
| Transport | Aldehyde dehydrogenase |
| Monocarboxylate & sugar transport | Aldo/kedo reductase |
| ABC transport | Multicopper oxidase |
| Amino acid permease | 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

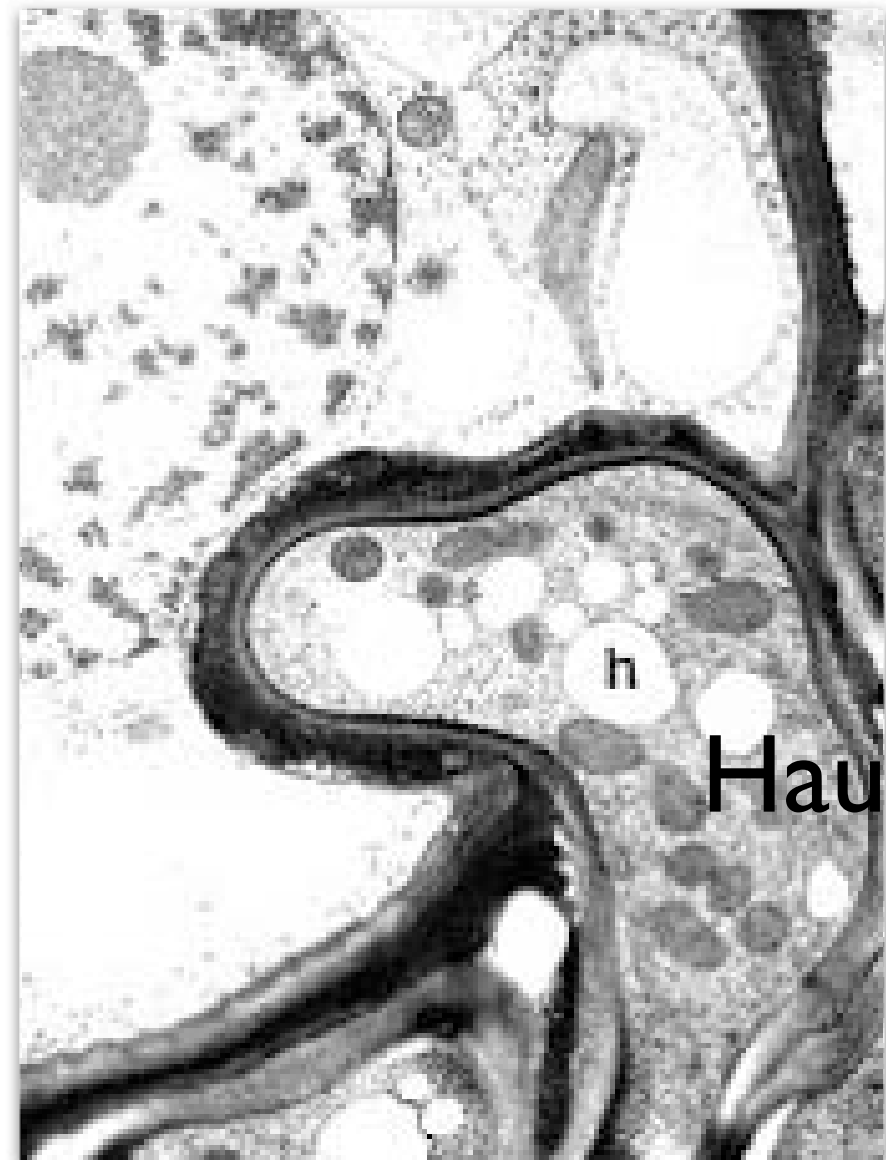- ABC transporters (ATP Binding Cassette)

# Branches with transporter expansions

- Sugar related, Drug pump, and MFS

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

  - Euascomycota

- Vitamin transport

  - *C.neoformans, Fusarium*

  - *A. nidulans* (Biotin)

# What do phytopathogens use transporters for?

- Sugar transporters are used to extract nutrients from host

  - Haustorium: specialized structure for plant parasitism

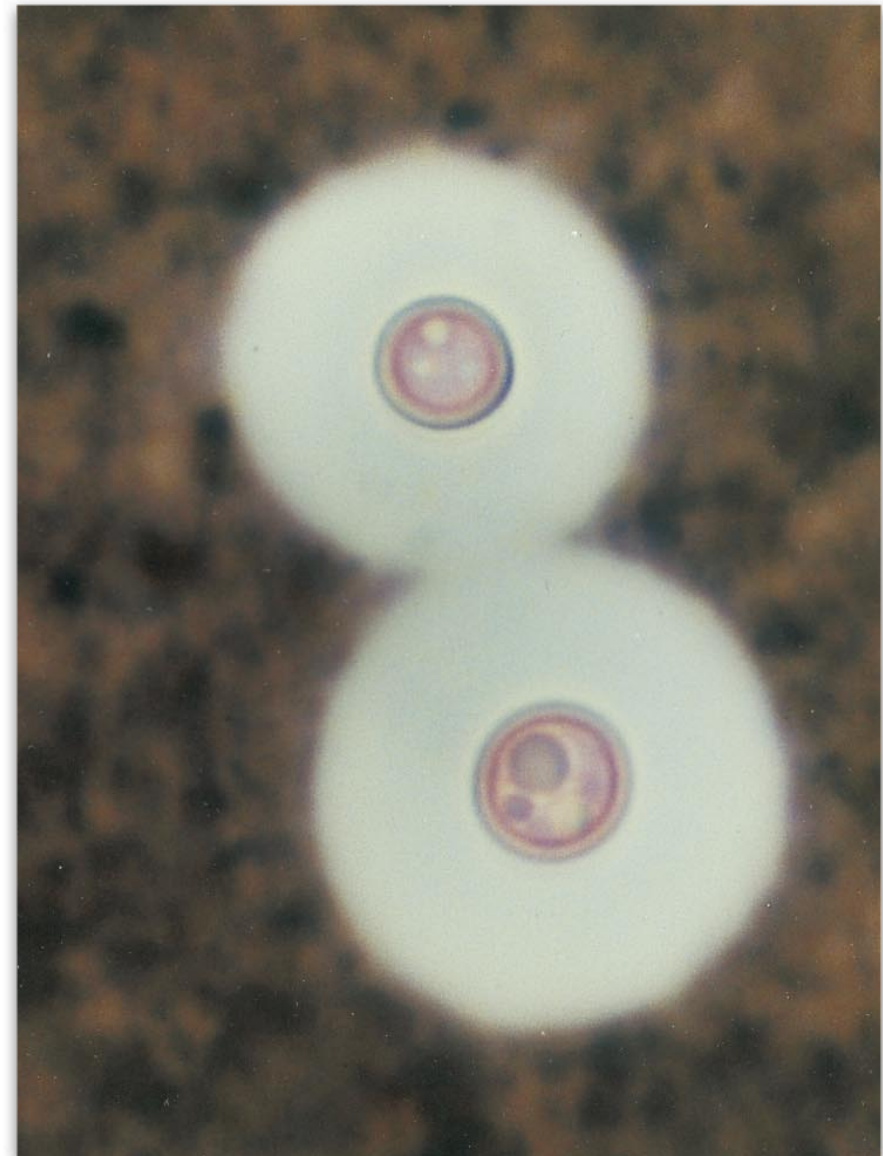  - Many sugar transporters highly and specifically expressed in haustoria



Haustorium

Robert Bauer   http://tolweb.org/

# Cryptococcus sugar transporters

- 3x as many sugar transporters in *C. neoformans* than other basidiomycetes

- "sugar coated killer"

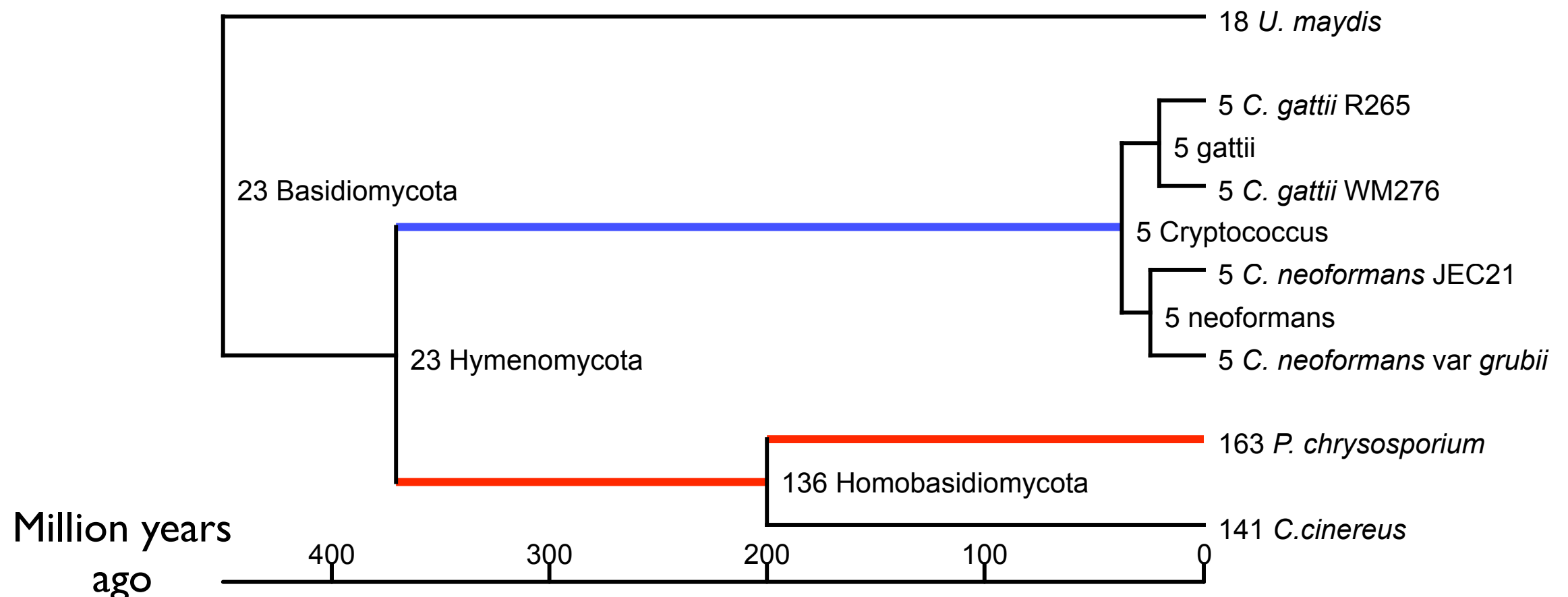- Capsule is mixture of glucose, xylose, and mannose.



Zerpa et al, 1996

# CYP64 was from independent duplication



C. cinereus expansion

P. chrysosporium expansion

Mario Cervini

Tom Volk

# Local duplications created CYP64 expansion



pchr_24

9k  10k  11k  12k  13k  14k  15k  16k  17k  18k  19k  20k  21k  22k  23k  24k

**GLEAN models**

GLEAN_02414
Probability 1

GLEAN_02415
Probability 0.999937

GLEAN_02416
Probability 0.646357

GLEAN_02417
Probability 0.990598

**Pfam domains**

p450
Cytochrome P450 evalue:1e-28

p450
Cytochrome P450 evalue:6e-26

p450
Cytochrome P450 evalue:6.3e-23

p450
Cytochrome P450 evalue:9e-07

http://fungal.genome.duke.edu

# Family size contractions

- *Histoplasma, Coccidioides* many families

- Hemiascomycetes - P450

- *C. neoformans* - P450

- *U. maydis* - Lactose transport

# 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?

# Overall conclusions

- Multiple genome sequences have helped resolve several outstanding questions in evolution introns

- Gene family expansions can be important in identifying molecular basis for adaptation

# Future directions

- UC Berkeley with John Taylor

- Adaptation and speciation in fungi

- Focus on pathogenic fungus *Coccidioides*

  - Signatures of adaptation among genomes of 12 sequenced strains

# Acknowledgments

Fred Dietrich

Greg Wray
Lincoln Stein
John McCusker
Alex Hartemink

Joseph Heitman
Marcy Uyenoyama
Rytas Vilglyas
Tim James
Pat Pukkilla

Matthew Hahn
DUMRU
John Perfect
Andy Alspaugh
Tom Mitchell

Ewan Birney

UPGG

Robert Cramer
James Fraser
Steven Giles
Alex Idnurm
Scott Roy

Dave Des Marais
Heath O'Brien
Matt Rockman
MDG
Dietrich Lab
Andria Allen
Stephanie Diezmann
Charles Hall
Shan Huang
Philippe Lüdi
Laura Kavanaugh
Sandra Reynolds
Mark DeLong