# Tree pattern matching
## Applied to phylogenetic trees

Jean-François Dufayard
Laurent Duret
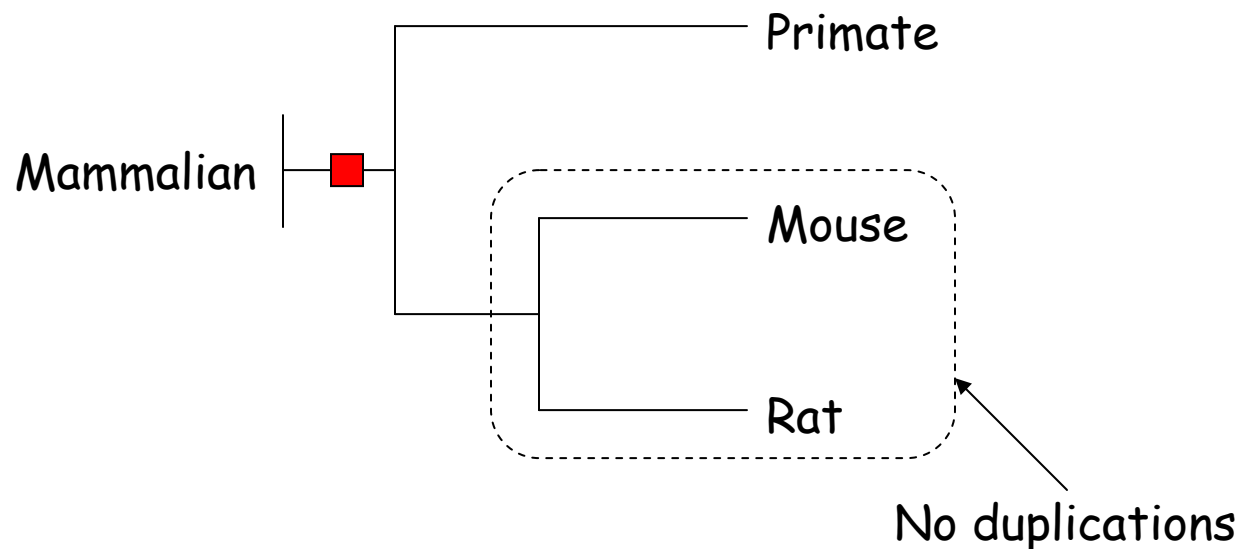Guy Perrière
François Rechenmann

INRIA
RHÔNE-ALPES

• Comparative sequence analysis is a powerful approach to understand genome evolution and is widely used to predict the function of genes.

• This approach requires a phylogenetic analysis to distinguish orthologous and paralogous genes.

• To simplify such phylogenomic analyses, we have developped two databases of homologous genes: HOVERGEN (vertebrates), HOBACGEN (bacteria and archea):

> Genes are classified into families (BLASTP).

> Multiple alignments and phylogenetic trees are computed for each family.

> Taxonomic data from NCBI.

> Protein sequences from SwissProt and TrEMBL.

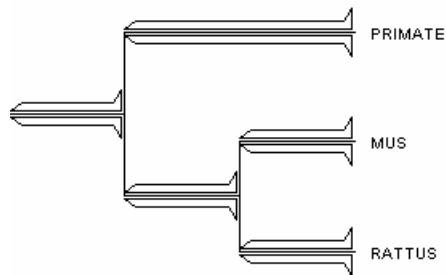> DNA sequences from EMBL.

## Example

Find 3 genes:
- ➢ 1 from primate
- ➢ 1 from mouse
- ➢ 1 from rat

• Primate gene must be paralogous to others
• Mouse gene and rat gene must be strictly orthologous
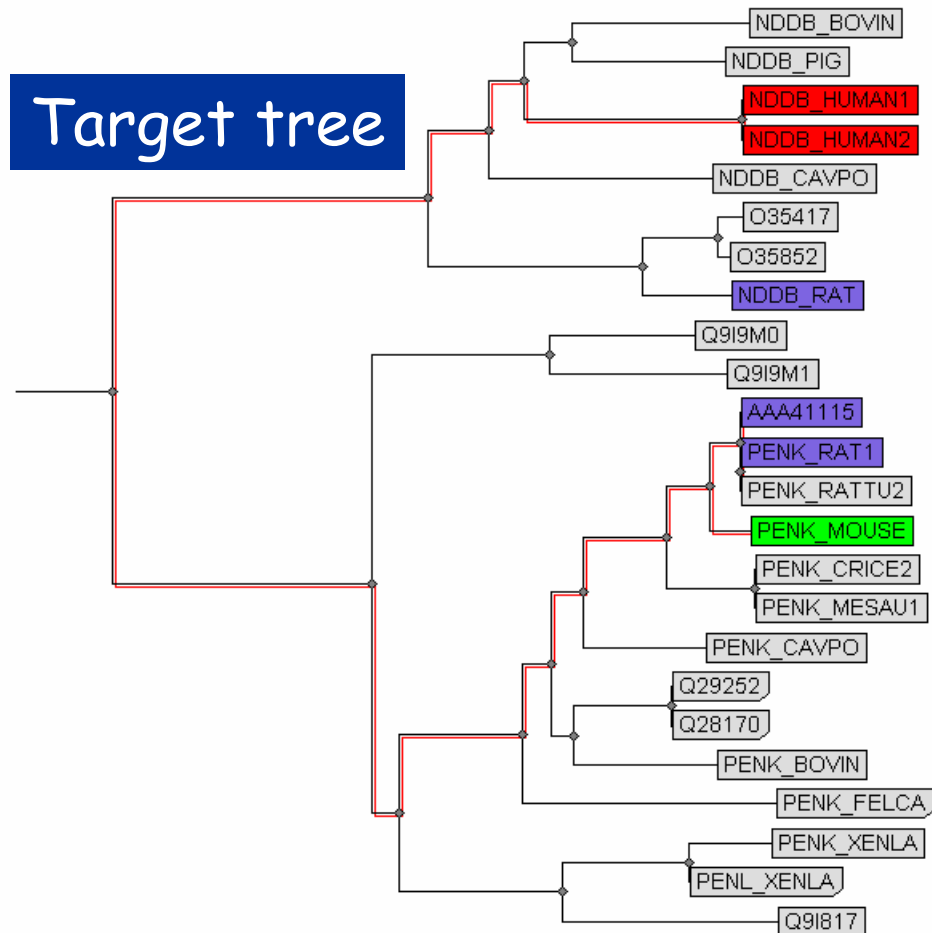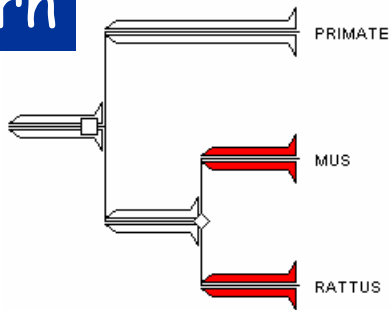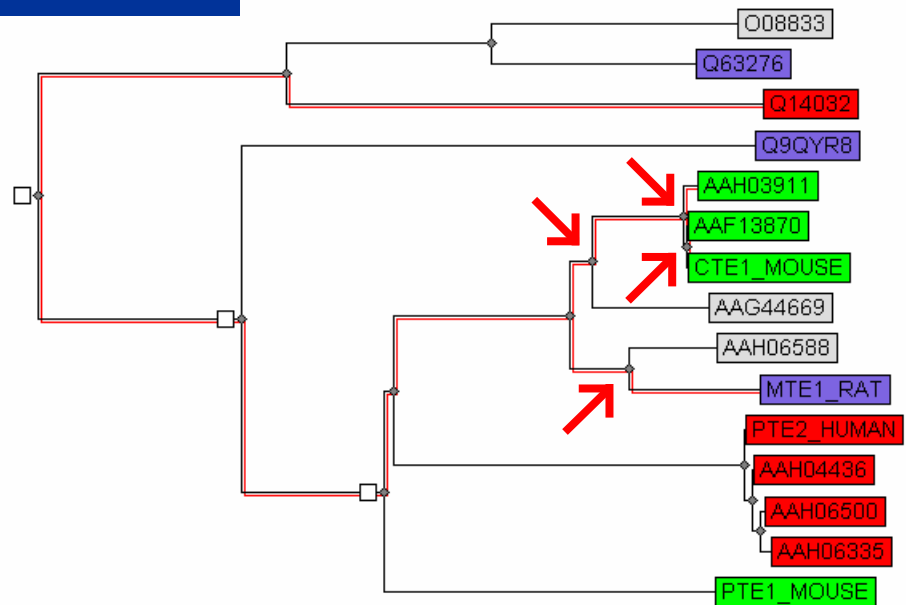• Gene duplication must be posterior to mammalian divergence

# Speciation and duplication nodes

# Constraints on subtrees



Pattern

Target tree

**Find each occurrence of the pattern in the database**

Find each occurrence of the pattern in the database

**Find each occurrence of the pattern in the database**

**Find each occurrence of the pattern in the database**

# CONCLUSION

• This system allows to request HOVERGEN and HOBACGEN using phylogenetic criteria

• 9800 trees in HOVERGEN, or 11500 tree in HOBACGEN can be explored with FAMFETCH, which remains an interactive application.

• Available in june 2002 as a new release of FAMFETCH:
     http://pbil.univ-lyon1.fr/hobacgen/client.html