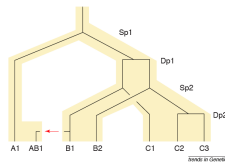


Exercise

- On the basis of Fitch's definitions (previous slide), qualify the relationships between each pair of genes in the illustrative schema.
 - P paralog
 - O ortholog
 - X xenolog
 - A analog

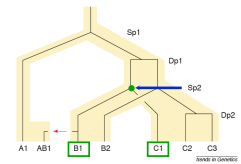


	A1	AB1	B1	B2	C1	C2	C3
A1							
AB1							
B1							
B2							
C1							
C2							
C3							

- Orthologs** can formally be defined as a pair of genes whose last common ancestor occurred immediately before a speciation event (ex: a₁ and a₂).
- Paralogs** can formally be defined as a pair of genes whose last common ancestor occurred immediately before a gene duplication event (ex: b₂ and b₁).
Source: Zvelebil & Baum, 2000

Exercise

- Example: B1 versus C1**
 - The two sequences (B1 and C1) were obtained from taxa B and C, respectively.
 - The cenancestor (blue arrow) is the taxon that preceded the second speciation event (Sp2).
 - The common ancestor gene (green dot) coincides with the cenancestor.
- > B1 and C1 are orthologs

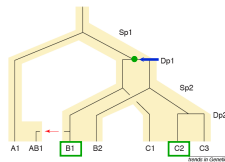


	A1	AB1	B1	B2	C1	C2	C3
A1							
AB1							
B1							
B2							
C1							
C2							
C3							

- Orthologs** can formally be defined as a pair of genes whose last common ancestor occurred immediately before a speciation event.
- Paralogs** can formally be defined as a pair of genes whose last common ancestor occurred immediately before a gene duplication event.
Source: Zvelebil & Baum, 2000

Exercise

- Example: B1 versus C2**
 - The two sequences (B1 and C2) were obtained from taxa B and C, respectively.
 - The common ancestor gene (green dot) is the gene that just preceded the duplication Dp1.
 - This common ancestor is much anterior to the cenancestor (blue arrow).
- > B1 and C2 are paralogs

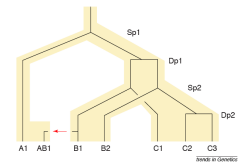


	A1	AB1	B1	B2	C1	C2	C3
A1							
AB1							
B1							
B2							
C1							
C2							
C3							

- Orthologs** can formally be defined as a pair of genes whose last common ancestor occurred immediately before a speciation event.
- Paralogs** can formally be defined as a pair of genes whose last common ancestor occurred immediately before a gene duplication event.
Source: Zvelebil & Baum, 2000

Solution to the exercise

- On the basis of Fitch's definitions (previous slide), qualify the relationships between each pair of genes in the illustrative schema.
 - P paralog
 - O ortholog
 - X xenolog
 - A analog



The shaded relation of a gene (tree) is shown from a common ancestor to an ancestral population (the grey background), descending to three populations labeled A, B and C. There are two speciation events (Sp1 and Sp2), each occurring at the junctions shown as an upside down Y. There are also two gene duplication events (Dp1 and Dp2), depicted by a horizontal bar. The genes whose common ancestor resides at a horizontal bar (gene duplication) are paralogous. Thus, C2 and C1 are paralogous to each other but are orthologous to B2. Both are paralogous to B1 but orthologous to A1. The red arrow denotes the transfer of the B1 gene from species B to species A. As a result, the B1 gene is orthologous to all other genes. All three subtop relationships are reflexive, that is, A1 => B1 implies B1 => A1, where => stands for "is orthologous to". However, the relationships are not transitive. Thus, C2 => A1 => C1 might be true, but it is not necessarily true that C2 => C1, as indeed it is not in the figure if => is read as "is orthologous to". A different non-transitivity occurs for paralogs: B2 => B1 => C1 => C2.

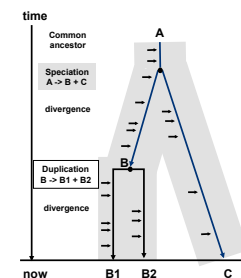
	A1	AB1	B1	B2	C1	C2	C3
A1	I						
AB1	X	I					
B1	O	X	I				
B2	O	X	P	I			
C1	O	X	O	P	I		
C2	O	X	P	O	P	I	
C3	O	X	P	O	P	P	I

How to detect orthology relationships ?

- Bidirectional best hits**
 - Concepts
 - Best hit (BH)
 - Reciprocal (RBH) or bidirectional (BBH) best hit.
 - Problem 1: non-reciprocity of the BH relationship, which may result from various effects
 - Multidomain proteins -> non-transitivity of the homology relationship
 - Detection: no paralogy
 - Paralogs in one genome correspond to the same ortholog in the other genome
 - Non-symmetry of the BLAST result
 - Can be circumvented by using dynamical programming (Smith-Waterman)
 - Problem 2: unequivocal but fake reciprocal best hit
 - Duplication followed by a deletion
 - Two paralogs can be BBH, but the true orthologs are not present anymore in the genome (due to duplication).
 - Ex: Hox genes
 - Conceptual problem: intrinsically unable to treat multi-orthology relationships
 - Ex: Fitch figure: B2 is ortholog to both C2 and C3, but only one of these will be its Best Hit.
- Conclusion:** the analysis of BBH is intrinsically unable to reveal the true orthology relationships

Conclusions - Orthology versus paralogy

- The shaded tree represents the history of the species, the thin black tree the history of the sequences.
- Homology can be subdivided in two subtypes, depending on the species/sequence history:
 - Orthologs** are sequences whose last common ancestor occurred immediately before a speciation event.
 - Paralogs** are sequences whose last common ancestor occurred immediately before a duplication event. (Fitch, 1970; Zvelebil & Baum, 2000)
- Example:
 - B and C are orthologs, because their last common ancestor lies just before the speciation A -> B + C
 - B1 and B2 are paralogs because the first event that follows their last common ancestor (B) is the duplication B -> B1 + B2
- Beware! These definitions are often misunderstood, even in some textbooks.
- Contrarily to a strong belief, orthology can be a 1 to N relationship.
 - B1 and C are orthologs, because the first event after their last common ancestor (A) was the speciation A -> B + C
 - B2 and C are orthologs because the first event after their last common ancestor (A) was the speciation A -> B + C
- The strategy to search reciprocal best hits (RBH) is thus a simplification that misses many true orthologs (it is essentially justified by pragmatic reasons).
- The orthology relationship is reciprocal but **not transitive**.
 - C <-[orthologous]-> B1
 - C <-[orthologous]-> B2
 - B1 <-[paralogous]-> B2
- The commonly used concept of **clusters of orthologs is thus an aberration**.



How to circumvent the weaknesses of RBH ?

- Solutions to the problems with RBH
 - Domain analysis: analyze the location of the hits in the alignments
 - Resolves the problems of gene fusion (two different fragments of a protein in genome A correspond to 2 distinct proteins of genome B)
 - Analysis of the evolutionary history : full phylogenetic inference + reconciliation of the sequence tree and the species tree
 - Resolves the cases of multiple orthology relationships (n to n)
 - Does not resolve the problems of differential deletions after regional duplications
 - Solving the problem of regional duplications followed by differential deletion
 - Analysis of synteny: neighbourhood relationships between genes across genomes
 - Analysis of pseudo-genes: allows to infer the presence of a putative gene in the common ancestor
 - This is OK when the duplication affects a regions sufficiently large to encompass multiple genes.
- These solutions require a case-by-case analysis -> this is not what you will find in the large-scale databases.
- Resources:
 - Ensembl database
 - SPRING database

Criteria for genome-wise detection of orthologs

- Criterion for detecting paralogy
 - Two genes from a given species (e.g. C) are more similar to each other than to their best hit in genome B.
- Pairs of orthologous genes
 - BeT (Best-scoring BLAST hit)
 - Insufficient to infer orthology
 - Bidirectional best hit (BBH)
 - Better approximation
 - Discuss the problem of gene loss
- Clusters of orthologous genes (COGs)
 - Triangular definition of COGs (Tatusov, 1997)
 - KOG: eukaryotic Orthologous Groups
 - Question: is there any interest of defining a new term for eukaryotes ?
- To discuss
 - theoretical weakness of the COG concept: since orthology is NOT a transitive relationship.
 - Pragmatic value of the concept

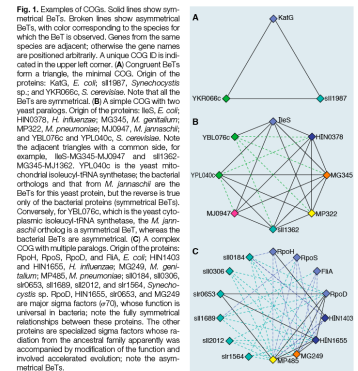


Figure from Tatusov, 1997

Example of pairwise alignment

- Example of alignment


```

TTTGCCTT--AAATCGTGTAGCAATTT  s=substitution
s|ss||||ggs|ggggg|||||s|      g=gap
ATGCCGTTTTTAA-----TAGCAATAT  |=identical residues
            
```
- Gaps, insertions and deletions
 - Gaps can reflect either an insertion in one of the sequences, or a deletion in the other one.
 - The simple observation of two aligned sequences is insufficient to decide whether a gap results from an insertion or a deletion.
 - The term *indel* is sometimes used in this case to designate the evolutionary event.

Global versus local alignment

- Global alignment
 - (e.g. Proteins having a common ancestor and being conserved over their whole sequences)


```

LQGPSKGTGKGS-SRSWDN
|---|---|---|---|
LN-ITKSAGKGAIMRLGDA
                    
```
- Local alignment
 - Example: proteins sharing a common domain


```

LQGPSKGTGKGS-SSRIWDN
|---|
LN-ITKAGKGAIMRLGDA
                    
```

Adapted from Didier Gonze

Some definitions

- Identity
 - The level of identity is a simple calculation of fraction of residues which are identical between the two aligned sequences.
- Similarity
 - Two residues are considered similar if their substitution does not affect the function of the protein.
 - The level of identity is a simple calculation of fraction of residues which are similar between the two aligned sequences.
 - We will see below the criteria to consider that two residues are similar.
- Homology
 - Homology indicates the fact that two sequences diverged from a common ancestor.

Suggested readings

- Homology
 - Tatusov, R. L., Koonin, E. V. & Lipman, D. J. (1997). A genomic perspective on protein families. *Science* 278, 631-7.
 - Fitch, W. M. (2000). Homology a personal view on some of the problems. *Trends Genet* 16, 227-31.
 - Koonin, E. V. (2005). Orthologs, paralogs, and evolutionary genomics. *Annu Rev Genet* 39, 309-38.
 - Zvelebil, M. J. and Baum, J. O. (2008). *Understanding Bioinformatics*. Garland Science: New York and London.
- Phylogenetic profiles
 - Tatusov, R. L., Koonin, E. V. & Lipman, D. J. (1997). A genomic perspective on protein families. *Science* 278, 631-7.
 - Pellegrini, M., Marcotte, E. M., Thompson, M. J., Eisenberg, D. & Yeates, T. O. (1999). Assigning protein functions by comparative genome analysis: protein phylogenetic profiles. *Proc Natl Acad Sci U S A* 96, 4285-8.
- Gene fusion
 - Marcotte, E. M., Pellegrini, M., Ng, H. L., Rice, D. W., Yeates, T. O. & Eisenberg, D. (1999). Detecting protein function and protein-protein interactions from genome sequences. *Science* 285, 751-3.