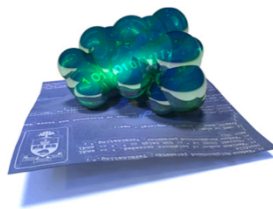# Analyzing Phylogenetic Trees
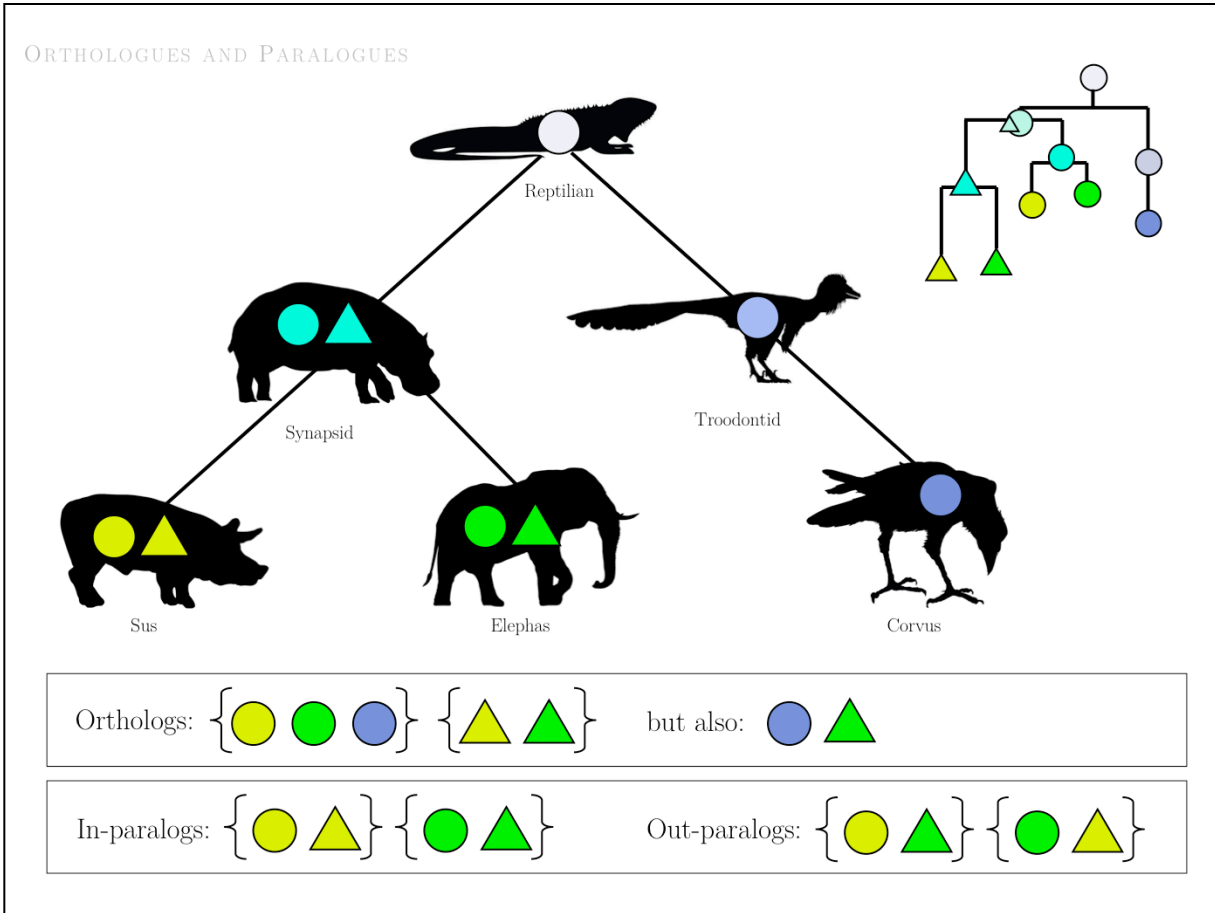
Boris Steipe

DEPARTMENT OF BIOCHEMISTRY – DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO

# Analyzing mixed gene trees

The **interpretation** of phylogenetic trees involves assigning all branchpoints to *speciation* or *duplication* events, and reconciling the inferred topology and assignment with the (known) speciation tree of the organisms that are represented.

**Definitions**:

- **orthologues** arise through **speciation;**

- **paralogues** arise through **duplicaion.**

cf. Koonin EV. (2005) Orthologs, paralogs, and evolutionary genomics. *Annu Rev Genet.* **39**:309-338.

The literature distinguishes between in-paralogues (within one species) and out paralogues (between species). However there is no distinct term for the **RBM gene** (or BBH – Bidirectional Best Hit), one just assumes that RBMs **are** orthologues although even in bacteria this is true in (only) about 95% of the cases.

cf. Wolf & Koonin (2012) A tight link between orthologs and bidirectional best hits in bacterial and archaeal genomes. *Genome Biol Evol.* **4**:1286-94.

## Orthologues are derived from **a single ancestral gene**[1] in **the last common ancestor**[2] of **two species**[3].
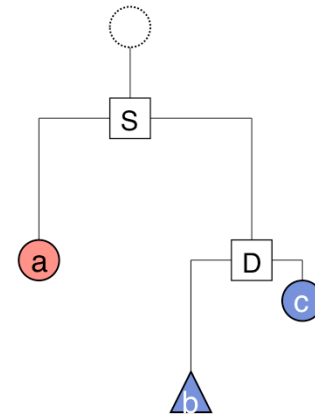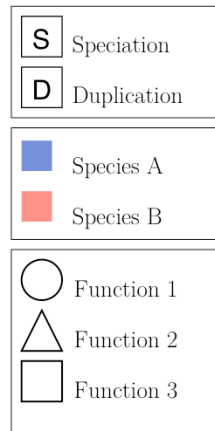
(1) If two genes are derived from paralogues, they are considered paralogues even if the orthologues have been lost. (*"pseudoorthologues"*)

(2) We consider the **last** common ancestor, not some more ancient one.

(3) Orthology is a property of pairs of genes, not multiples.

(NOT) Orthology is not (necessarily) a one-to-one relationship.

(NOT) Orthology is not a transitive relationship, the equivalence relation of homology does not (necessarily) hold true!

| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ◼ | Species A |
| ◼ | Species B |

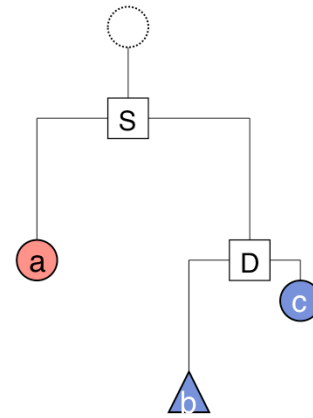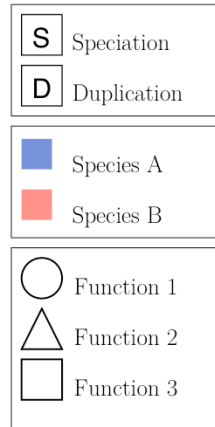| | |
|---|---|
| ○ | Function 1 |
| △ | Function 2 |
| □ | Function 3 |

c and a are orthologues, a and b are orthologues, but b and c are NOT.

To analyze mixed gene trees, we evaluate each branch point and determine whether it represents a speciation or a duplication event.

**Speciation events give rise to orthologues, duplication events give rise to paralogues.**

4

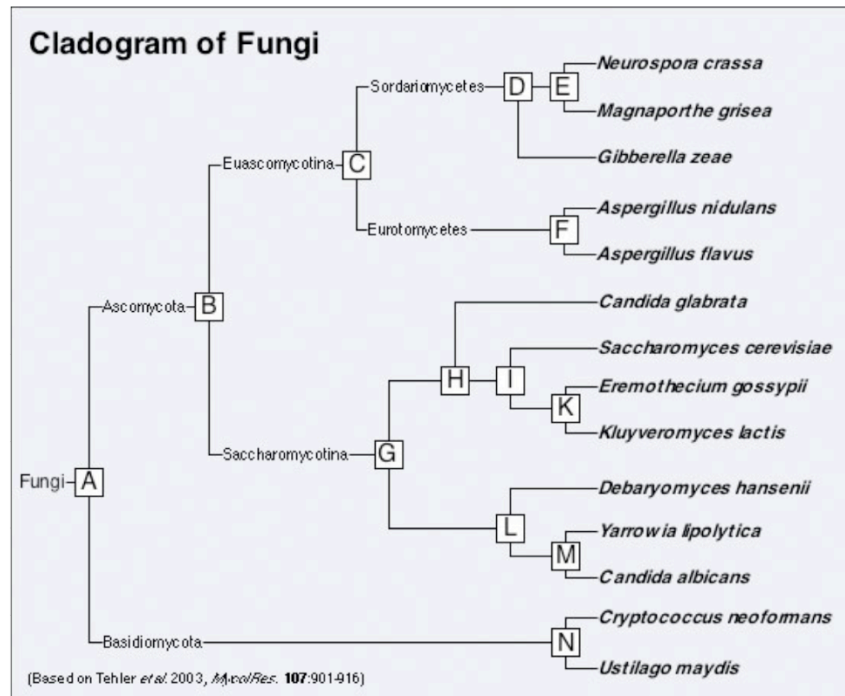# Paralogues are genes within a species that are descendants of one gene that has duplicated.

Note that both paralogues duplicating AFTER speciation are orthologues to a gene in another species that has not duplicated. These have been called **Inparalogs** and they are collectively **Co-orthologous** to a comparison gene in a different species. However, only one of these fulfills the *reciprocal best match* criterion.

| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ■ | Species A |
| ■ | Species B |

| | |
|---|---|
| ○ | Function 1 |
| △ | Function 2 |
| ▢ | Function 3 |

b and c are paralogues (inparalogs) and they are co-orthologues to a. Only a and c will be RBM.

For orthologous genes, the gene tree should recapitulate the sequence of speciation events.

**Cladogram of Fungi**
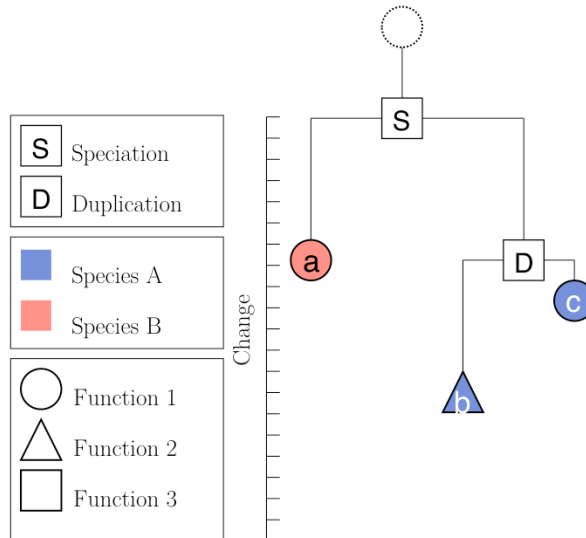
- Fungi–A
  - Ascomycota–B
    - Euascomycotina–C
      - Sordariomycetes–D–E
        - Neurospora crassa
        - Magnaporthe grisea
        - Gibberella zeae
      - Eurotomycetes–F
        - Aspergillus nidulans
        - Aspergillus flavus
    - Saccharomycotina–G
      - H–I
        - Candida glabrata
        - Saccharomyces cerevisiae
        - Eremothecium gossypii–K
        - Kluyveromyces lactis
      - L–M
        - Debaryomyces hansenii
        - Yarrowia lipolytica
        - Candida albicans
  - Basidiomycota–N
    - Cryptococcus neoformans
    - Ustilago maydis

(Based on Tehler *et al.* 2003, *Mycol. Res.* **107**:901-916)

*Reciprocal best match* finds the respectively most similar genes in a genome. If the acquisition of a new function involves a period of accelerated evolution, duplicated genes that retain their function will be more similar to each other than those that change.
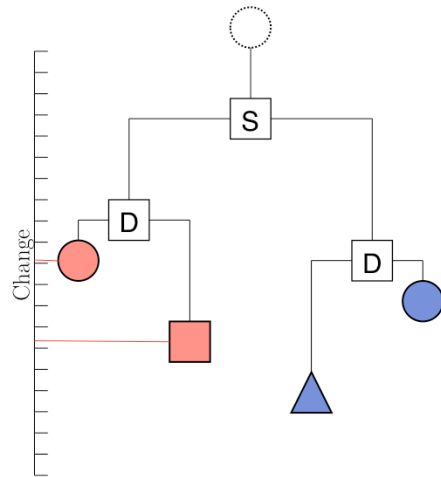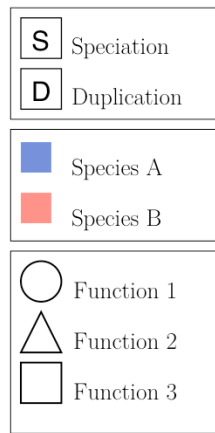
## Why does *reciprocal best match* find orthologues?

Because we assume that after gene duplication the two descendants evolve at different rates. One paralogue will evolve at a faster rate from the ancestor than the other.

Reciprocal best match finds the respectively most similar genes in a genome.

| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ■ | Species A |
| ■ | Species B |

| | |
|---|---|
| ○ | Function 1 |
| △ | Function 2 |
| □ | Function 3 |

# Why does *reciprocal best match* find functionally most similar orthologues?

We assume that the faster evolving of two paralogues undergoes neo- and/or subfunctionalization.

| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ▦ (blue) | Species A |
| ▦ (red) | Species B |

| | |
|---|---|
| ◯ | Function 1 |
| △ | Function 2 |
| ▢ | Function 3 |

Change

Simply add up the vertical path lengths in this tree between all pairs across the speciation node to show that RBM indeed finds the two genes with conserved Function 1, and not the pairs with either of the novel Functions 2 or 3.

8

Organism A

Organism B

Organism C

Organism D

Protein 1$^A$

Protein 1$^B$

Protein 1$^C$

Protein 1$^D$

If you consider a single, orthologous gene, you should expect that the gene tree recapitulates at least the topology of the cladogram, hopefully also the branch lengths of the phylogram.
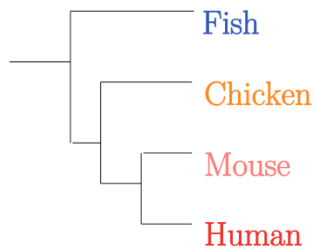
Sometimes we calculate phylogenetic trees from both orthologues and paralogues, to compile all available information into the same framework. What do we expect?
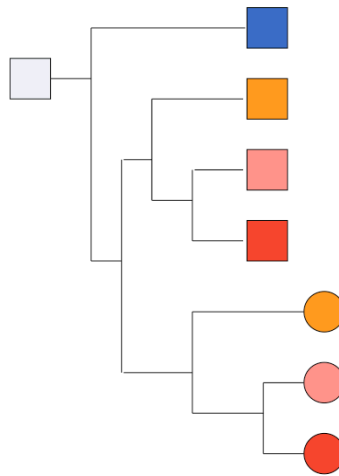
If there is a gene duplication somewhere in the tree:
- we expect all descendants to inherit the duplicated gene
- we expect orthologues within the duplicated section of the tree to pattern like the species tree.
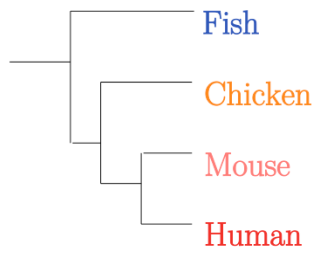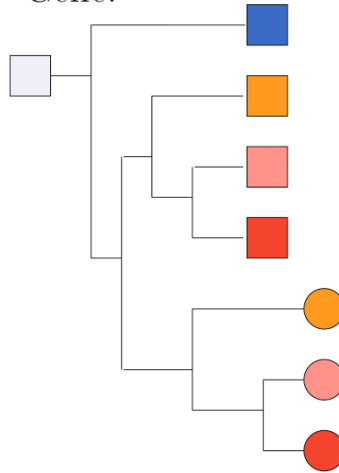
Species:

Gene:



Fish

Chicken

Mouse

Human

The duplication point lies **before the cenancestor of the included species** and **after the cenancestor of included/excluded species.**
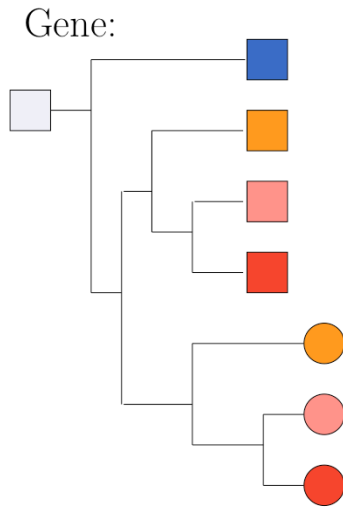
Species:

Gene:

Fish

Chicken

Mouse

Human

Gene-loss, lack of resolution of branching points and HGT can complicate the analysis.
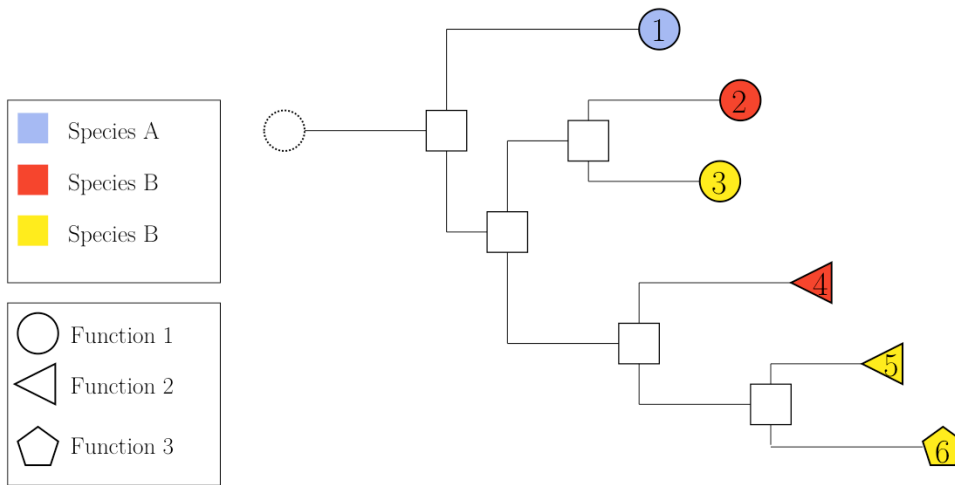
11

To analyse a mixed gene tree:     Gene:
- we need a species tree for
  reference
- we need to keep track of
  species in the tree
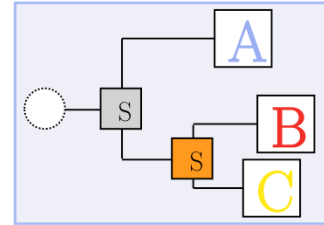- we need to keep track of
  groups of orthologues

Assume we have the following mixed gene tree.

How did it evolve? What Speciations and duplications led to this tree?
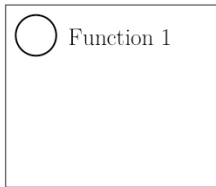


Species A
Species B
Species B

Function 1
Function 2
Function 3
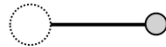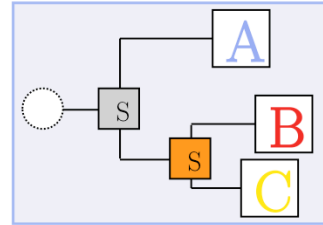
Always start from a reference model - the phylogram (or at least cladogram) of species. Whatever else happens over time, all variation is constrained by the periods of *joint* or *independent* evolution that is due to the *sequence of speciation events*, represented by the branching topology of the tree.

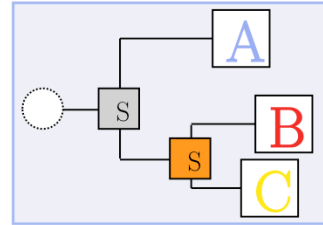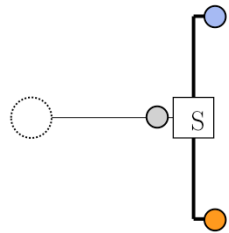| | |
|---|---|
| ■ (blue) | Species A |
| ■ (red) | Species B |
| ■ (yellow) | Species B |

Some gene evolves through the mists of time, part of–
and passenger in– an organism that contains it.

| | |
|---|---|
| ■ (blue) | Species A |
| ■ (red) | Species B |
| ■ (yellow) | Species B |

| | |
|---|---|
| ◯ | Function 1 |

A Speciation generates two orthologues that evolve further by *separately* accumulating mutations.

| S | Speciation |

| | Species A |
| | Species B |
| | Species B |

| ◯ | Function 1 |



16

Further speciation events generate all OTUs .
Obviously we expect the phylogram of orthologues to
closely resemble the species tree.

| S | Speciation |

| | Species A |
| | Species B |
| | Species B |

| ◯ | Function 1 |

Hoever, if there is a duplication event in the tree – in this case *before* the speciation event we just discussed, the story changes: a duplication event generates a **copy** of a gene *in one organism*. Typically, one of the copies may acquire a new function.



| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ■ | Species A |
| ■ | Species B |
| ■ | Species B |

| | |
|---|---|
| ○ | Function 1 |
| ◁ | Function 2 |

Consequently, after the duplication, nothing is changed for *one of the duplicates* as the species separate...

| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ⬛ | Species A |
| 🟥 | Species B |
| 🟨 | Species B |

| | |
|---|---|
| ◯ | Function 1 |
| ◁ | Function 2 |

... while the *other duplicate* evolves according to the same speciation pattern as its sibling (albeit, typically with different rates). Every duplication creates a copy of the branching pattern of the species tree.



| S | Speciation |
|---|------------|
| D | Duplication |

| | |
|---|------------|
| ■ | Species A |
| ■ | Species B |
| ■ | Species B |

| ○ | Function 1 |
|---|------------|
| ◁ | Function 2 |

The final tree can be decomposed into full or partial copies of the species tree. In real life, *uncertain branching order* and *gene loss* complicates the tree. But we can always look for the best match to the model.



| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ▪ | Species A |
| ▪ | Species B |
| ▪ | Species B |

| | |
|---|---|
| ◯ | Function 1 |
| ◁ | Function 2 |
| ⬠ | Function 3 |

Note: the number of *duplication nodes* corresponds exactly to the number of duplication events. The number of *speciation nodes* is determined by the number and size of branches that have been inserted.

| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ■ | Species A |
| ■ | Species B |
| ■ | Species B |

| | |
|---|---|
| ○ | Function 1 |
| ◁ | Function 2 |
| ⬠ | Function 3 |

different *nodes* may represent the same speciation *event*

To *analyze* a tree, you perform this process in reverse!

Label all speciation and duplication events in this tree. Define sets of most similar orthologues and of paralogues.



| | |
|---|---|
| S | Speciation |
| D | Duplication |

| | |
|---|---|
| ■ | Species A |
| ■ | Species B |
| ■ | Species B |

| | |
|---|---|
| ◯ | Function 1 |
| ◁ | Function 2 |
| ⬠ | Function 3 |

To *analyze* a tree, you perform this process in reverse!

Label all speciation and duplication events in this tree. Define sets of most similar orthologues and of paralogues.



| | Speciation |
|---|---|
| S | |
| D | Duplication |

| | |
|---|---|
| ▪ | Species A |
| ▪ | Species B |
| ▪ | Species B |

| | |
|---|---|
| ◯ | Function 1 |
| ◁ | Function 2 |
| ⬠ | Function 3 |

Most similar orthologues:
$\{1, 2, 3\}, \{4, 5\}$

Paralogues:
$\{2, 4\}, \{3, 5, 6\}$

# gene tree (mixed)

```
                                +-GIBZE A
                        +-16+-1+-MAGGR C
                     +-19  +-NEUCR B
                  +-24  +-ASPNI C
                     +-YARLI A
            +-25      +-2+CANAL E
                  +-22   +DEBHA C
                            +-KLULA C
                  +-23  +-21+-1+-CANGL C
      +-----27          +-SACCE PHD1
                     +EREGO A
                        +-CANAL C
   +--30        +-2+-18+DEBHA E
                  +---CANGL D
            +------1+---CANAL A
      +-37        +---DEBHA D
35-43
            +---3+CANAL B
                 +-DEBHA A
         +-53      +--CANGL A
               +-3+-SACCE SWI4
         +--3+      +---EREGO C
               +-3+----KLULA A
                     +ASPFU MBP1
               +-4+ASPNI MBP1
            +-49  +-GIBZE MBP1
         +-50   +-MAGGR MBP1
               +-4+-NEUCR MBP1
         +-54    +CANAL MBP1
               +-3+-DEBHA MBP1
         +-51      +-CANGL MBP1
               +-4+SACCE MBP1
            +-46   +-EREGO MBP1
   +-56        +-4+-KLULA MBP1
         +-57    +---CRYNE MBP1
               +---YARLI MBP1
         +-5+      +ASPFU A
               +-3+ASPNI D
            +-42      +GIBZE B
         +-58    +-4+-3+-MAGGR D
            +-47  +-40  +--NEUCR C
               +-52   +--YARLI B
         +-60  +--CRYNE A
            +-5+---USTMA B
               +---USTMA MBP1
```

A selection of sequences (e.g. the result of
a PSI-BLAST search contains orthologues
and paralogues in the same tree. In order
to interpret such a tree, we MUST
distinguish between branch points that
represent *duplication events* and branch
points that represent *speciation events*.

# gene tree (mixed)

```
                              +-1+GIBZE A
                         +-16+-1+MAGGR C
                     +-19   +NEUCR B
                 +-24  +ASPNI C
                 +-24 +-YARLI A
          +-25        +-20+CANAL E
          +-25    +-22   +DEBHA C
                  +-22    +-1+KLULA C
          +-23   +-21+-17+CANGL C
  +-----27  +-21    +-SACCE PHD1
           +EREGO A
  +--30    +-18+-CANAL C
          +-26+-18+DEBHA E
           +---CANGL D
       +------13+---CANAL A
                +---DEBHA D
   +-37
35-43
          +---3+CANAL B
          +---3+-DEBHA A
   +-53   +-32+--CANGL A
          +-34+-32+-SACCE SWI4
          +-34 +-33+---EREGO C
                +-33+----KLULA A
          +-4+ASPFU MBP1
       +-49+-4+ASPNI MBP1
   +-50   +-GIBZE MBP1
       +-48+-MAGGR MBP1
   +-56   +-48+-NEUCR MBP1
       +-54  +-38+CANAL MBP1
       +-51  +-38+-DEBHA MBP1
       +-51  +-44+-CANGL MBP1
            +-46+-44+SACCE MBP1
   +-57     +-46+-EREGO MBP1
            +-45+-KLULA MBP1
       +-55+---CRYNE MBP1
          +---YARLI MBP1
                 +-36+ASPFU A
               +-36+ASPNI D
          +-42     +GIBZE B
   +-58   +-47 +-40+-39+-MAGGR D
          +-47 +-40+-39  +--NEUCR C
          +-52 +--YARLI B
   +-60 +---CRYNE A
          +-59+---USTMA B
          +-59+---USTMA MBP1
```
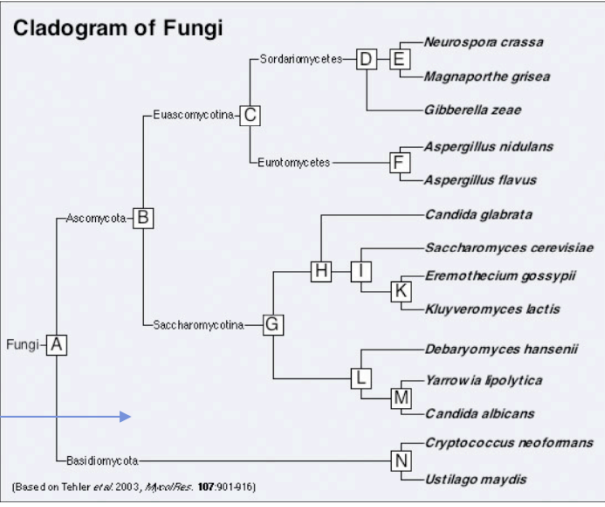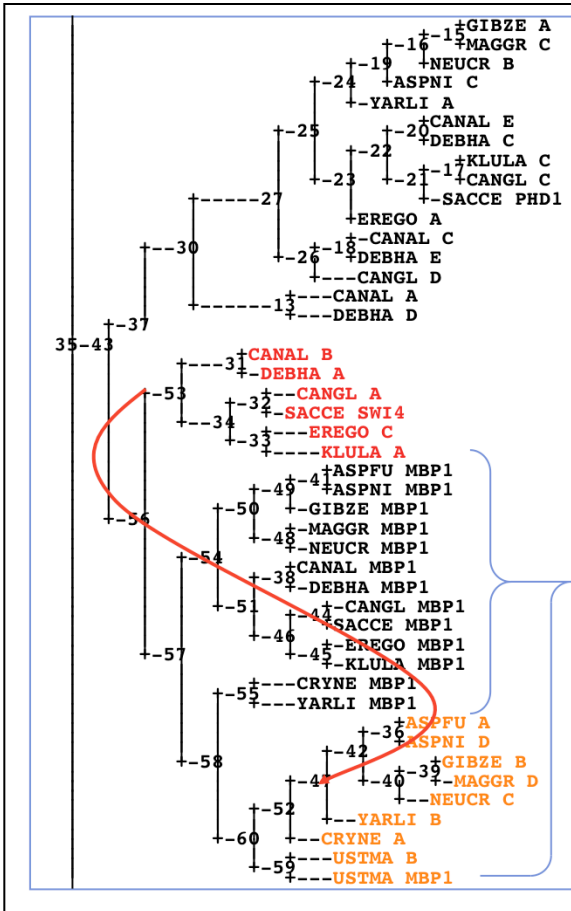
Cladogram of Fungi

Mbp1 orthologues pattern (mostly) according to the species tree in this example.

# gene tree (mixed)

**Cladogram of Fungi**

Euascomycotina — C
Sordariomycetes — D — E — Neurospora crassa / Magnaporthe grisea
Gibberella zeae
Eurotomycetes — F — Aspergillus nidulans / Aspergillus flavus
Ascomycota — B
Saccharomycotina — G — H — Candida glabrata
I — Saccharomyces cerevisiae
K — Eremothecium gossypii / Kluyveromyces lactis
L — Debaryomyces hansenii
M — Yarrowia lipolytica / Candida albicans
Fungi — A
Basidiomycota — N — Cryptococcus neoformans / Ustilago maydis

(Based on Tehler *et al.* 2003, *Mycol Res.* **107**:901-916)

Violations of the species tree may suggest
to re-evaluate the annotations as well as
the topology of the gene tree.

GIBZE A
MAGGR C
-16 -15
-19 NEUCR B
-24 ASPNI C
YARLI A
-25 CANAL E
-20 DEBHA C
-22 KLULA C
-23 -21 -17 CANGL C
SACCE PHD1
-27 EREGO A
-30 CANAL C
-18 DEBHA E
-26 CANGL D
-13 CANAL A
DEBHA D
-37
35-43
-31 CANAL B
DEBHA A
-53 CANGL A
-32 SACCE SWI4
-34 EREGO C
-33 KLULA A
-49 -47 ASPFU MBP1
ASPNI MBP1
-50 GIBZE MBP1
-56 -48 MAGGR MBP1
NEUCR MBP1
-54 CANAL MBP1
-38 DEBHA MBP1
-51 CANGL MBP1
-44 SACCE MBP1
-46 EREGO MBP1
-57 -45 KLULA MBP1
-55 CRYNE MBP1
YARLI MBP1
-36 ASPFU A
ASPNI D
-42 GIBZE B
-58 -41 -40 -39 MAGGR D
NEUCR C
-52 YARLI B
-60 CRYNE A
-59 USTMA B
USTMA MBP1

Errors arise from:

      Sampling (indels!)

      HGT

      Methodology

      Long branch attraction

      ...

Probably not from convergence though.

*Long Branch Attraction* is a pervasive problem of molecular phylogenies.

**Problem:** highly divergent sequences may group together in a tree regardless of their true relationship. This is due to the fact that the number of states is limited, and widely divergent sequences will pick up mutual similarities to the average distribution.
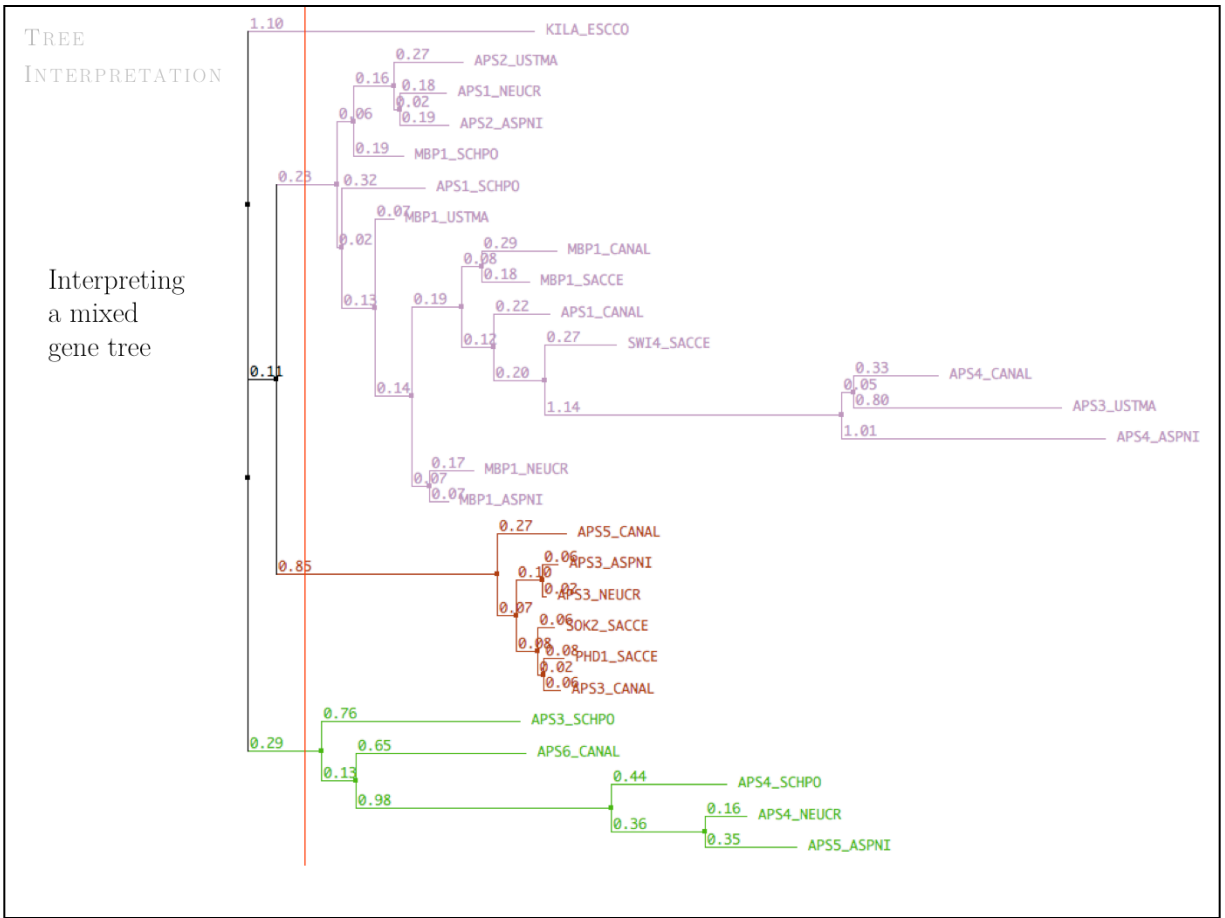
**Symptom:** unexpected grouping patterns and poor bootstrap values. Try to remove, recalculate, re-insert. (The information contributed by a highly diverged sequence to the tree is not very critical anyway.)

**Cures:** ML methods are a bit less sensitive. Correct for multiple substitutions. Try to use slower-evolving traits. Best approach: add intermediary OTU's sequences (always better to add information than to massage the algorithm).

**Even if you are interested only in a few members of a clade, it is good to include as many OTUs as feasible for the tree building.**

See also: http://en.wikipedia.org/wiki/Long_branch_attraction (this article deserves to be rewritten though, sounds a bit like a high-school project).
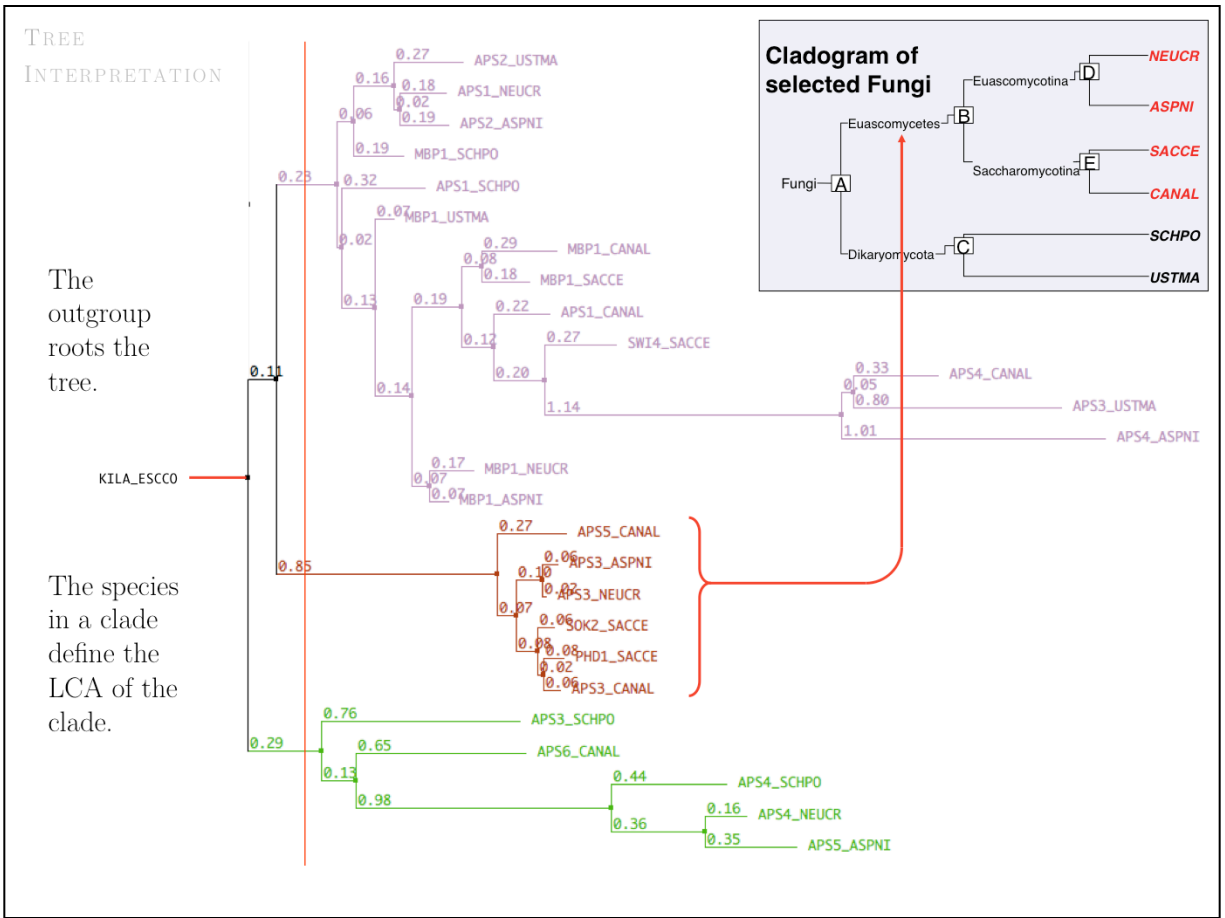
Recent paper: Kück P, Mayer C, Wägele J-W, Misof B (2012) Long Branch Effects Distort Maximum Likelihood Phylogenies in Simulations Despite Selection of the Correct Model. PLoS ONE 7(5): e36593. doi:10.1371/journal.pone.0036593.
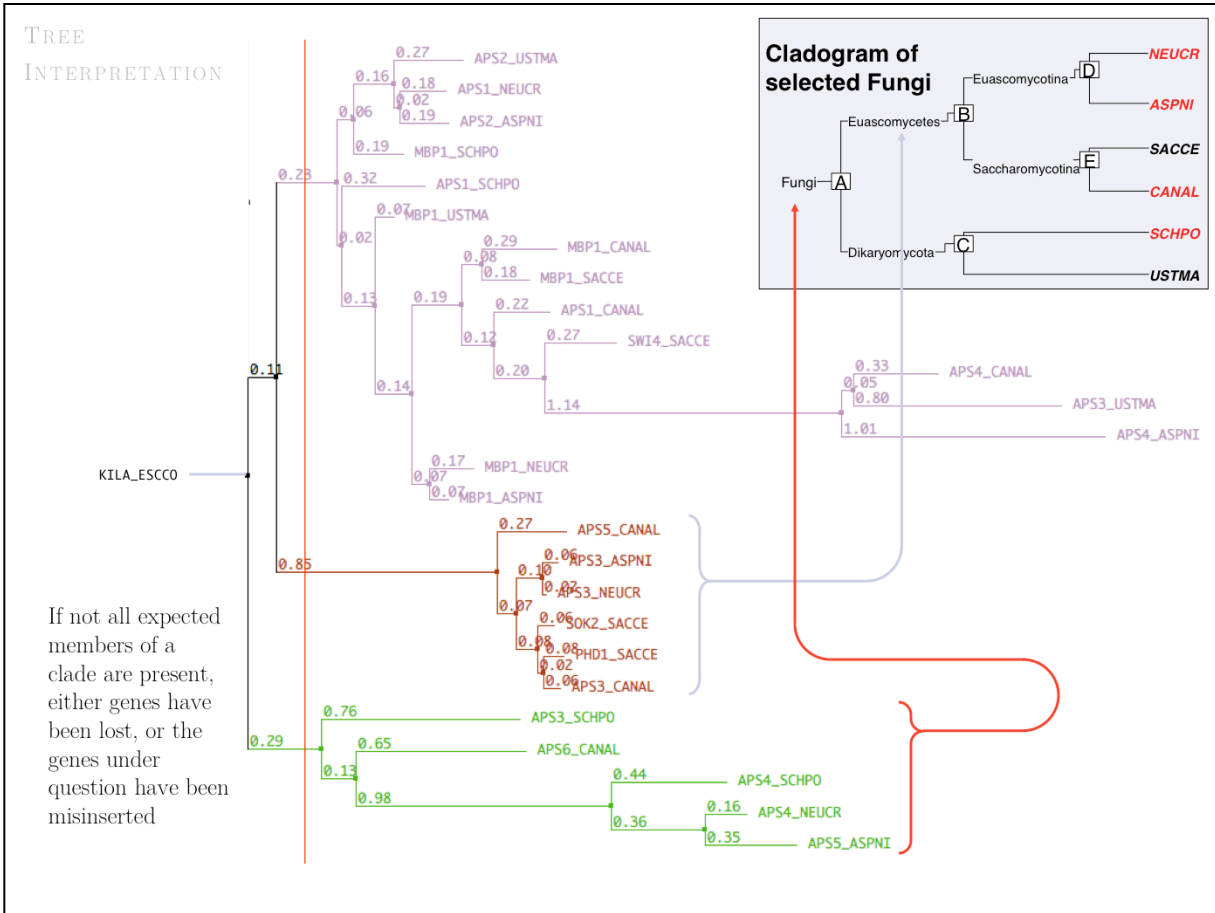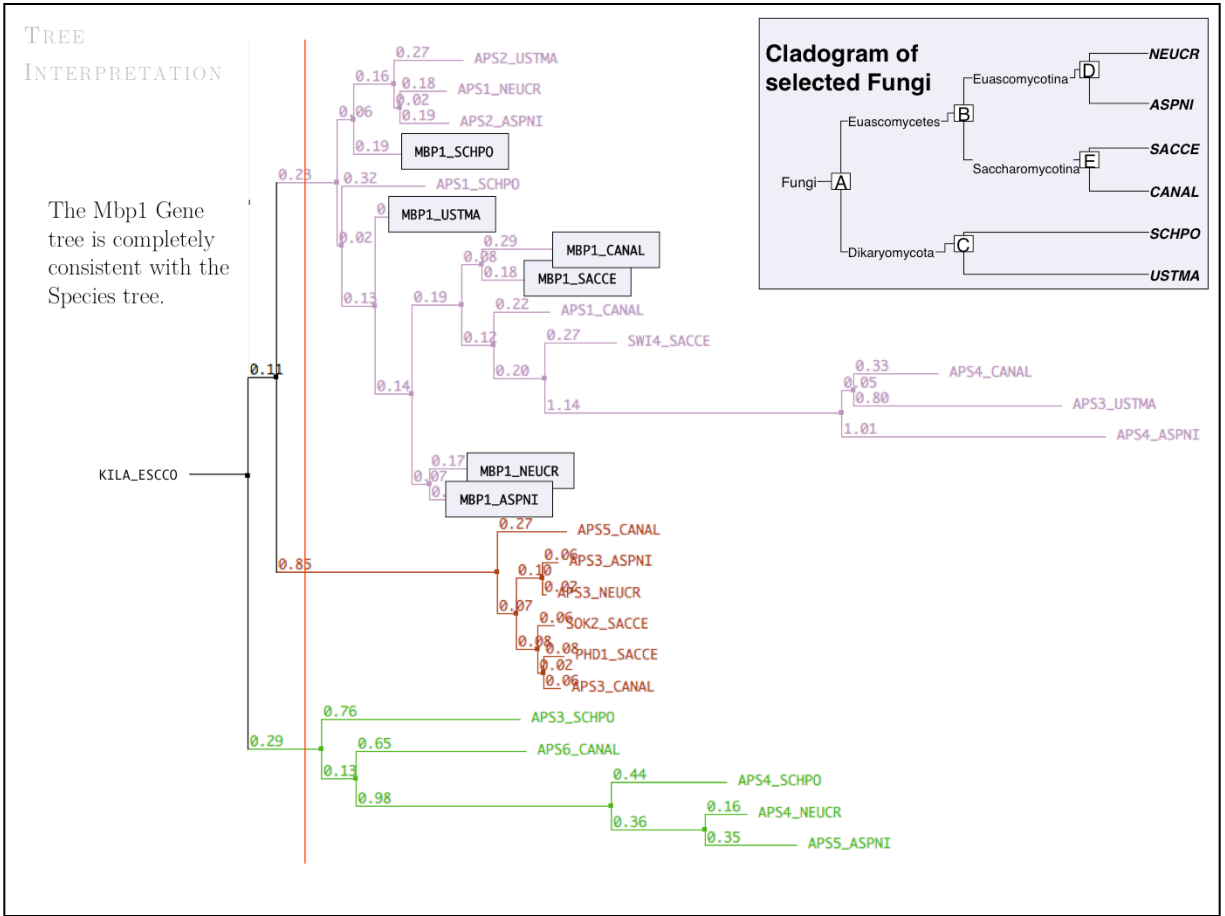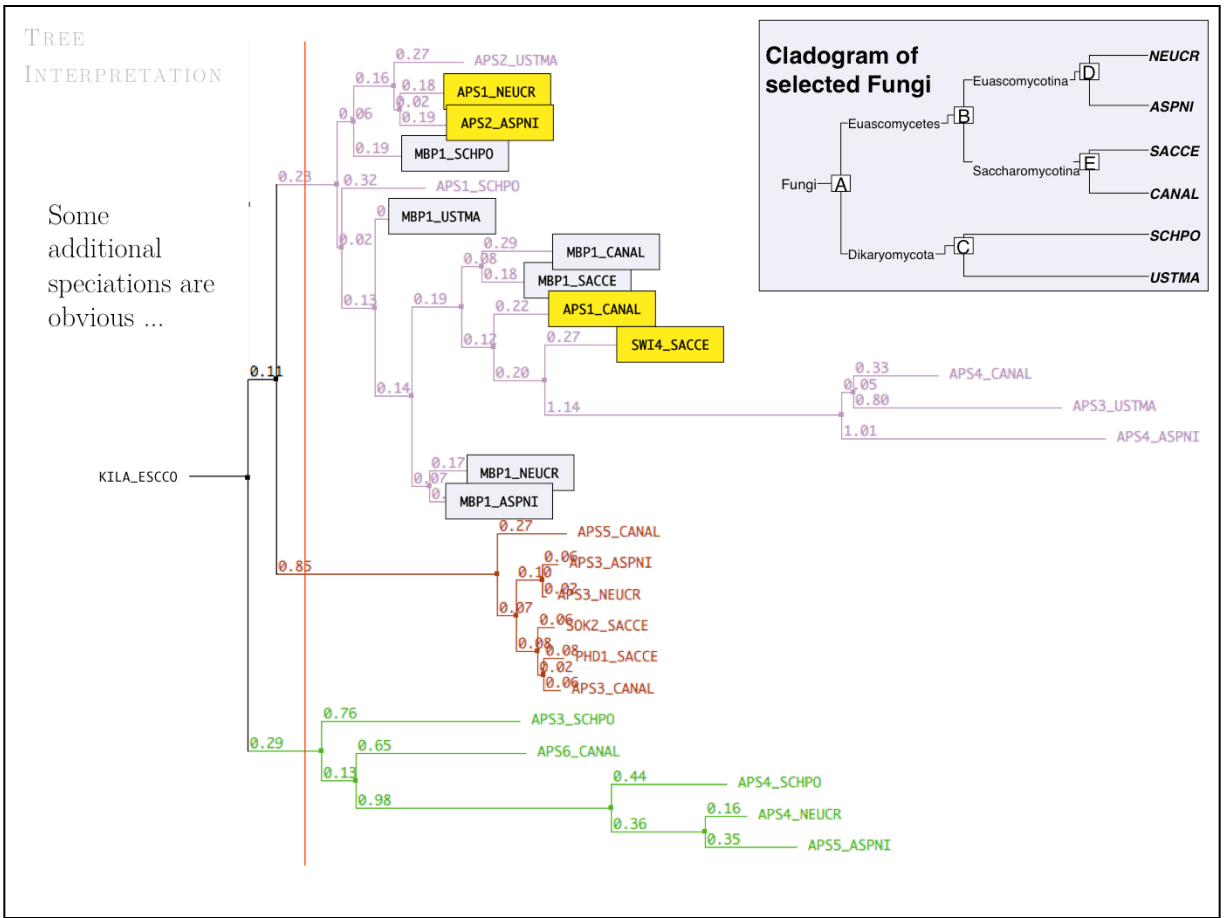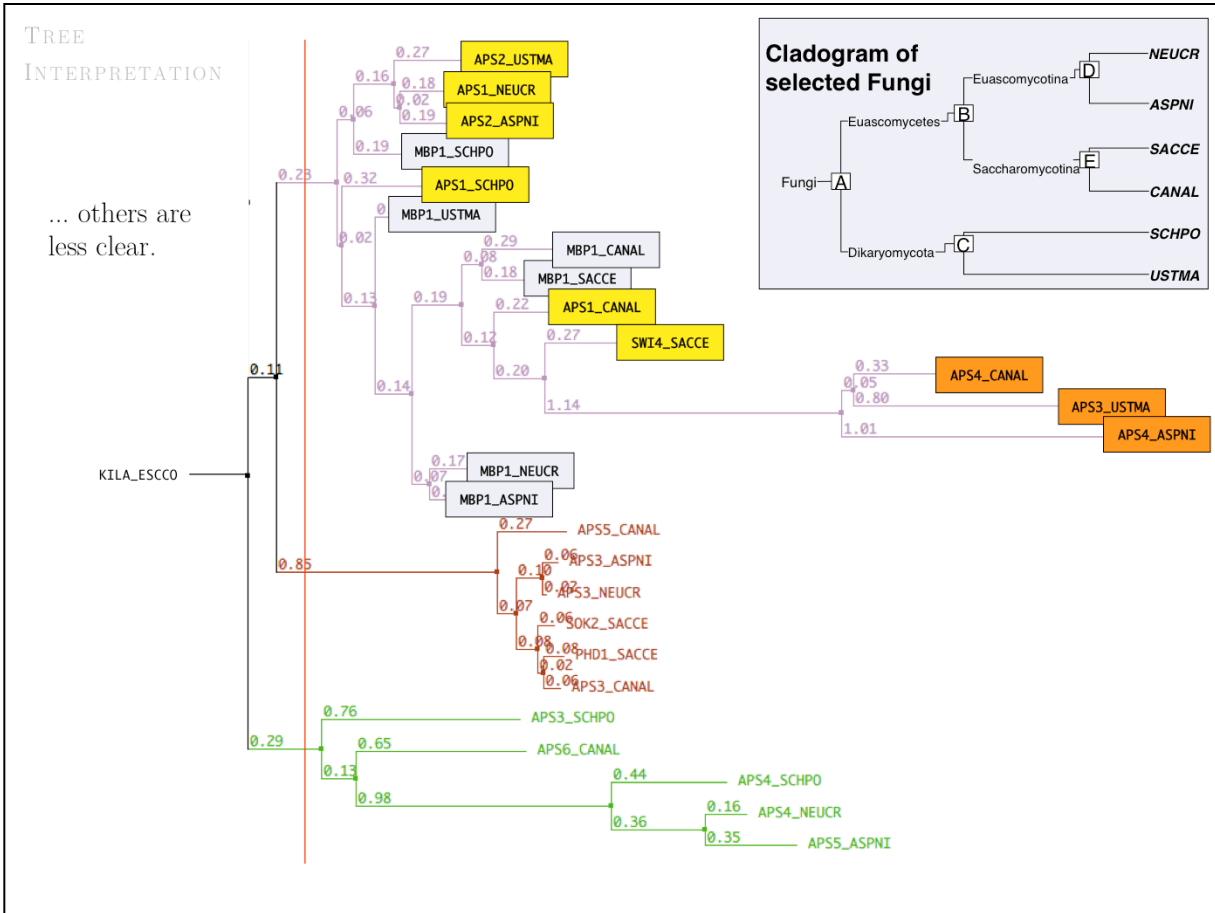
A worked example from course data, displayed in the Jalview tree window.

A worked example from course data, displayed in the Jalview tree window.

A worked example from course data, displayed in the Jalview tree window.
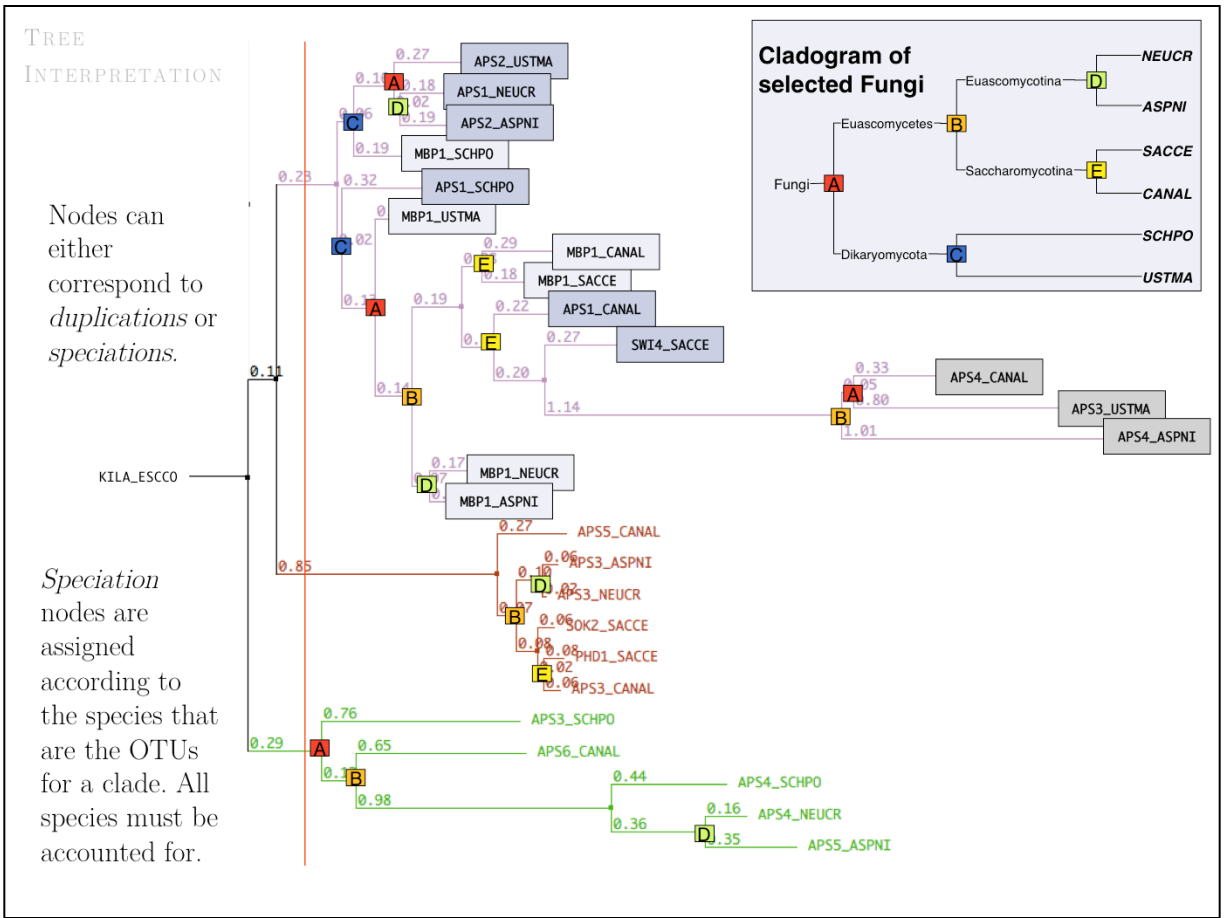
A worked example from course data, displayed in the Jalview tree window.

A worked example from course data, displayed in the Jalview tree window.

A worked example from course data, displayed in the Jalview tree window.
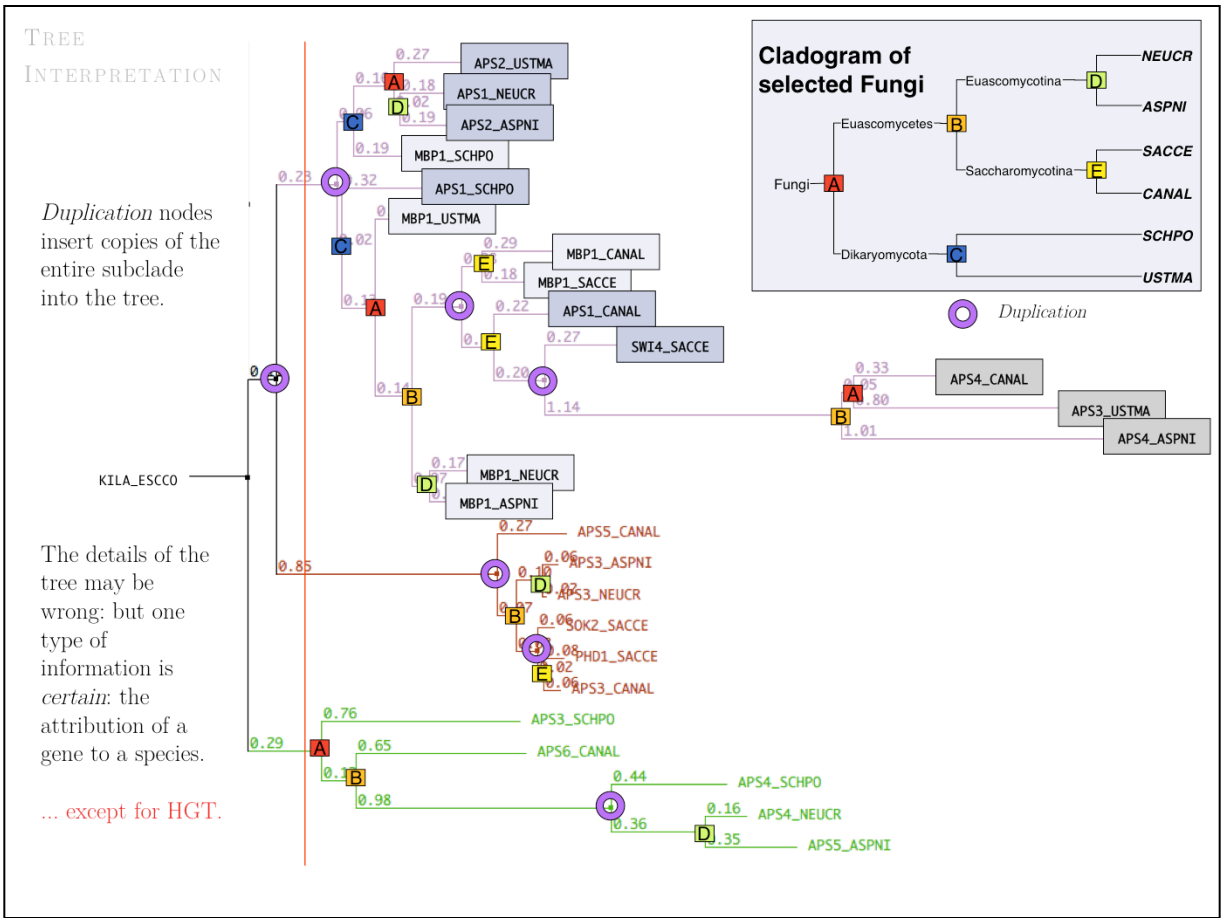
A worked example from course data, displayed in the Jalview tree window.

A worked example from course data, displayed in the Jalview tree window.

A worked example from course data, displayed in the Jalview tree window.

http://steipe.biochemistry.utoronto.ca/abc

B O R I S . S T E I P E @ U T O R O N T O . C A

DEPARTMENT OF BIOCHEMISTRY  &  DEPARTMENT OF MOLECULAR GENETICS
UNIVERSITY OF TORONTO,  CANADA