

Spring 2025
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) 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) 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. $d_N(C)$

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

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

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

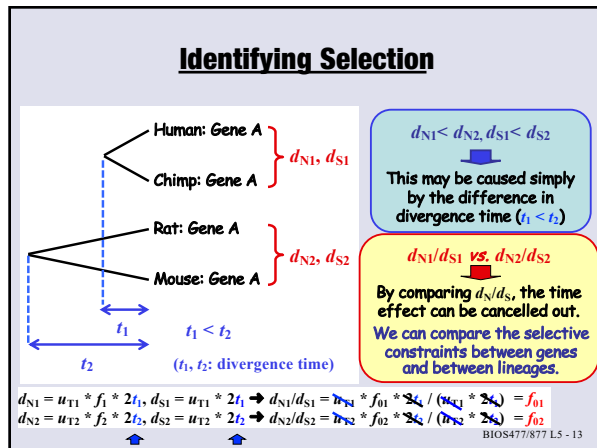
$d_{N1} < d_{N2}, d_{S1} < d_{S2}$

This may be caused simply by the difference in divergence time ($t_1 < t_2$)

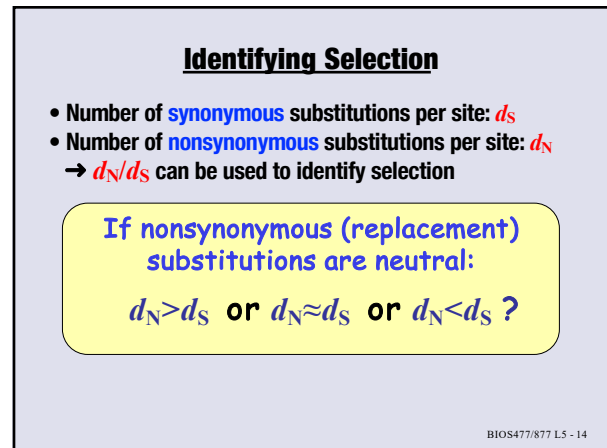
$d_{N1} = u_{T1} * f_1 * 2t_1, d_{S1} = u_{T1} * 2t_1$
 $d_{N2} = u_{T2} * f_2 * 2t_2, d_{S2} = u_{T2} * 2t_2$

BIOS477/877 L5 - 12

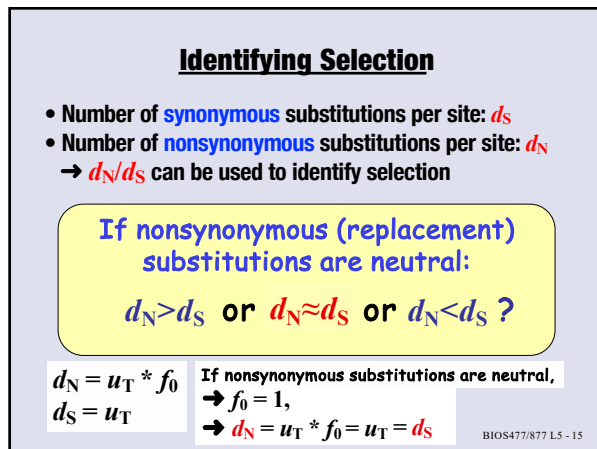
12



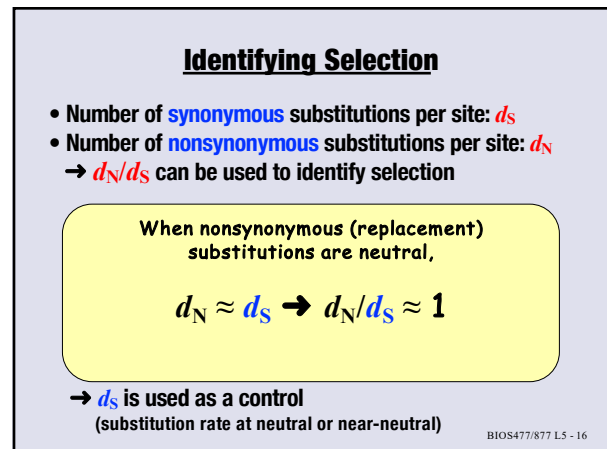
13



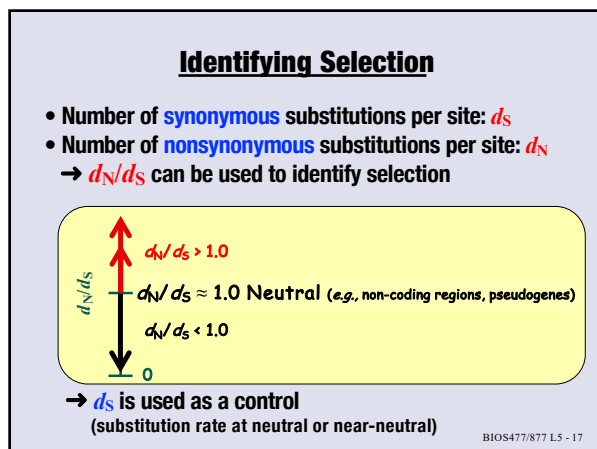
14



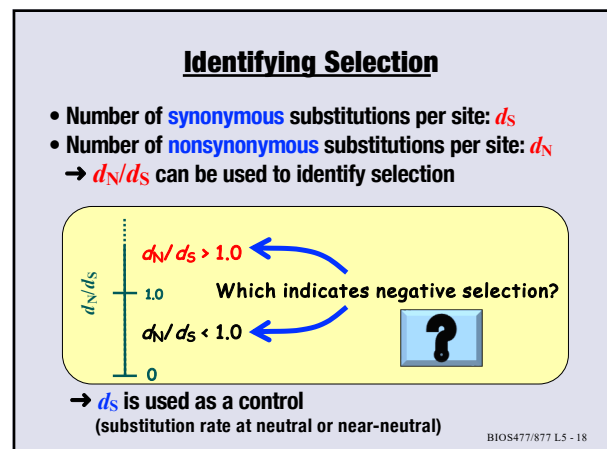
15



16



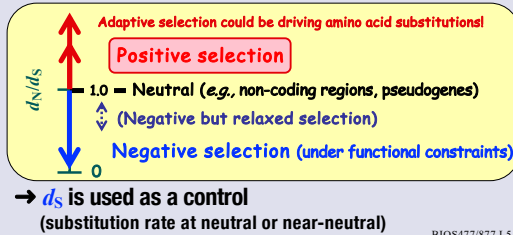
17



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



BIOS477/877 L5 - 20

20

Identifying Selection

Sequence 1: **exon** ATGCGCTCTACCAACAGAACATCTCTTTGTGCGCGGCTGCGCGGCTTGGTCTGAC

Sequence 2: ATGCGCTCTACCAACAGAACATCTCTTTGTGCGCGGCTGCGCGGCTTGGTCTGAC

Sequence 1: ACCAGTCGGGAGTTCTCAGCGGSAITCTCAGGTTTGTCACCTTATTTATTTGTT

Sequence 2: ACCAGTCGGGAGTTCTCAGCGGSAITCTCAGGTTTGTCACCTTATTTATTTGTT

Sequence 1: TTGTTTTCGAAATCTTACTTTGTTTTCGCGGTGGTTAG

Sequence 2: TTGTTTTCGAAATCTTACTTTGTTTTCGCGGTGGTTAG

Sequence 1: TTCCAGAGATTTACCTTTTCTTTTCTTTTCTTTTCTAG

Sequence 2: TTCCAGAGATTTACCTTTTCTTTTCTTTTCTTTTCTAG

intron

	Exon			Intron		
[Codon position]	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) ≤ or ≈ Intron(3rd)

Exon(2nd) / Exon(3rd) << Intron(2nd) / Intron(3rd)

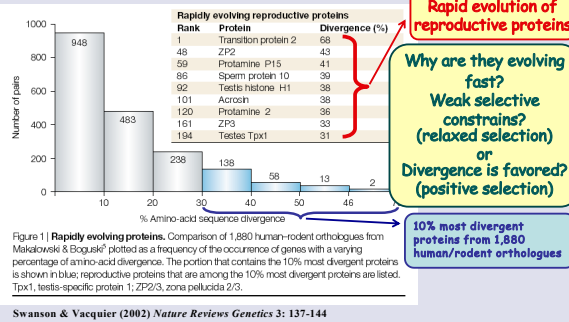
0 << 0.86 (≈1)

(Similar to d_N/d_s analysis)

BIOS477/877 L5 - 21

21

Identifying Selection

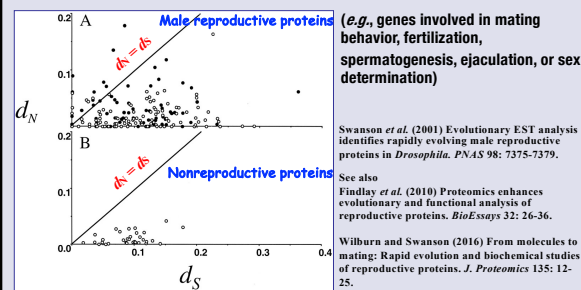


BIOS477/877 L5 - 22

22

Selection and functions

- Elevated d_N/d_s ratios are often found in reproductive genes

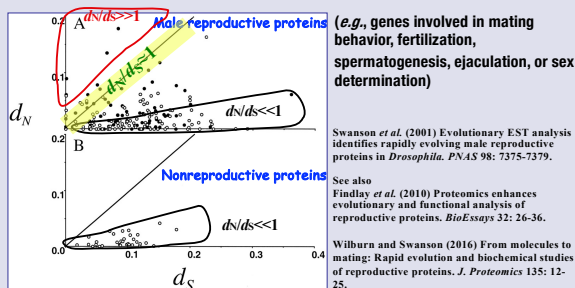


BIOS477/877 L5 - 23

23

Selection and functions

- Elevated d_N/d_s ratios are often found in reproductive genes



BIOS477/877 L5 - 24

24

Selection and functions

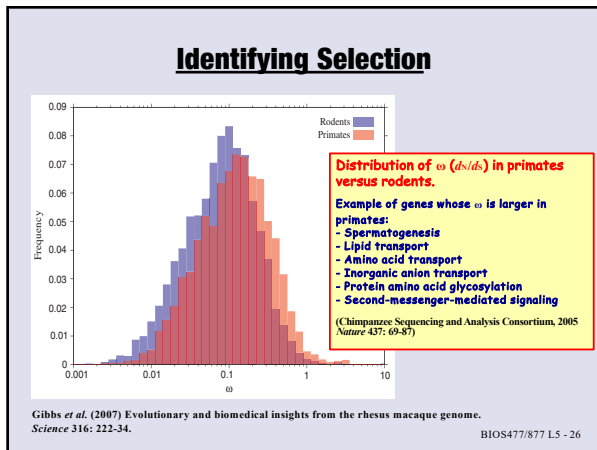
Table 1. Selected examples of protein-coding genes in which positive selection was detected by using the d_N/d_s ratio

Gene	Ratio
Genes involved in defense	
Os11L1 gene	61
Colicin gene	62
Defensin genes	63
FcγR	64
Interleukin-1α genes	65
MHC genes	66
Polypeptide chain inhibitor genes	67
Rh blood group and RH50 genes	68
Ribonuclease genes	69
Transferrin gene	70
Type I interferon gene	71
α ₁ Protease inhibitor gene	72
Genes involved in evolving defense systems or immunity	
Capid gene	73
CSP, TRAP, AEA-2 and P9B3	74
Defensin coding region	75
E gene	76
Enolase gene	77
gH glycoprotein	78
Hemagglutinin	79
Invasion gene	80
Mucosin	81
mpg 1a	82
mpg 1b	83
Outer membrane protein gene	84
Polypeptide chain inhibitor genes	85
Port protein 1 gene	86
S and H glycoprotein genes	87
Sigma 1 protein gene	88
Variation determinant gene	89
Defensive systems or immunity	
Escherichia coli	90
Escherichia coli	91
Escherichia coli	92
Escherichia coli	93
Escherichia coli	94
Escherichia coli	95
Escherichia coli	96
Escherichia coli	97
Escherichia coli	98
Escherichia coli	99
Escherichia coli	100
Escherichia coli	101
Escherichia coli	102
Escherichia coli	103
Escherichia coli	104
Escherichia coli	105
Escherichia coli	106
Escherichia coli	107
Escherichia coli	108
Escherichia coli	109
Escherichia coli	110
Escherichia coli	111
Escherichia coli	112
Escherichia coli	113
Escherichia coli	114
Escherichia coli	115
Escherichia coli	116
Escherichia coli	117
Escherichia coli	118
Escherichia coli	119
Escherichia coli	120
Escherichia coli	121
Escherichia coli	122
Escherichia coli	123
Escherichia coli	124
Escherichia coli	125
Escherichia coli	126
Escherichia coli	127
Escherichia coli	128
Escherichia coli	129
Escherichia coli	130
Escherichia coli	131
Escherichia coli	132
Escherichia coli	133
Escherichia coli	134
Escherichia coli	135
Escherichia coli	136
Escherichia coli	137
Escherichia coli	138
Escherichia coli	139
Escherichia coli	140
Escherichia coli	141
Escherichia coli	142
Escherichia coli	143
Escherichia coli	144
Escherichia coli	145
Escherichia coli	146
Escherichia coli	147
Escherichia coli	148
Escherichia coli	149
Escherichia coli	150
Escherichia coli	151
Escherichia coli	152
Escherichia coli	153
Escherichia coli	154
Escherichia coli	155
Escherichia coli	156
Escherichia coli	157
Escherichia coli	158
Escherichia coli	159
Escherichia coli	160
Escherichia coli	161
Escherichia coli	162
Escherichia coli	163
Escherichia coli	164
Escherichia coli	165
Escherichia coli	166
Escherichia coli	167
Escherichia coli	168
Escherichia coli	169
Escherichia coli	170
Escherichia coli	171
Escherichia coli	172
Escherichia coli	173
Escherichia coli	174
Escherichia coli	175
Escherichia coli	176
Escherichia coli	177
Escherichia coli	178
Escherichia coli	179
Escherichia coli	180
Escherichia coli	181
Escherichia coli	182
Escherichia coli	183
Escherichia coli	184
Escherichia coli	185
Escherichia coli	186
Escherichia coli	187
Escherichia coli	188
Escherichia coli	189
Escherichia coli	190
Escherichia coli	191
Escherichia coli	192
Escherichia coli	193
Escherichia coli	194
Escherichia coli	195
Escherichia coli	196
Escherichia coli	197
Escherichia coli	198
Escherichia coli	199
Escherichia coli	200

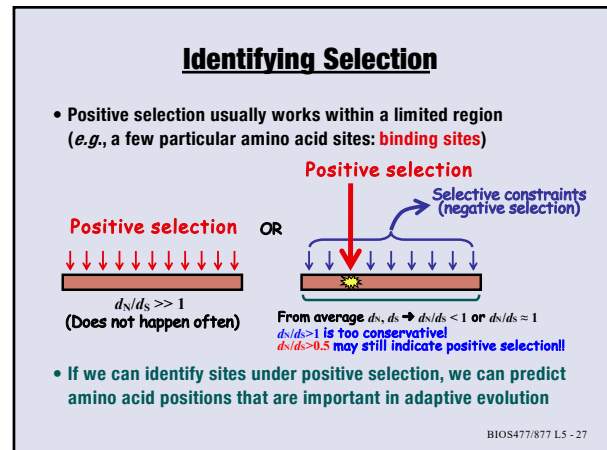
Yang & Bielawski (2000) *Statistical methods for detecting molecular adaptation*. *TREE* 15: 496-503.

BIOS477/877 L5 - 25

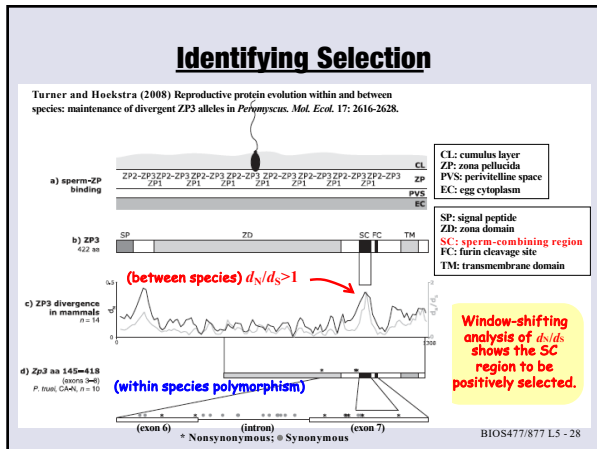
25



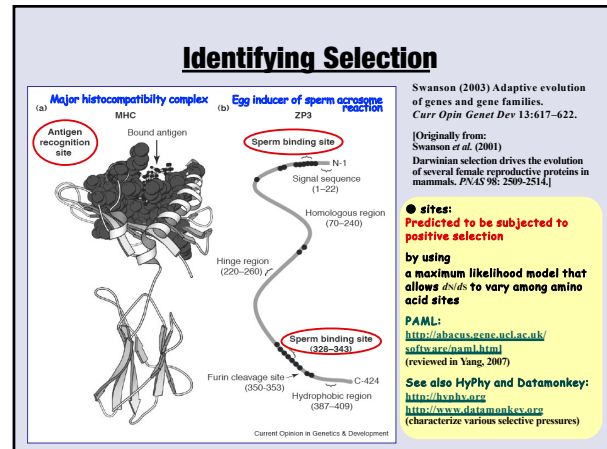
26



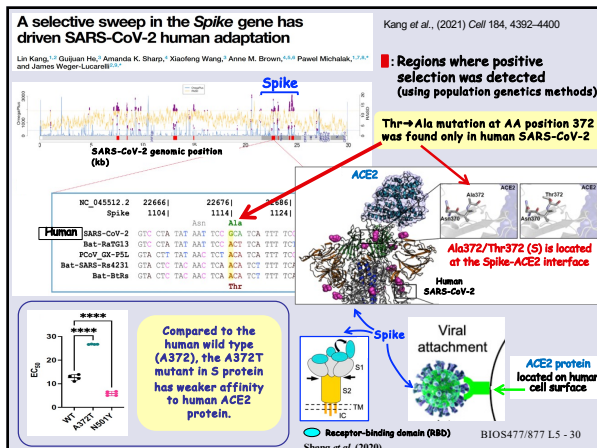
27



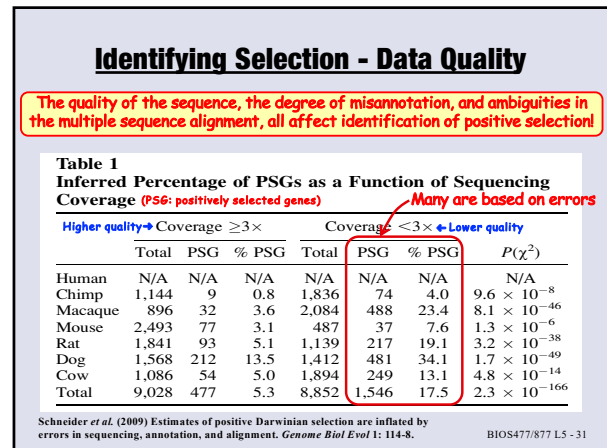
28



29

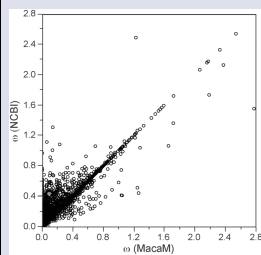


30

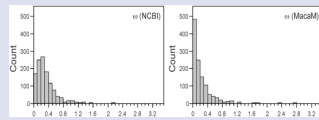


31

Identifying Selection - Data Quality



Low quality annotation of the rhesus macaque genome inflated estimated ω (ds/ds).



- Human vs. rhesus genes
- NCBI annotation: draft quality rhesus genome
- MacAM annotation: new high quality rhesus genome (same assembly but different gene annotation)

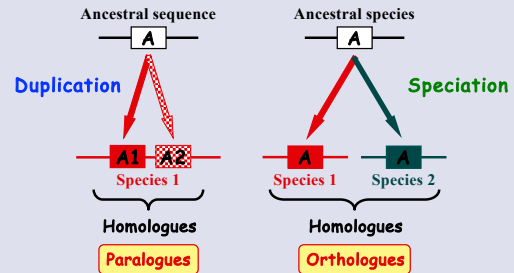
Gratino *et al.* (2016) Advantages of an improved rhesus macaque genome for evolutionary analyses. *PLoS ONE* 11: e0167376.

BIOS477/877 L5 - 32

32

Homologue, Orthologue, Parologue

- **Homologues:** sequences that share a common ancestor



Read Gabaldon & Koonin (2013)

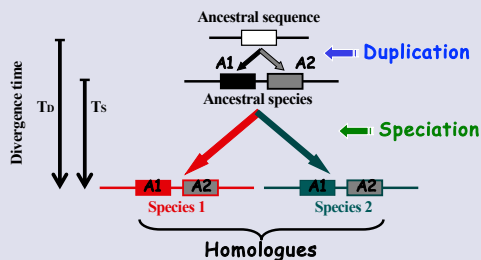
Mini-reviews: Koonin (2001)/Jensen (2001)/Theißen (2002)

BIOS477/877 L5 - 33

33

Homologue, Orthologue, Parologue

- **Homologues:** sequences that share a common ancestor

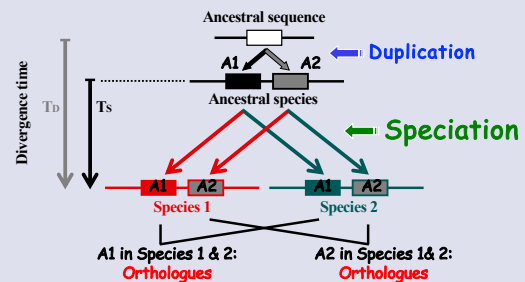


BIOS477/877 L5 - 34

34

Homologue, Orthologue, Parologue

- **Orthologues:** derived from a speciation event

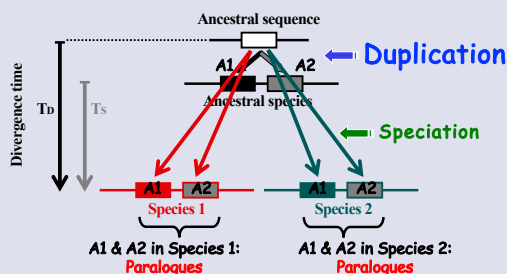


BIOS477/877 L5 - 35

35

Homologue, Orthologue, Parologue

- **Paralogues:** derived from a duplication event



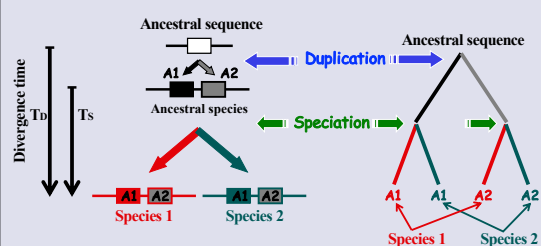
A1:Species 1 & A2:Species 2, A2:Species 1 & A1:Species 2 are also paralogues (paralogues can be in different genomes)

BIOS477/877 L5 - 36

36

Homologue, Orthologue, Parologue

- **Phylogenetic relationship of gene families**



A diagram to clarify duplication vs. speciation (a visual help, informal)

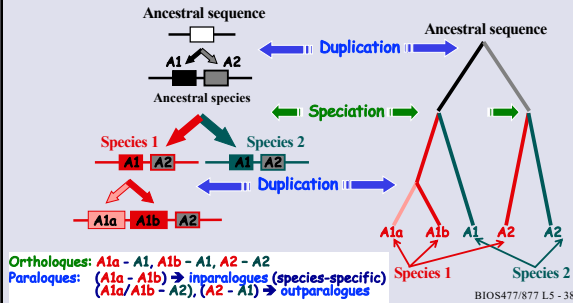
Duplication: joined by a horizontal line
Speciation: joined by an inverted Y
See Jensen (2001)

BIOS477/877 L5 - 37

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

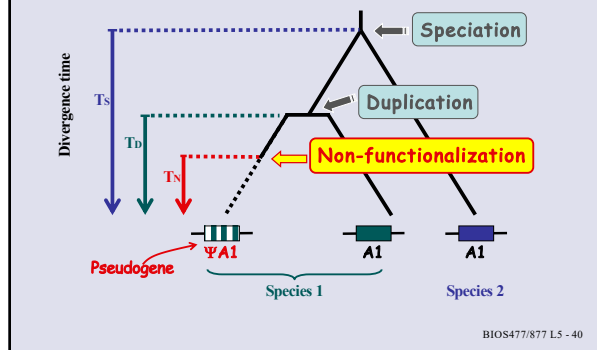
3) Sub-functionalization?

- Functional complementation between duplicated genes

BIOS477/877 L5 - 39

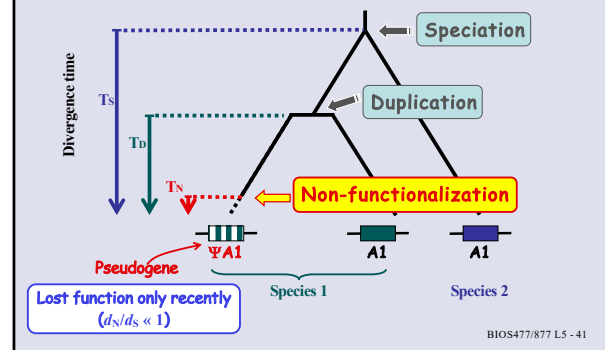
39

Pseudogene



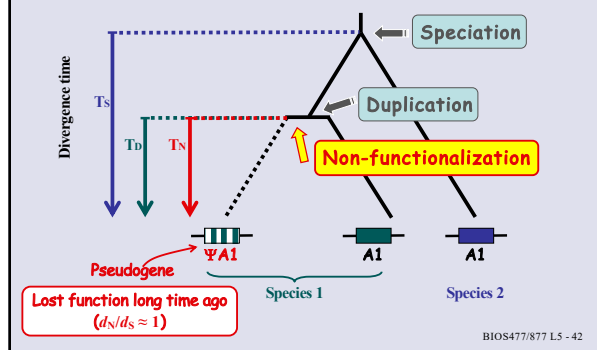
40

Pseudogene



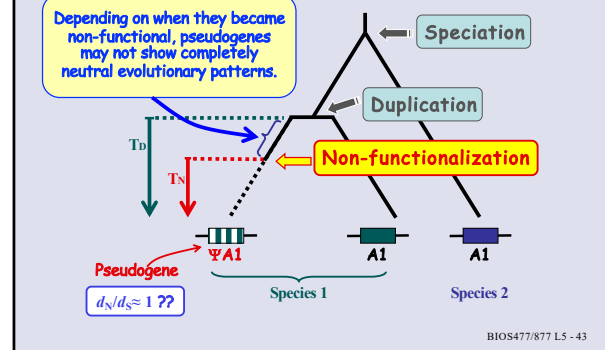
41

Pseudogene

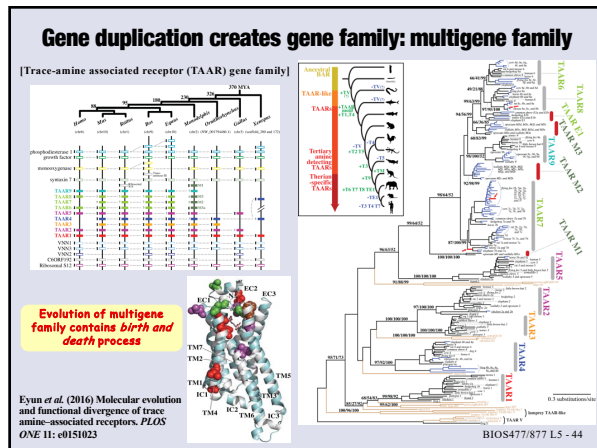


42

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.ipk-gatersleben.de/jdotter/>
 - **YASS: Visualization of local pairwise alignments**
<https://bioinfo.univ-lille.fr/yass/index.php>

See the Course website for more programs

BIOS477/877 L5 - 46

46