

Spring 2024  
**BIOS 477/877**  
*Bioinformatics and Molecular Evolution*  
**Lecture 5**

BIOS477/877 L5 - 1

1

**TODAY'S TOPICS**

➤ **Molecular Evolution - part 3**

BIOS477/877 L5 - 2

2

**Identifying Selection**

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$d_N(AB) < d_N(C)$

$d_N(AB) = u_T(AB) * f_0(AB)$   
 $d_N(C) = u_T(C) * f_0(C)$

Both mutation rates ( $u_T$ ) and selective constraints ( $f_0$ ) affect nonsynonymous rates ( $d_N$ )

Note: In literatures,  $K_A$  and  $K_S$  are also used instead of  $d_N$  and  $d_S$

(Based on nonsynonymous substitution rates estimated from human-rat comparison) BIOS477/877 L5 - 3

3

**Identifying Selection**

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$d_N(AB) < d_N(C)$

$d_N(AB) = u_T(AB) * f_0(AB)$   
 $d_N(C) = u_T(C) * f_0(C)$

Does it mean:  
 $u_T(AB) < u_T(C)$  ?  
 $f_0(AB) < f_0(C)$  ?

Both mutation rate ( $u_T$ ) and selective constraints ( $f_0$ ) can affect  $d_N(AB) < d_N(C)$ .

(Based on nonsynonymous substitution rates estimated from human-rat comparison) BIOS477/877 L5 - 4

4

**Identifying Selection**

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$d_N(AB) < d_N(C)$

$u_T(AB) * f_0(AB) < u_T(C) * f_0(C)$

If mutation rates are constant within a gene,  
 $\rightarrow u_T(AB) = u_T(C)$   
 $\rightarrow d_N(AB) < d_N(C)$  can be explained by  
 $f_0(AB) < f_0(C)$

(Based on nonsynonymous substitution rates estimated from human-rat comparison) BIOS477/877 L5 - 5

5

**Identifying Selection**

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$d_N(AB) < d_N(C)$

$u_T(AB) * f_0(AB) < u_T(C) * f_0(C)$

If mutation rates are different  
 $\rightarrow u_T(AB) \neq u_T(C)$   
 $\rightarrow d_N(AB) < d_N(C)$  cannot be simply explained by  
 $f_0(AB) < f_0(C)$

(Based on nonsynonymous substitution rates estimated from human-rat comparison) BIOS477/877 L5 - 6

6

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$d_N(AB) < d_N(C)$   
 $d_N(AB) = u_T(AB) * f_0(AB)$   
 $d_N(C) = u_T(C) * f_0(C)$

Synonymous substitutions are assumed to be neutral or near-neutral.

$d_S = u_T f_0 = u_T$   
where  $f_0 \approx 1$  (neutral)

(Based on nonsynonymous substitution rates estimated from human-rat comparison) BIOS477/877 L5 - 7

7

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$d_N(AB) < d_N(C)$   
 $d_N(AB) = u_T(AB) * f_0(AB)$   
 $d_N(C) = u_T(C) * f_0(C)$

Synonymous substitutions are assumed to be neutral or near-neutral.

$d_S(AB) = u_T(AB)$   
 $d_S(C) = u_T(C)$

(Based on nonsynonymous substitution rates estimated from human-rat comparison) BIOS477/877 L5 - 8

8

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$d_N(AB)$  vs.  $d_N(C)$

$d_N(AB) = u_T(AB) * f_0(AB)$	$d_N(C) = u_T(C) * f_0(C)$
$d_S(AB) = u_T(AB)$	$d_S(C) = u_T(C)$
$\frac{d_N(AB)}{d_S(AB)} = \frac{u_T(AB) * f_0(AB)}{u_T(AB)}$	$\frac{d_N(C)}{d_S(C)} = \frac{u_T(C) * f_0(C)}{u_T(C)}$
$\frac{d_N(AB)}{d_S(AB)} = f_0(AB)$	$\frac{d_N(C)}{d_S(C)} = f_0(C)$

BIOS477/877 L5 - 9

9

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

$\frac{d_N(AB)}{d_S(AB)} = f_0(AB)$	$\frac{d_N(C)}{d_S(C)} = f_0(C)$
-------------------------------------	----------------------------------

$d_N/d_S$  can be used to identify selection

- Shows only the level of selective constraints ( $f_0$ )
- We don't need to worry about mutation rates ( $u_T$ )
- Even if  $u_T(AB) \neq u_T(C)$ ,  $d_N/d_S$  can be compared

BIOS477/877 L5 - 10

10

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$  (or  $K_S$ )
- Number of **nonsynonymous** substitutions per site:  $d_N$  (or  $K_N$ )

→  $d_N/d_S$  can be used to identify selection

~~$d_N(AB) < d_N(C)$~~   
 $d_N(AB)/d_S(AB) ? d_N(C)/d_S(C)$

This is the better comparison

Regardless of the mutation rates, we can compare:  $f_0(AB)$  vs.  $f_0(C)$

BIOS477/877 L5 - 11

11

### Identifying Selection

$d_{N1} < d_{N2}, d_{S1} < d_{S2}$   
 This may be caused simply by the difference in divergence time ( $t_1 < t_2$ )

$d_{N1} = u_{T1} * f_1 * 2t_1, d_{S1} = u_{T1} * 2t_1$   
 $d_{N2} = u_{T2} * f_2 * 2t_2, d_{S2} = u_{T2} * 2t_2$

BIOS477/877 L5 - 12

12

### Identifying Selection

Human: Gene A }  $d_{N1}, d_{S1}$   
 Chimp: Gene A }  
 Rat: Gene A }  $d_{N2}, d_{S2}$   
 Mouse: Gene A }

$t_1 < t_2$   
 ( $t_1, t_2$ : divergence time)

$d_{N1} < d_{N2}, d_{S1} < d_{S2}$   
 This may be caused simply by the difference in divergence time ( $t_1 < t_2$ )

$d_{N1}/d_{S1}$  vs.  $d_{N2}/d_{S2}$   
 By comparing  $d_N/d_S$ , the time effect can be cancelled out. We can compare the selective constraints between genes and between lineages.

$d_{N1} = u_{T1} * f_1 * 2t_1, d_{S1} = u_{T1} * 2t_1 \rightarrow d_{N1}/d_{S1} = f_1 * 2t_1 / (u_{T1} * 2t_1) = f_1$   
 $d_{N2} = u_{T2} * f_2 * 2t_2, d_{S2} = u_{T2} * 2t_2 \rightarrow d_{N2}/d_{S2} = f_2 * 2t_2 / (u_{T2} * 2t_2) = f_2$

BIOS477/877 L5 - 13

13

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$
- Number of **nonsynonymous** substitutions per site:  $d_N$

→  $d_N/d_S$  can be used to identify selection

If nonsynonymous (replacement) substitutions are neutral:  
 $d_N > d_S$  or  $d_N \approx d_S$  or  $d_N < d_S$  ?

BIOS477/877 L5 - 14

14

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$
- Number of **nonsynonymous** substitutions per site:  $d_N$

→  $d_N/d_S$  can be used to identify selection

If nonsynonymous (replacement) substitutions are neutral:  
 $d_N > d_S$  or  $d_N \approx d_S$  or  $d_N < d_S$  ?

$d_N = u_T * f_0$   
 $d_S = u_T$

If nonsynonymous substitutions are neutral,  
 →  $f_0 = 1$ ,  
 →  $d_N = u_T * f_0 = u_T = d_S$

BIOS477/877 L5 - 16

15

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$
- Number of **nonsynonymous** substitutions per site:  $d_N$

→  $d_N/d_S$  can be used to identify selection

When nonsynonymous (replacement) substitutions are neutral,  
 $d_N \approx d_S \rightarrow d_N/d_S \approx 1$

→  $d_S$  is used as a control  
 (substitution rate at neutral or near-neutral)

BIOS477/877 L5 - 16

16

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$
- Number of **nonsynonymous** substitutions per site:  $d_N$

→  $d_N/d_S$  can be used to identify selection

$d_N/d_S > 1.0$   
 $d_N/d_S \approx 1.0$  Neutral (e.g., non-coding regions, pseudogenes)  
 $d_N/d_S < 1.0$

→  $d_S$  is used as a control  
 (substitution rate at neutral or near-neutral)

BIOS477/877 L5 - 17

17

### Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$
- Number of **nonsynonymous** substitutions per site:  $d_N$

→  $d_N/d_S$  can be used to identify selection

$d_N/d_S > 1.0$   
 Which indicates negative selection?  
 $d_N/d_S < 1.0$

→  $d_S$  is used as a control  
 (substitution rate at neutral or near-neutral)

BIOS477/877 L5 - 18

18

## Identifying Selection

- Number of **synonymous** substitutions per site:  $d_S$
- Number of **nonsynonymous** substitutions per site:  $d_N$
- $d_N/d_S$  can be used to identify selection

Adaptive selection could be driving amino acid substitutions!

Positive selection

1.0 = Neutral (e.g., non-coding regions, pseudogenes)

(Negative but relaxed selection)

Negative selection (under functional constraints)

0

→  $d_S$  is used as a control (substitution rate at neutral or near-neutral)

BIOS477/877 L5 - 20

20

## Identifying Selection

exon

Sequence 1 ATGGCTTCACCAACAGAACATCATCTTTGTGCCGGTCTGGCGCATTGGTCTGAC

Sequence 2 ATGGCTTCACCAACAGAACATCATCTTTGTGCCGGTCTGGCGCATTGGTCTGAC

Sequence 1 ACCAGTCCGAAATGCTCAACCGGATCTCAAGSTTGTGTCACCTTATTTATTTGTT

Sequence 2 ACCAGTCCGAAATGCTCAACCGGATCTCAAGSTTGTGTCACCTTATTTATTTGTT

intron

Sequence 1 TTCTTTCCAAAATCTACTTTGTTTCCCGGTGGTTAG

Sequence 2 TTCTTTCCAAAATCTACTTTGTTTCCCGGTGGTTAG

[Codon position]	Exon			Intron		
	1st	2nd	3rd	1st	2nd	3rd
# nucleotide sites	31	31	31	23	22	22
# nucleotide substitutions	2	0	4	12	12	14
Nucleotide substitutions/site	0.06	0	0.13	0.52	0.55	0.64

Exon(2nd) < Exon(3rd)

Intron(2nd)  $\lesssim$  or  $\approx$  Intron(3rd)

Exon(2nd) / Exon(3rd)  $\ll$  Intron(2nd) / Intron(3rd)

(Similar to  $d_N/d_S$  analysis)

BIOS477/877 L5 - 21

21

## Identifying Selection

Rank	Protein	Divergence (%)
1	Transition protein 2	66
48	ZP2	43
59	Protamine P15	41
86	Spam protein 10	39
92	Testis histone H1	38
101	Acrosin	38
120	Protamine 2	36
161	ZP3	33
194	Testes Tpx1	31

Rapid evolution of reproductive proteins

Why are they evolving fast?  
Weak selective constraints? (relaxed selection) or Divergence is favored? (positive selection)

10% most divergent proteins from 1,880 human/rodent orthologues

Figure 1 | **Rapidly evolving proteins.** Comparison of 1,880 human-rodent orthologues from Makalowski & Boguski<sup>2</sup> plotted as a frequency of the occurrence of genes with a varying percentage of amino-acid divergence. The portion that contains the 10% most divergent proteins is shown in blue; reproductive proteins that are among the 10% most divergent proteins are listed. Tpx1, testes-specific protein 1; ZP2/3, zona pellucida 2/3.

Swanson & Vacquier (2002) *Nature Reviews Genetics* 3: 137-144

BIOS477/877 L5 - 22

22

## Selection and functions

- Elevated  $d_N/d_S$  ratios are often found in reproductive genes

Male reproductive proteins (e.g., genes involved in mating behavior, fertilization, spermatogenesis, ejaculation, or sex determination)

Nonreproductive proteins

Swanson *et al.* (2001) *Evolutionary EST analysis identifies rapidly evolving male reproductive proteins in Drosophila*. *PNAS* 98: 7375-7379.

See also Findlay *et al.* (2010) *Proteomics enhances evolutionary and functional analysis of reproductive proteins*. *BioEssays* 32: 26-36.

Wilburn and Swanson (2016) *From molecules to mating: Rapid evolution and biochemical studies of reproductive proteins*. *J. Proteomics* 135: 12-25.

BIOS477/877 L5 - 23

23

## Selection and functions

- Elevated  $d_N/d_S$  ratios are often found in reproductive genes

Male reproductive proteins (e.g., genes involved in mating behavior, fertilization, spermatogenesis, ejaculation, or sex determination)

Nonreproductive proteins

Swanson *et al.* (2001) *Evolutionary EST analysis identifies rapidly evolving male reproductive proteins in Drosophila*. *PNAS* 98: 7375-7379.

See also Findlay *et al.* (2010) *Proteomics enhances evolutionary and functional analysis of reproductive proteins*. *BioEssays* 32: 26-36.

Wilburn and Swanson (2016) *From molecules to mating: Rapid evolution and biochemical studies of reproductive proteins*. *J. Proteomics* 135: 12-25.

BIOS477/877 L5 - 24

24

## Selection and functions

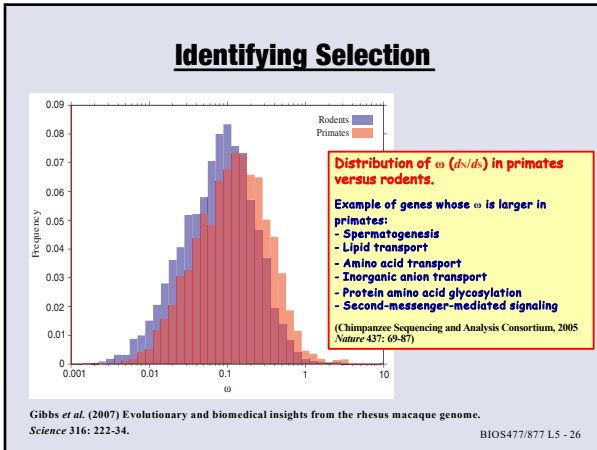
Table 1. Selected examples of protein-coding genes in which positive selection was detected by using the  $d_N/d_S$  ratio

Gene	Species	Function	Refs
<b>Genes involved in defense</b>			
Class I chitinase gene	<i>Drosophila</i>		61
Colicin genes	<i>Escherichia coli</i>		62
Defensin genes	Roberts		63
P17	Mus		64
Immunoglobulin V <sub>H</sub> genes	Mammals		65
MHC genes	Mammals		66
Polylactonase inhibitor genes	Legume and dicots		67
R1 blood group and RH50 genes	Primates and rodents		68
Ribonuclease genes	Primates		69
Transferrin gene	Salmonid fishes		70
Type I interferon gene	Mammals		71
$\alpha$ -Proteinase inhibitor genes	Roberts		72
<b>Genes involved in evading defensive systems or immunity</b>			
Capsid gene	FluV virus		73
CSP, TSP, A5A2 and P102	<i>Plasmodium falciparum</i>		74
DNA-methylase coding region	Hepatitis D virus		75
F gene	Phages G1, $\phi$ X174, $\phi$ 29		76
Evolvec			77
gr1/gr2			78
Hereditary fructose intolerance			79
Melanin			80
major lig	<i>Anguilla marmorata</i>		81
Oral membrane protein gene	Chlamydia		82
Polylactonase genes	Fungal pathogens		83
Porn protein 1 gene	Nematodes		84
S and H <sub>2</sub> glycoprotein genes	Murine coronavirus		85
Sigma 1 protein gene	Reovirus		86
Viral defense gene	Tetrahymena		87

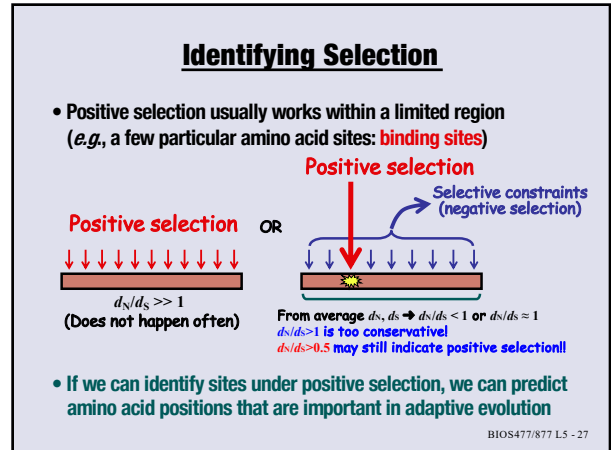
Yang & Bielawski (2000) *Statistical methods for detecting molecular adaptation*. *TREE* 15: 496-503.

BIOS477/877 L5 - 25

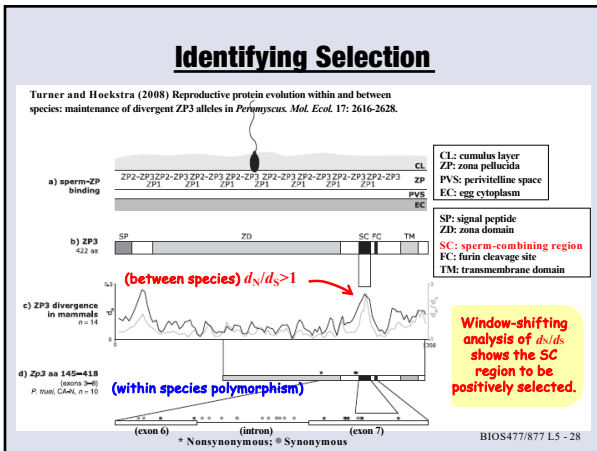
25



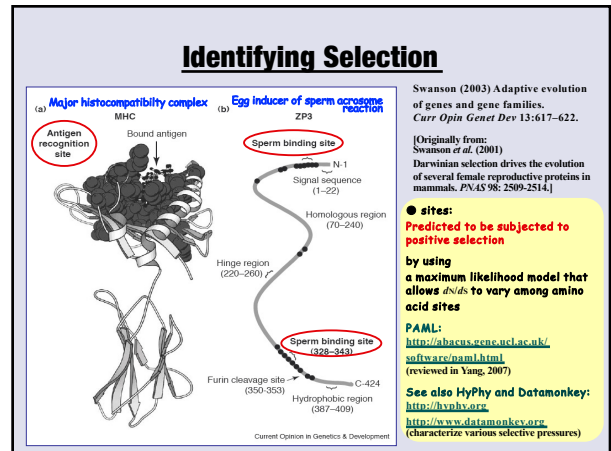
26



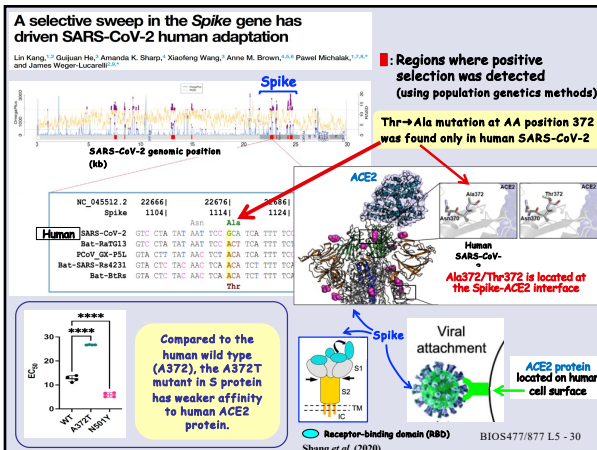
27



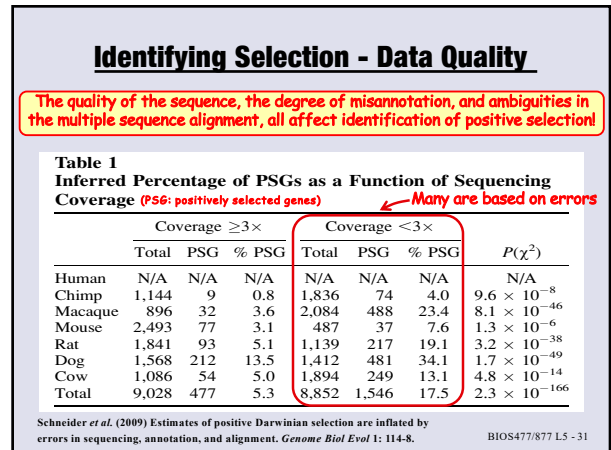
28



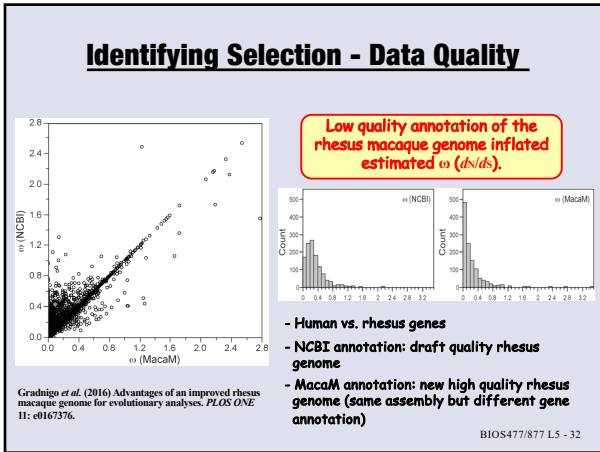
29



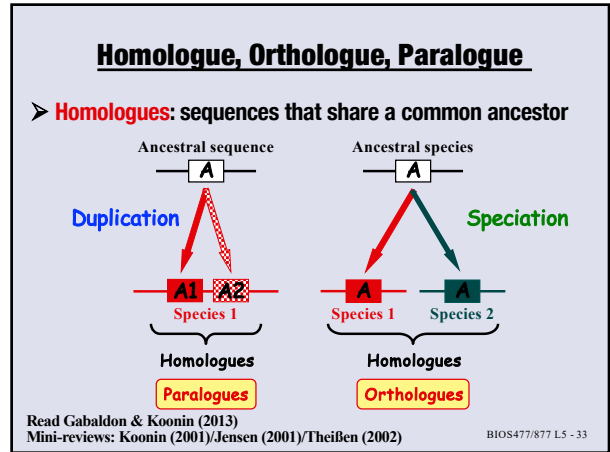
30



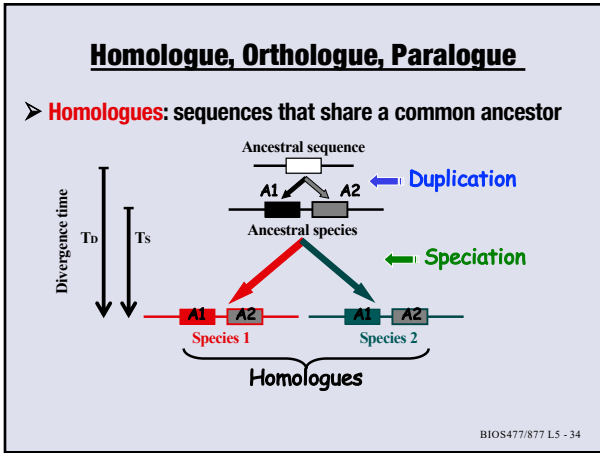
31



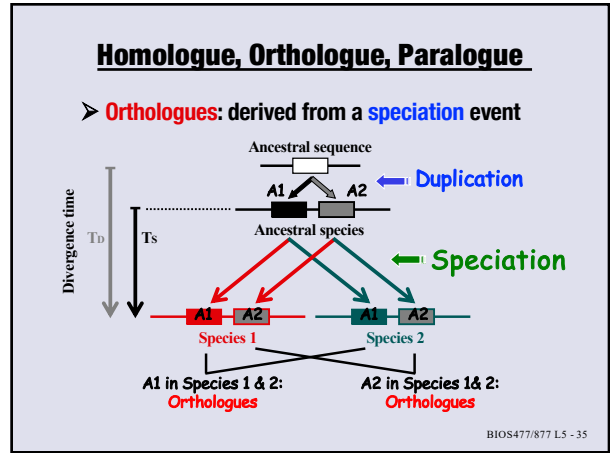
32



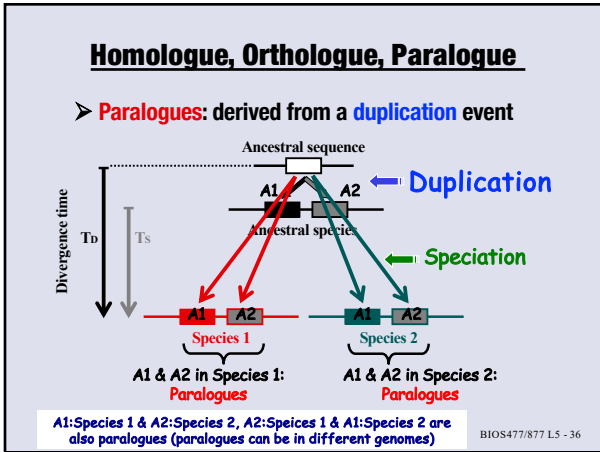
33



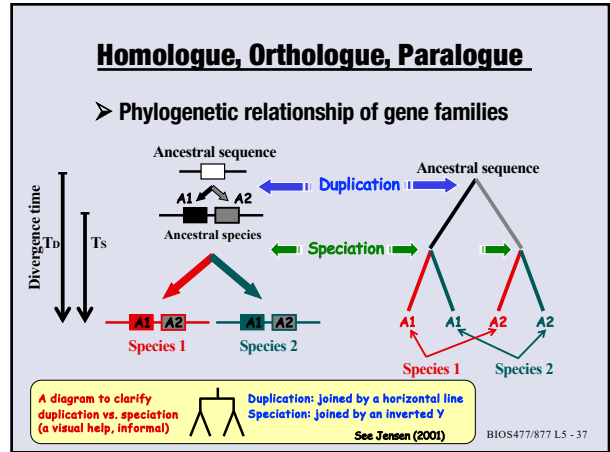
34



35



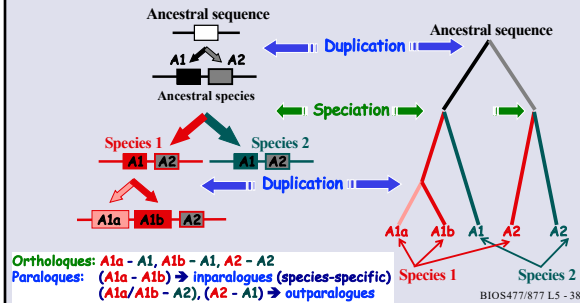
36



37

## Homologue, Orthologue, Parologue

➤ Duplications can happen anytime



38

## Gene Duplication

➤ Gene duplication creates **large gene families**  
 ➤ Important means of providing a substrate on which evolution can work

- 15% or more of human genes are duplicates!

1) Duplicated copies: **Redundancy**

- Weaker functional constraint (neutral or near neutral)
- One copy can be degenerated to become a **pseudogene**
- **Non-functionalization** ... likely to be lost

2) One copy can acquire a novel gene function

- **Neo-functionalization**, rare but important

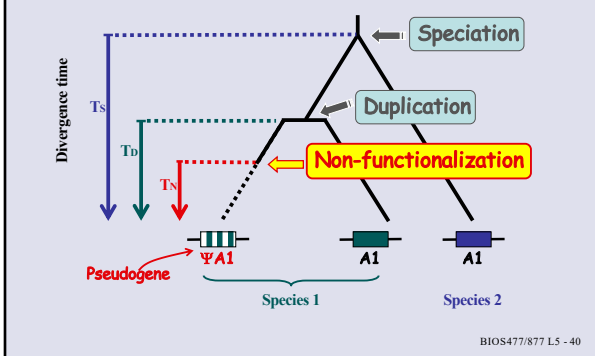
3) **Sub-functionalization?**

- Functional complementation between duplicated genes

BIOS477/877 L5 - 39

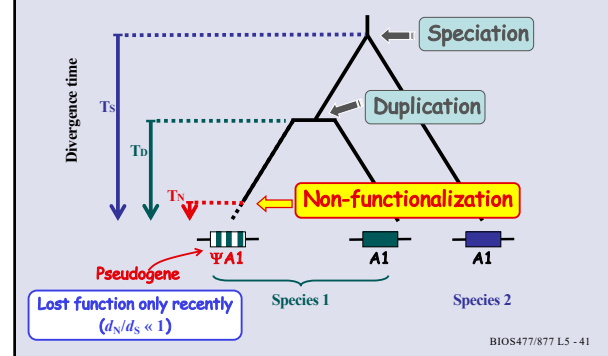
39

## Pseudogene



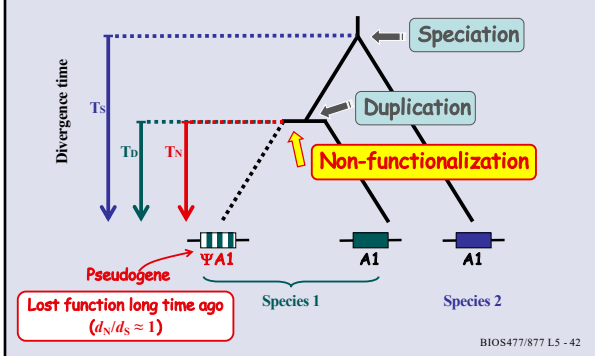
40

## Pseudogene



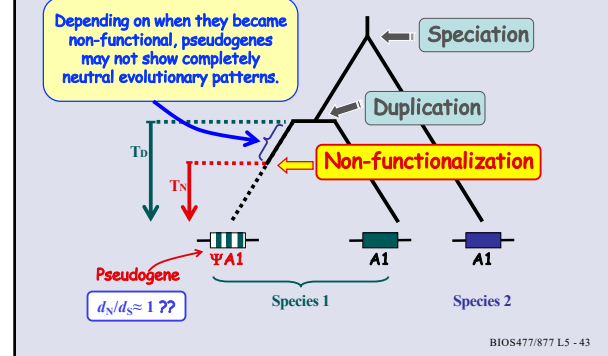
41

## Pseudogene



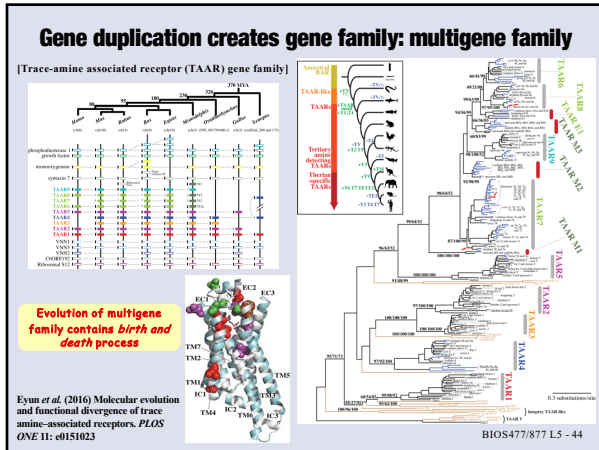
42

## Pseudogene



43





44

### Homology vs. Similarity

- **Similarity:** the extent to which sequences are related.  
→ makes no statement about descent from a common ancestor
- **Homology:** sequence similarity that can be attributed to **descent from a common ancestor**

**Homology ≠ Similarity !!**

- Sequences can be either **homologous** or non-homologous, but not in between (*e.g.*, you cannot say two genes are 10% homologous!)
- **Homology** is not directly measurable or observable.
- **Similarity** is a direct measurement.

Higgs and Attwood (2005) Chapter 1 page 8  
Read also Petsko (2001), Koonin (2001), Jensen (2001)

BIOS477/877 L5 - 45

45

### Next time ...

- **Pairwise sequence comparison by dotplot**
  - **Dotlet JS** <https://dotlet.vital-it.ch>
  - **Dotlet** (old Java version) <https://myhits.isb-sib.ch/cgi-bin/dotlet>  
[For security reasons, old Java programs are not available within UNL network]
  - **DotMatcher** (A program in EMBOSS)  
<http://emboss.bioinformatics.nl/cgi-bin/emboss/dotmatcher>  
(see course Web for other EMBOSS links)
  - **JDotter: Java Dot Plot Alignments**  
<http://pgrc.ink-gatersleben.de/jdotter/>
  - **YASS: Visualization of local pairwise alignments**  
<https://bioinfo.univ-lille.fr/yass/index.php>
  - **Dotter (part of SeqTools)** See the Course website for more programs  
<http://www.sanger.ac.uk/science/tools/seqtools>

BIOS477/877 L5 - 46

46