# Today's Lecture

- PhastCons

# PhastCons PhyloHMM
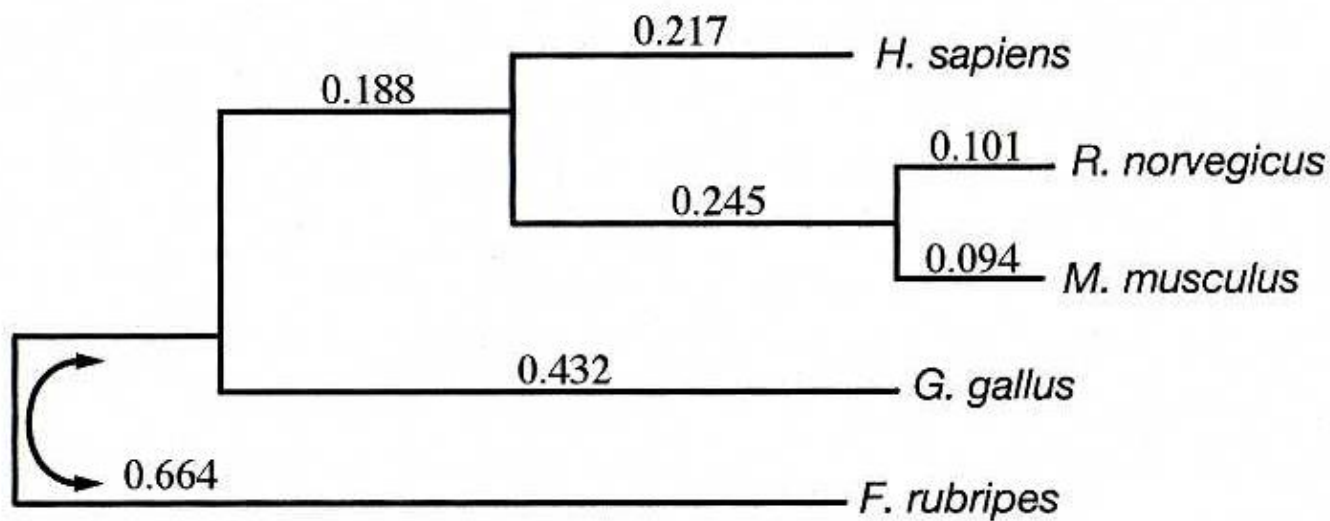


$\mu = a_{cn}$
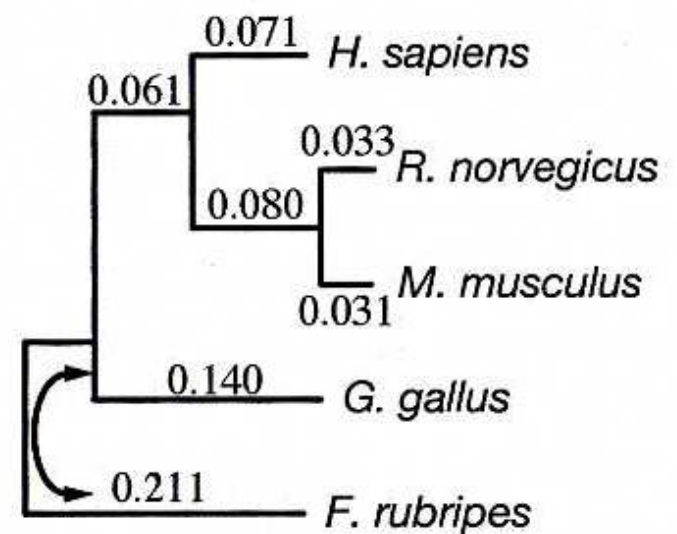
$\nu = a_{nc}$

from Siepel A. *et al.* (2005). Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* 15:1034-50.

**Nonconserved**

0.217 — H. sapiens

0.188

0.101 — R. norvegicus

0.245

0.094 — M. musculus

0.432 — G. gallus

0.664 — F. rubripes

**Conserved**

0.071 — H. sapiens

0.061

0.033 R. norvegicus

0.080

0.031 M. musculus

0.140 — G. gallus

0.211 — F. rubripes
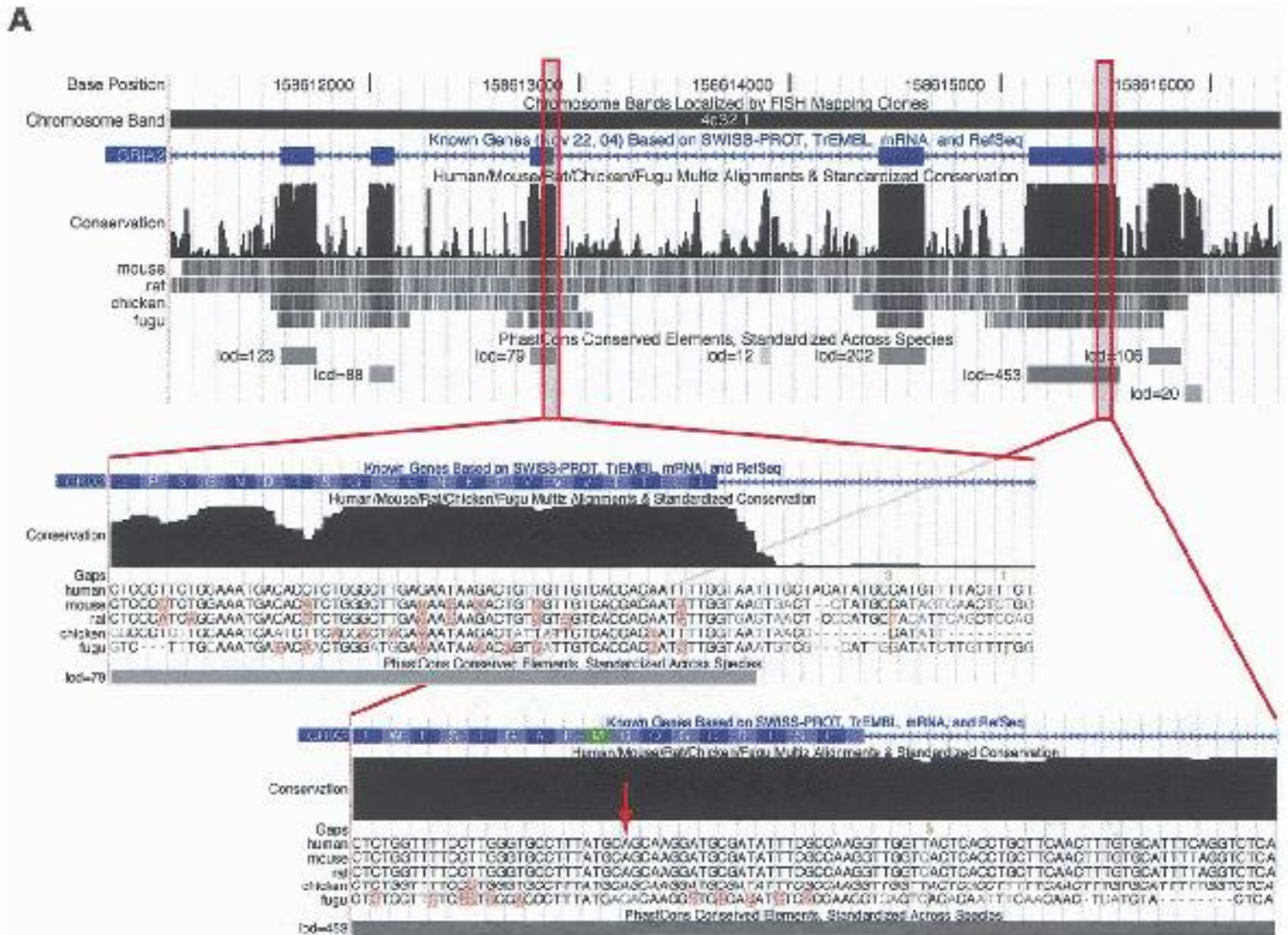
- branch lengths:
  - Expected # substitutions/site over corresponding evolutionary time period
  - for neutral state, should reflect underlying mutation rate
  - for conserved state: mutation rate $\times$ scaling factor $\rho$
    - $\rho$ = frac of mutations that escape purifying selection
    - $\rho \approx .33$ (for vertebrates)

from Siepel A. *et al.* (2005). Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* 15:1034-50.

# Some general issues in applying probability models, in the PhyloHMM context

- Is the model computable?

- Is the model 'reasonable'?
  - 2 states enough?
  - Markov condition on transition probabilities

- How good is the input data?
  - Alignability of neutral sequence
  - Accuracy of genome sequence alignments

- Are results reliable?
  - No true 'test set' – instead, putative false positive rate, and 'biological plausibility' of findings

# Alignment issues

- Multiz: progressive pairwise alignments
- accurate multiple genome alignment *not* a solved problem!
  - statistical assessment: Prakash & Tompa (2005, 2007, 2009)
  - ENCODE region alignment analyses: Margulies EH *et al.* 2007
  - major issues:
    - accurate gap placement (even for close species!!)
    - discrimination among paralogous sequences (e.g. repeats, duplications)
- inaccurate alignments cause
  - neutral rate to be *overestimated*
  - conserved segments to be *overidentified*
    - because more slowly mutating (or better aligned) neutral segments may be called conserved