

Comparative Genomics

- Multiple alignments
- Synteny
- Homologs
- Gene models

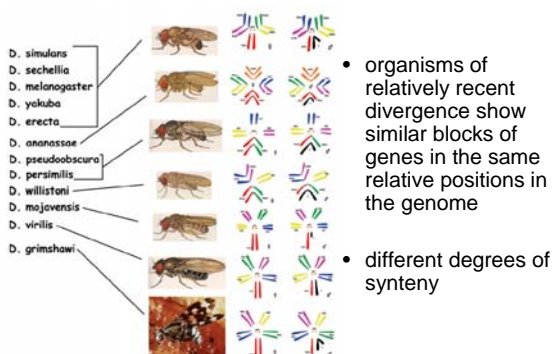
Evolution:

- Phylogeny
- Gene expression

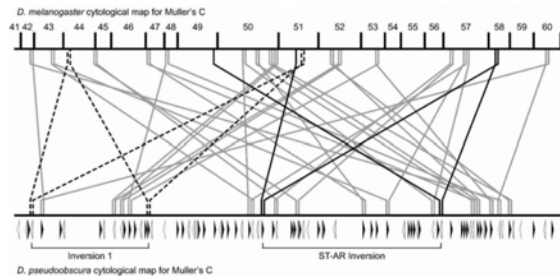
Multiple sequence alignment

- Heuristic vs. global optimisation
- DP – v. v. slow
- Progressive alignment construction – e.g. Clustal family
- Iterative methods – e.g. MUSCLE
- Consensus methods
- HMMs e.g. HMMer
- Motif finding e.g. MEME – see Regulation lectures
- Not practical on large scale

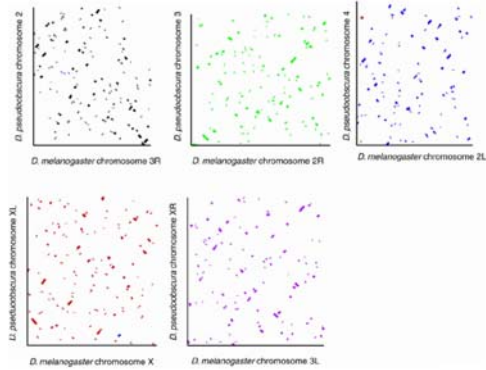
Synteny



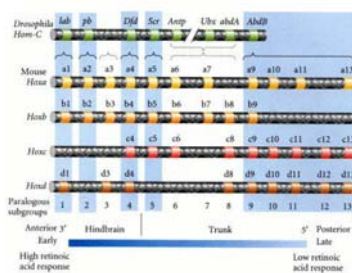
Synteny plot



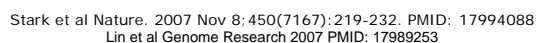
Synteny – dot plots

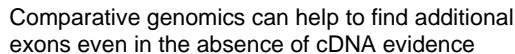


Synteny – Hox cluster



The most prominent syntenic group may be the Hox cluster which is even conserved between flies and mammals.





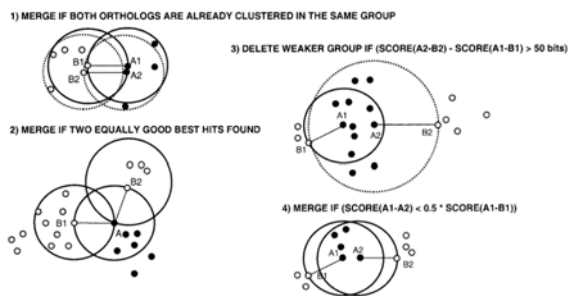
	Total	Confirmed	Unclear	Rejected ^a
Named genes	4711	4566 (96.9%)	105 (2.2%)	40 (0.8%)
Well-studied genes	893	882 (98.8%)	8 (0.9%)	3 (0.3%)
Other named genes	3,818	3,684 (96.5%)	97 (2.5%)	37 (1.0%)
CDS-only genes	9322	7879 (84.5%)	229 (2.4%)	1,414 (15.1%)
CDS-annotated	4373	3877 (88.7%)	728 (16.6%)	198 (4.5%)
Unclassified	12,233	10,445 (85.4%)	451 (3.7%)	2,137 (17.5%)
All genes	12,233	10,445 (85.4%)	451 (3.7%)	4,544 (37.1%)
Noncoding regions	15,564	3 (0.0%)	13,430 (86.3%)	2,131 (13.7%)



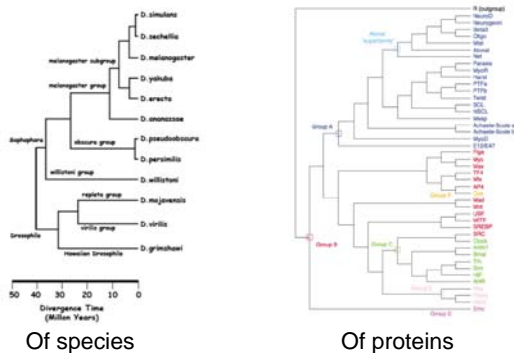
Finding homology

- BLAST bi-directional best hit of protein sequences and minimal sequence identity of 30%
- Protein family clusters (looser criterion, includes paralogy within the family), implemented in
 - ENSEMBL Compara
 - Inparanoid
 - OrthoMCL / MCL

Inparanoid

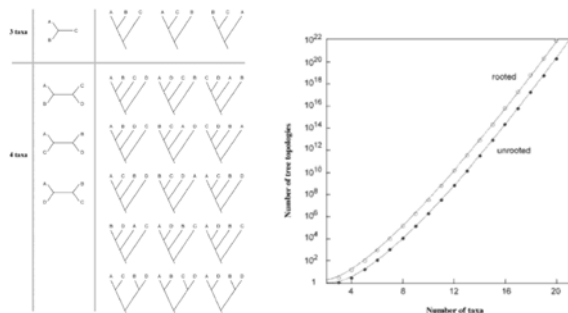


Phylogeny



- Phylogenetic trees are usually based on DNA or protein sequences.
- Comparisons are possible even between animals with no physical resemblance.
- Ideally, a phylogenetic reconstruction is unambiguous...

Graph theory shows how difficult tree growing is



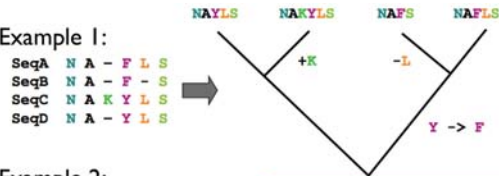
General principle



1. Make a multiple sequence alignment.
 2. Determine the distance between the sequences.
 3. Use these differences to infer the phylogenetic relationship.
- closely related species: use DNA.
 - in most other cases: use protein

Which hypothesis is true?

Example 1:



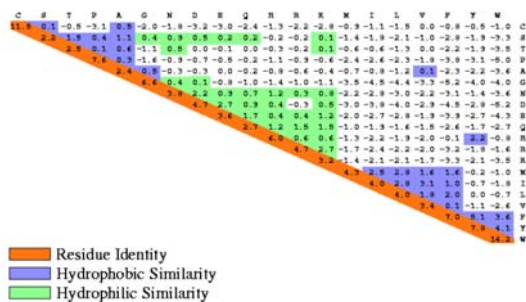
Example 2:

seqA TCAGACGATTG (11)
seqB TCGGAGCTG (9)

In both cases, you're making an assumption about how evolution has taken place!

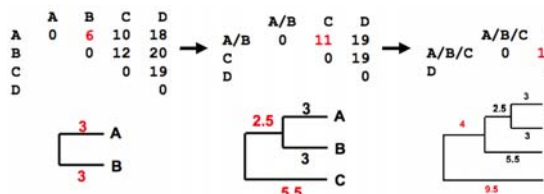
- I. TCAG-ACG-ATTG
TC-GGA-GC-T-G perfectly matching alignment
- II. TCAGACGATTG
TCGGAGCTG-- no internal gaps
- III. TCAG-ACGATTG
TC-GGA--GCTG a mixture of both

Gonnet Pam250 Matrix



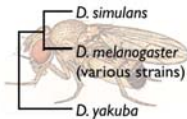
Tree building

- Huge range of different methods around, here only the most simple one: UPGMA.
- Unweighted pair group method using arithmetic means => this is similar to average distance clustering of gene expression data.



Gene expression

Half of the genes show developmental dynamics

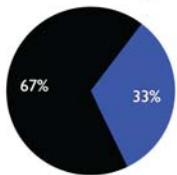


A phylogenetic tree showing the relationships between three species of Drosophila: *D. simulans*, *D. melanogaster* (various strains), and *D. yakuba*. The tree is rooted on the left, with *D. simulans* at the top, *D. melanogaster* in the middle, and *D. yakuba* at the bottom. The branches are labeled with the species names.

● conserved
● changes

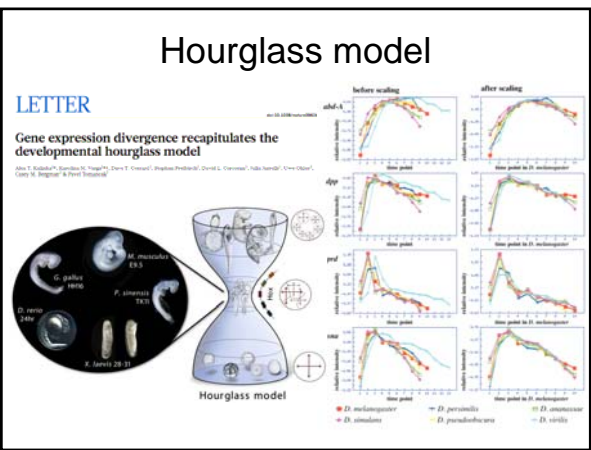
gene expression between species

Based on microarray data:



A pie chart showing the distribution of gene expression between species. The chart is divided into two segments: a black segment representing 67% (conserved) and a blue segment representing 33% (changes). The segments are labeled with their respective percentages.

Adapted from:
Rifkin et al., Nature Genetics 33: 138-144 (2003)



References

- Multiple sequence alignment - <http://www.ebi.ac.uk/Tools/msa/>
Durbin *et al.* - Biological Sequence Analysis -Cambridge University Press
- Inparanoid
 - Remm *et al.* - Automatic clustering of orthologs and in-paralogs from pairwise species comparisons - J. Mol. Biol. 2001
<http://inparanoid.sbc.su.se/cgi-bin/index.cgi>
- OrthoMCL / MCL
 - Chen *et al.* - OrthoMCL-DB: querying a comprehensive multi-species collection of ortholog groups - Nucl Acids Res. 2006
<http://orthomcl.org/orthomcl/>
 - Enright *et al.* - An efficient algorithm for large-scale detection of protein families - Nucl Ac Res. 2002
<http://micans.org/mcl/>
- Kalinka *et al.* - Gene expression divergence recapitulates the developmental hourglass model – Nature 2010

Networks References

- Buchanan et al (eds.) - Networks in Cell Biology - Cambridge University Press
- Barabási and Oltvai - Network biology: understanding the cell's functional organization - Nature Reviews Genetics 2004
- Boone *et al.* - Exploring genetic interactions and networks with yeast. - Nat Rev Genet. 2007
- Kelley and Ideker - Systematic interpretation of genetic interactions using protein networks. - Nat Biotechnol. 2005
- Markowitz and Spang - Inferring cellular networks – a review - BMC Bioinformatics 2007
- Schuster *et al.* - A general definition of metabolic pathways useful for systematic organization and analysis of complex metabolic networks - Nat Biotechnol. 2000
