

# Large-scale selection analyses

Maria Anisimova

Institute of Computational Life Sciences  
Zurich University of Applied Sciences - ZHAW

# Two decades of large-scale selection scans

## Inferring Nonneutral Evolution from Human-Chimp-Mouse Orthologous Gene Trios

Andrew G. Clark,<sup>1</sup> Stephen Glanowski,<sup>3</sup> Rasmus  
Paul D. Thomas,<sup>4</sup> Anish Kejariwal,<sup>4</sup> Melissa A.  
David M. Tanenbaum,<sup>5</sup> Daniel Civello,<sup>6</sup> Fu Lu,<sup>5</sup> Brian  
Steve Ferreira,<sup>3</sup> Gary Wang,<sup>3</sup> Xianqun Zh  
Thomas J. White,<sup>6</sup> John J. Sninsky,<sup>6</sup> Mark D. A  
Michele Cargill<sup>6†</sup>

2003, Science

## A Scan for Positively Selected Genes in the Genomes of Humans and Chimpanzees

Rasmus Nielsen<sup>1,2\*</sup>, Carlos Bustamante<sup>1</sup>, Andrew G. Clark<sup>3</sup>, Stephen Glanowski<sup>4</sup>, Timothy B. Sackton<sup>3</sup>,  
Melissa J. Hubisz<sup>1</sup>, Adi Fledel-Alon<sup>1</sup>, David M. Tanenbaum<sup>5</sup>, Daniel Civello<sup>6</sup>, Thomas J. White<sup>6</sup>,  
John J. Sninsky<sup>6</sup>, Mark D. Adams<sup>5a</sup>, Michele Cargill<sup>6</sup>

Open access, freely available online PLOS BIOLOGY

2005

OPEN ACCESS Freely available online

## Patterns of Positive Selection in Six Mammalian Genomes

Carolin Kosiol<sup>1</sup>, Tomáš Vinař<sup>1</sup>, Rute R. da Fonseca<sup>2</sup>, Melissa J. Hubisz<sup>3</sup>, Carlos D. Bustamante<sup>1</sup>, Rasmus  
Nielsen<sup>2</sup>, Adam Siepel<sup>1\*</sup>

PLOS GENETICS

2008

Research article

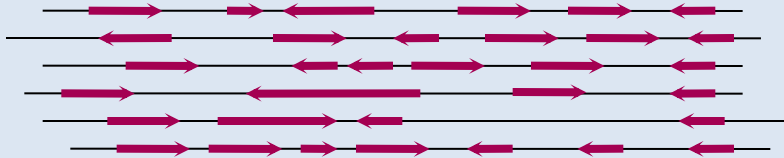
## A systematic search for positive selection in higher plants (Embryophytes)

Christian Roth<sup>1,2,3</sup> and David A Liberles<sup>\*1,3</sup>

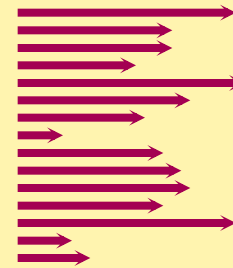
2006

# Large-scale selection scans step-by-step

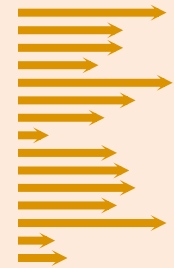
CDs from complete genomes /databases



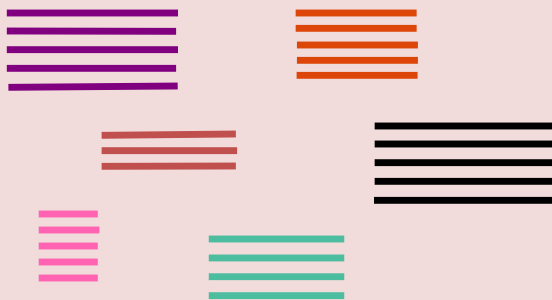
DNA CDs library



AA library



MSAs (both AA and CDs)



Homologous gene clusters  
(unaligned AA data)



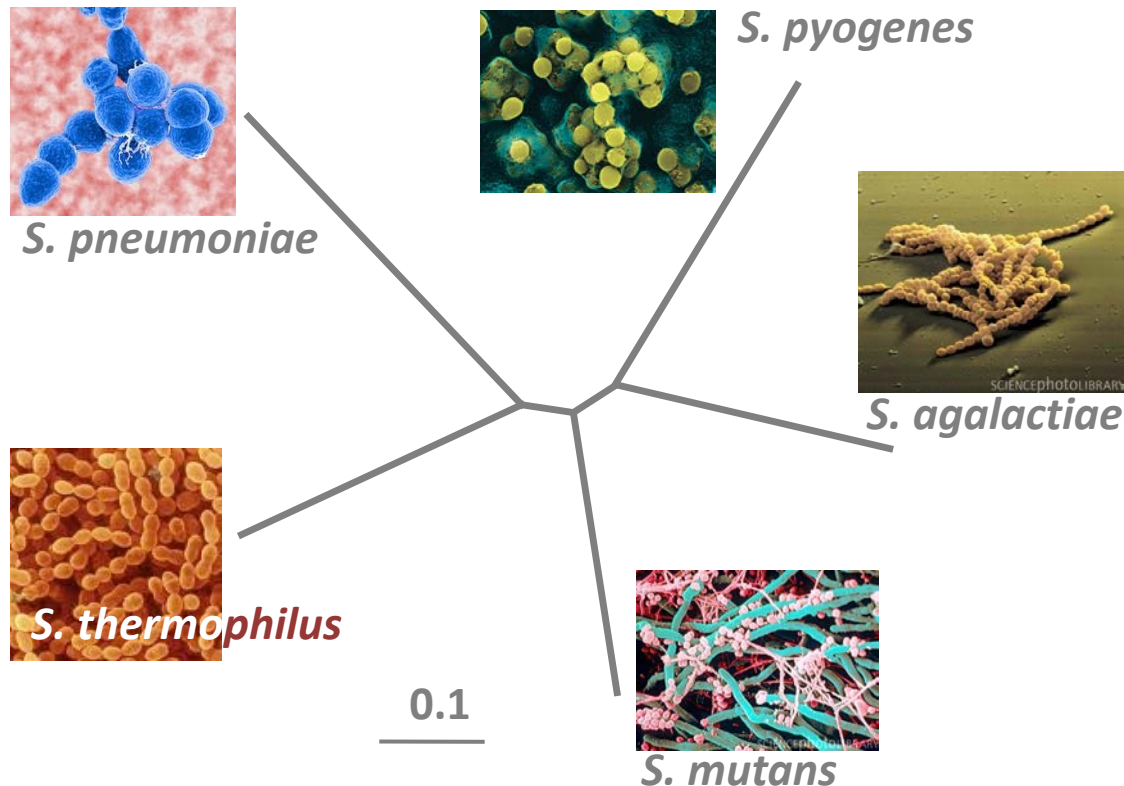
Statistical  
analysis



0.236251 9929 1.1542 0.00012 1.1542 0.03029 1.1542 0.03029  
0.6767251 9945 1.1542 0.03029 1.1542 0.03029  
0.1111251 9981 1.1542 0.01998 0.1111251 9981  
0.230004 9139 1.1542 0.00912 0.1111251 9981 1.1542 ..

What do we do with all these numbers?

# Natural selection in Streptococcus



Anisimova et al 2007 BMC Evol Biol

**12 complete genomes**

**Positive selection in 136 genes:**

29% connected to virulence

10% no ascribable function

7% essential to *S. pneumoniae*

19% with body-site specific

patterns of gene expression

during invasive disease in *S.*

*pyogenes* (infected blood,

cerebrospinal fluid,

epithelial cell contact)

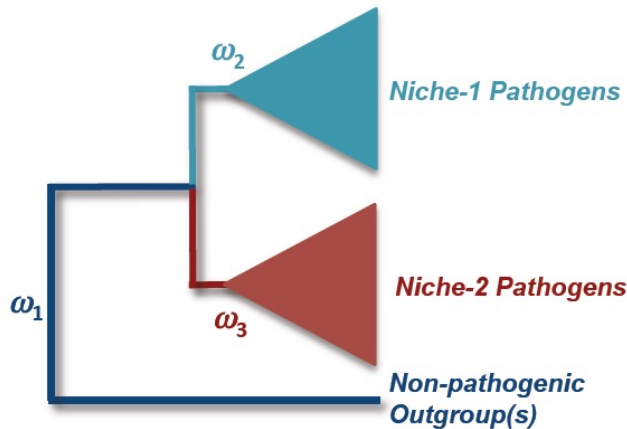
Positive selection affects both core and accessory genes,  
most likely due to the antagonistic interaction between host and parasite.

Products of both core and auxiliary genes participate in complex networks that  
comprise the molecular basis of virulence.

# Listeria phylogenomics

## Mapping selection to phenotype

### A: Gene-level data analysis



Null model (1 parameter):

$$H_0 : \omega_1 = \omega_2 = \omega_3$$

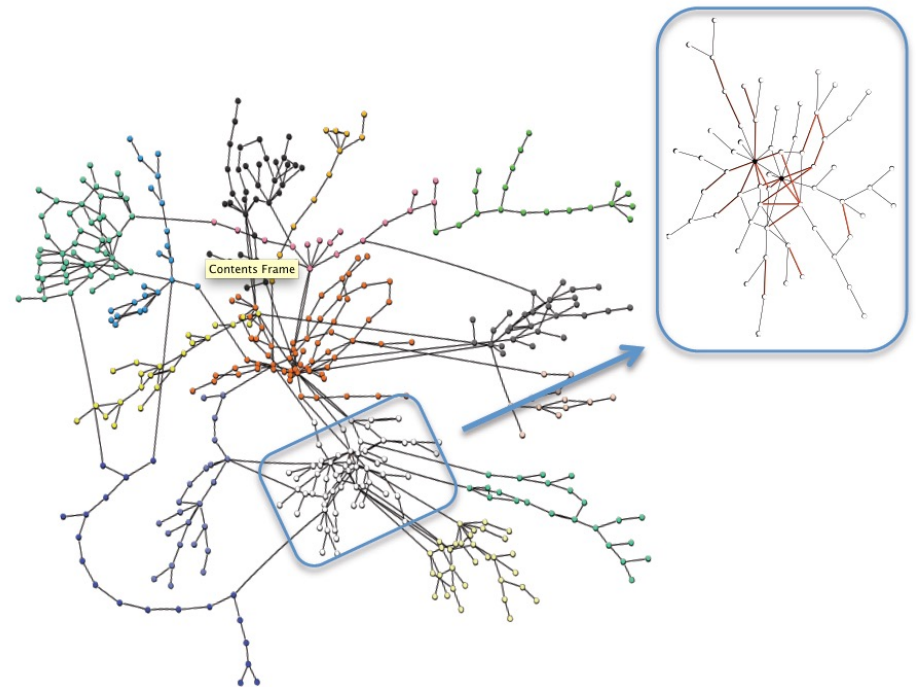
Alternatives (2 parameters):

$$H_1 : \omega_1 \neq \omega_2 = \omega_3$$

$$H_2 : \omega_1 = \omega_2 \neq \omega_3$$

$$H_3 : \omega_1 = \omega_3 \neq \omega_2$$

### B: Phenotype-level data analysis

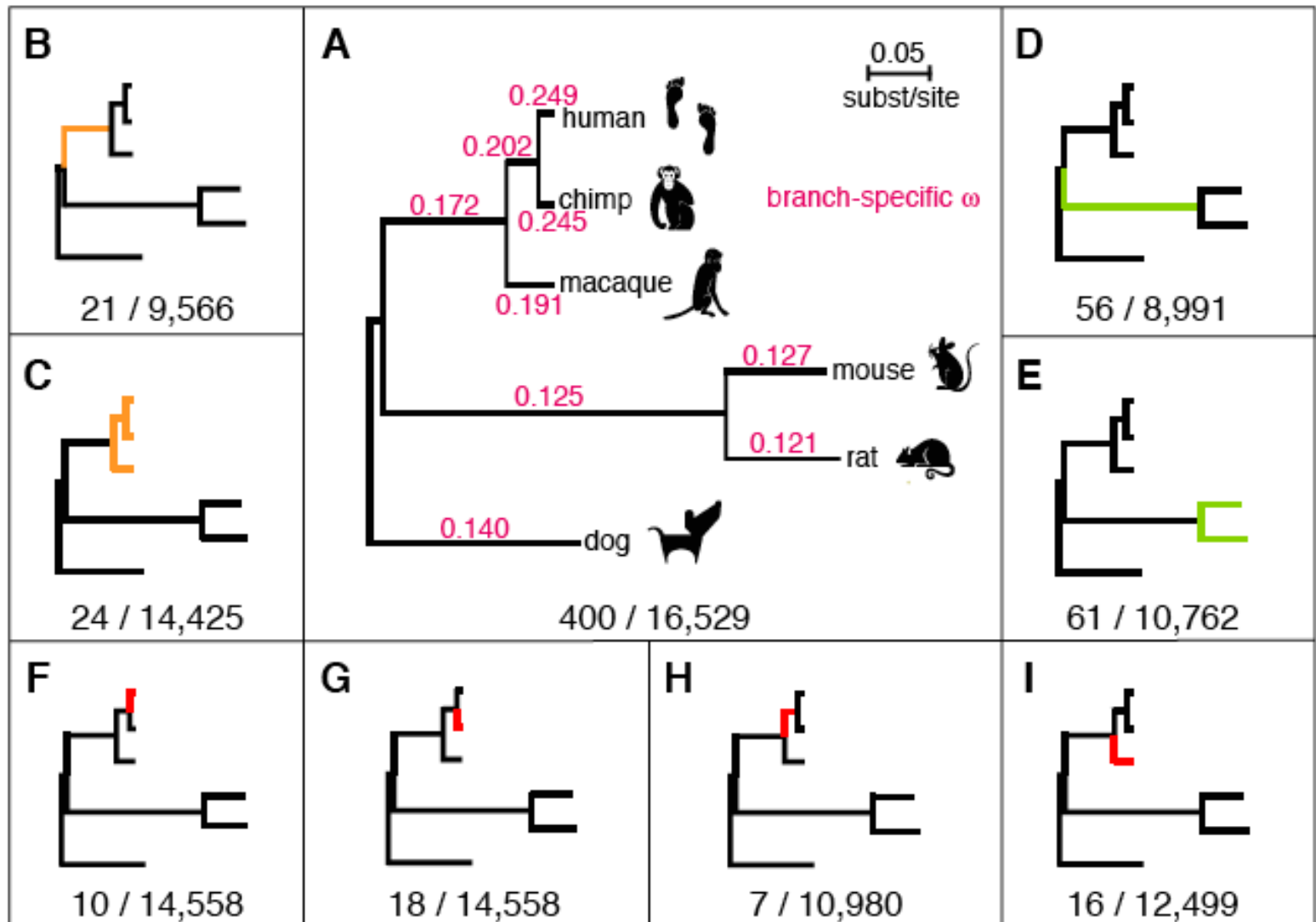


Blue box identifies a module in the metabolic network. Red links in the expanded view of this module indicate a significant cluster of genes subject to niche specific selection in "lineage I" of *L. monocytogenes*.

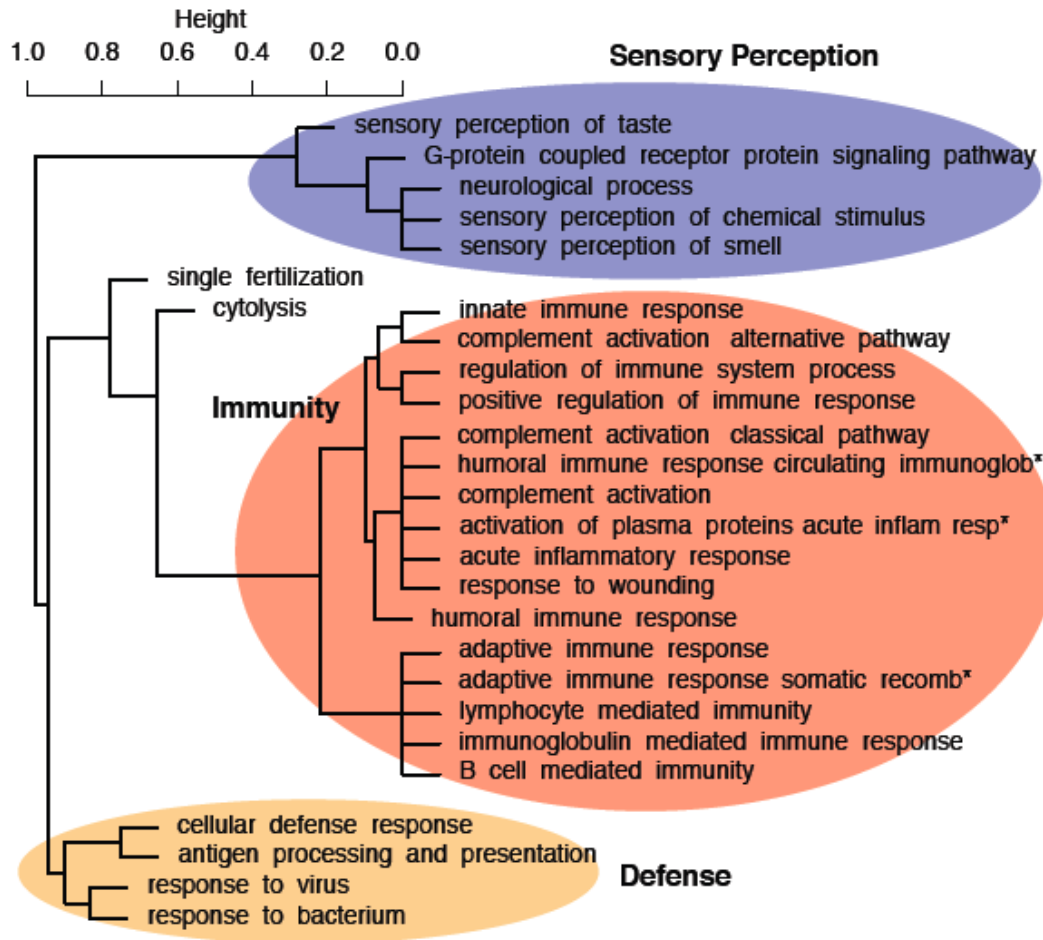


# Multiple LRTs: scan of mammalian genomes

*Kosiol et al (2008)*



# Multiple LRTs: scan of mammalian genomes



Dissimilarity measure:

$$d_{AB} = 1 - \frac{|N(A) \cap N(B)|}{\min\{|N(A)|, |N(B)|\}}$$

Hierarchical clustering of GO categories (biological process) over-represented with genes under positive selection

# Which proteins are under positive selection?

- Host proteins involved in defence or immunity against viral, bacterial, fungal or parasite attacks (MHC, immunoglobulin VH, class 1 chitinas).
- Viral or pathogen proteins involved in evading host defence (HIV env, nef, gap, pol, etc., capsid in FMD virus, flu virus hemagglutinin gene).
- Proteins or pheromones involved in reproduction (abalone sperm lysin, sea urchin bindin, proteins in mammals)
- Proteins that acquired new functions after gene duplication.
- Miscellaneous (diet, globins, etc. )





Conclusions



# Detecting positive selection

- Pairwise methods – very low power
- Branch models allow variation over time but assume one  $\omega$  for all sites - low power
- Site models allow variation among sites but assume selection pressure does not change over time – have higher power if positive selection is long term
- Branch-site models may be more successful at detecting episodic selection but are more difficult to fit, require more data and often have multiple sub-optimal peaks (caution with genome scans!)

# Testing for positive selection

- LRT is accurate even for small datasets
- Power of LRT is better for larger datasets
- Watch out for recombination
- Accurate parameter estimation is more difficult, depends on model assumptions
- Bayesian site prediction is even more difficult than LRTs and parameter estimation
- There is an optimal window of sequence divergence (sequences should be not too similar and not saturated)
- Robustness of results: Use several models & tests
- Check for local optima, especially for complex models

# Weaknesses of methods based on codon models

- Model assumptions may be unrealistic (but some assumptions matter more than others)
- The method detects positive selection only if it generates excessive nonsynonymous substitutions. It may lack power in detecting one-off directional selection or when the sequences are highly similar or highly divergent. Little power with population data.
- Do not work for noncoding DNA (but see Wong & Nielsen 2003 Genetics)
- Sensitive to sequence and alignment errors (Fletcher & Yang 2010 Mol Biol Evol 27; Privman et al. 2011 Mol Biol Evol 29; Jordan & Goldman 2012 Mol Biol Evol 29)

# Criticisms on codon models

by M. Nei, Y. Suzuki, & A.L. Hughes

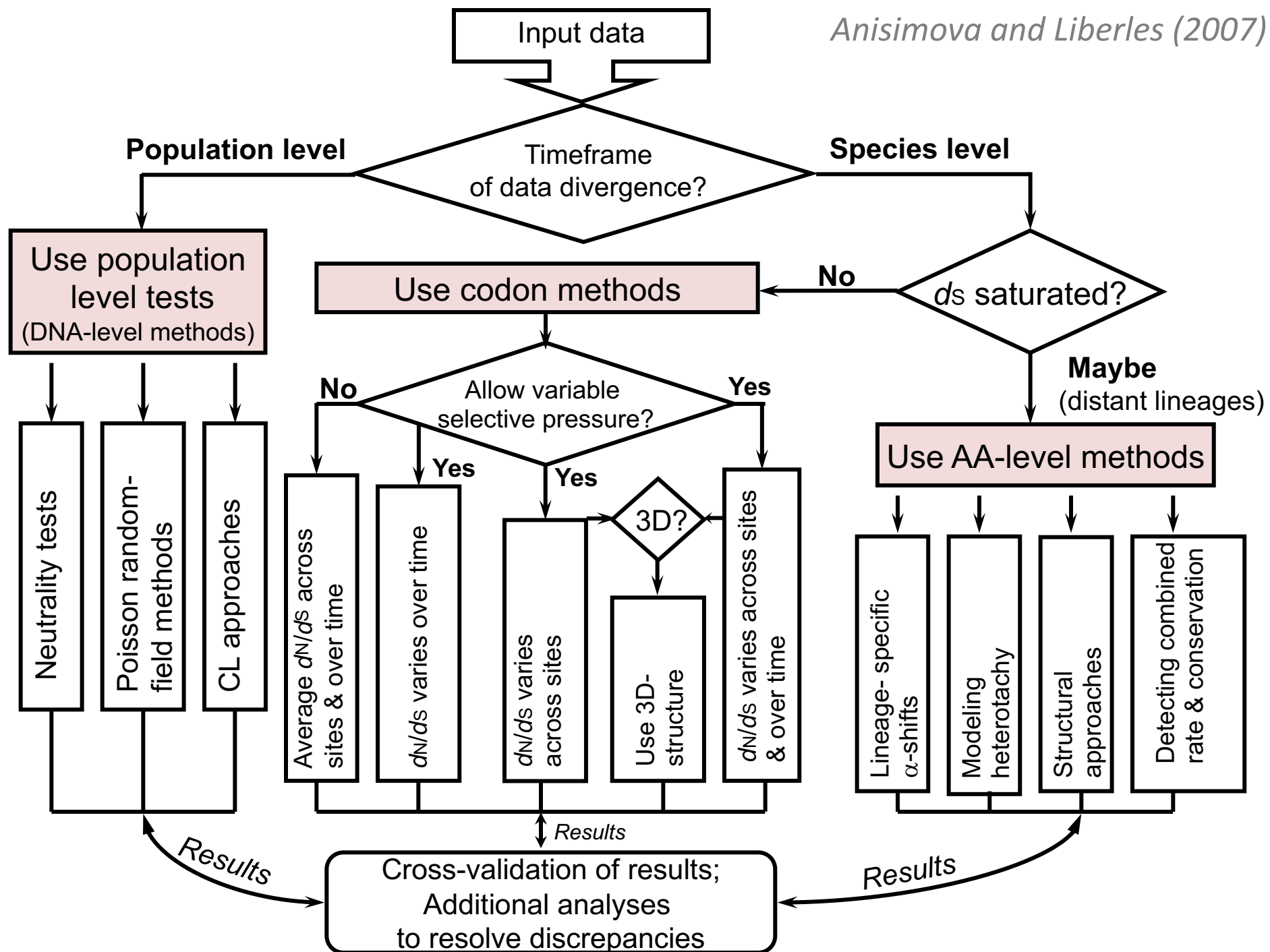
Hughes AL. 2007. Looking for Darwin in all the wrong places: the misguided quest for positive selection at the nucleotide sequence level. *Heredity* 99

Nozawa, Suzuki & Nei. 2009. *PNAS* 106

Yang Z, dos Reis M. 2011. Statistical properties of the branch-site test of positive selection. *Mol Biol Evol* 28

Zhai W, Nielsen R, Goldman N, Yang Z. 2012. Looking for Darwin in genomic sequences - validity and success of statistical methods. *Mol Biol Evol* 29

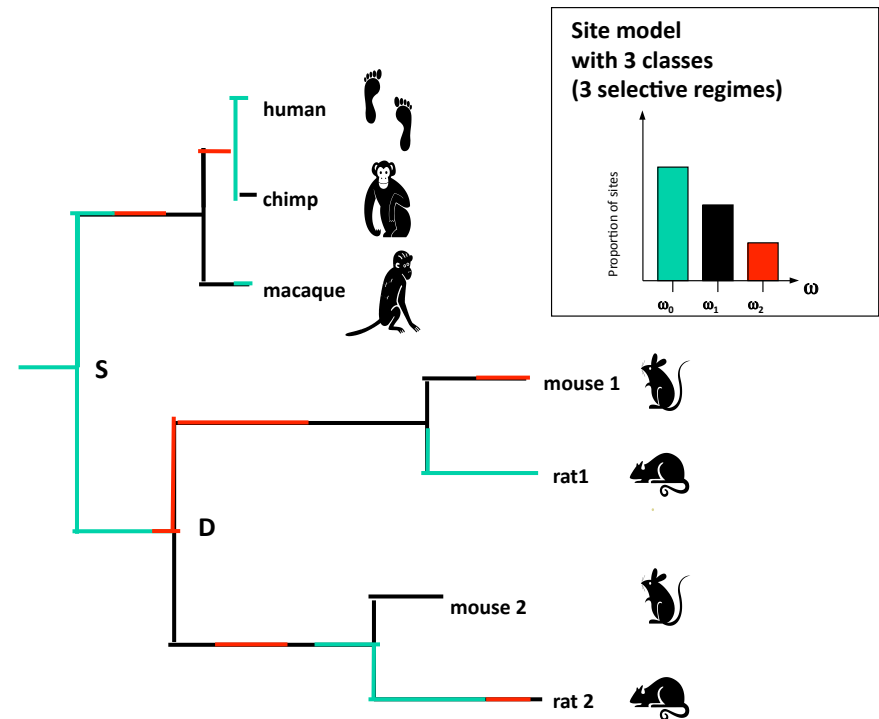
MacCallum, C. & Hill, E. 2006 Being positive about selection. *PLoS Biol* 4, e87





# The many faces of codon models

- Detecting selection
- Studying codon bias
- Inferring phylogenies
- Dating speciation events
- Ancestral reconstruction
- Changes in time & space
- Predicting coding regions
- Improved alignment
- Inferring gene features (phyloHMM, netHMM)
- Simulation of data



**Markov modulated model:**  
Guindon et al. 2004

**Reviews of codon models:**  
Kosiol and Anisimova 2012  
Anisimova and Kosiol 2009