# Multiple Alignments and Phylogenies

Mark Voorhies

3/31/2011

- Computer "lab notebook"
  - Dates
  - Input/output/parameters

- Computer "lab notebook"
  - Dates
  - Input/output/parameters (screenshots)

- Computer "lab notebook"
    - Dates
    - Input/output/parameters (screenshots)
    - E-mail result synopses (*e.g.*, to PI)

- Computer "lab notebook"
  - Dates
  - Input/output/parameters (screenshots)
  - E-mail result synopses (*e.g.*, to PI)
- BLAST
  - Fast, heuristic local alignment with statistical foundation

# Review

- Computer "lab notebook"
  - Dates
  - Input/output/parameters (screenshots)
  - E-mail result synopses (*e.g.*, to PI)
- BLAST
  - Fast, heuristic local alignment with statistical foundation
  - Gosh, we should have saved our BLAST reports!

- We have a bunch of sequences that look similar to our query

- We have a bunch of sequences that look similar to our query
- We infer that they are homologous to each other

- We have a bunch of sequences that look similar to our query
- We infer that they are homologous to each other
- What does that mean, anyway?

Homologs heritable elements with a common evolutionary origin.

Homologs heritable elements with a common evolutionary origin.

Orthologs homologs arising from speciation.

Paralogs homologs arising from duplication and divergence within a single genome.

# Nomenclature

Homologs heritable elements with a common evolutionary origin.

Orthologs homologs arising from speciation.

Paralogs homologs arising from duplication and divergence within a single genome.

Xenologs homologs arising from horizontal transfer.

Onologs homologs arising from whole genome duplication.

Graph

Node

Cycle

Edge

Connected
Subgraph

Tree = Connected Graph with no Cycles

Distances
(Pairwise relationships)

Topology
(Evolutionary history)

# Measure all pairwise distances by dynamic programming



$$\frac{(N\text{-}1)^2}{2}$$

# Measure all pairwise distances by dynamic programming

- Multiple Alignment
    - T-Coffee
    - MUSCLE
    - COBALT
- Tree building
    - MrBayes (Bayesian MCMC)
    - PhyML (maximum likelihood)

- Play with CLUSTALX, JALVIEW, and PSI-BLAST
- Read PLoS Comp. Biol. 4:e1000069