

Tests of Positive Selection based on the Comparison of Polymorphism and Divergence

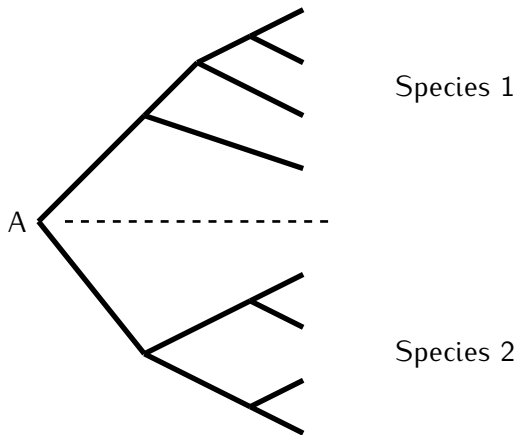
Julien Dutheil

`dutheil@evolbio.mpg.de`

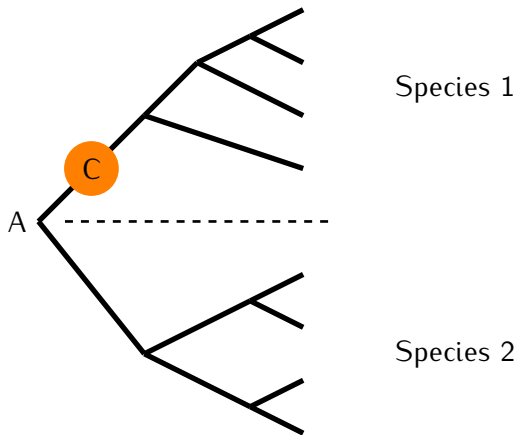
Max Planck Institute for Evolutionary Biology

June 22nd 2015

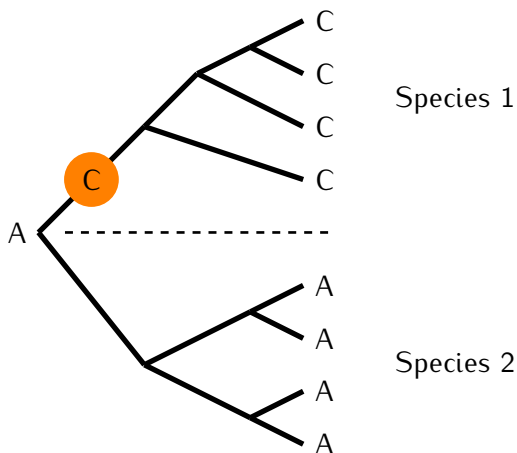
Within vs. between species



Within vs. between species

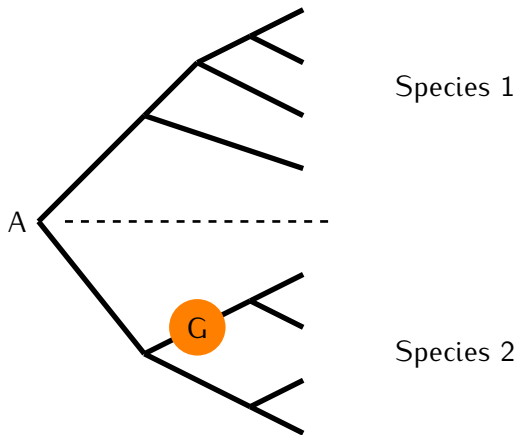


Within vs. between species

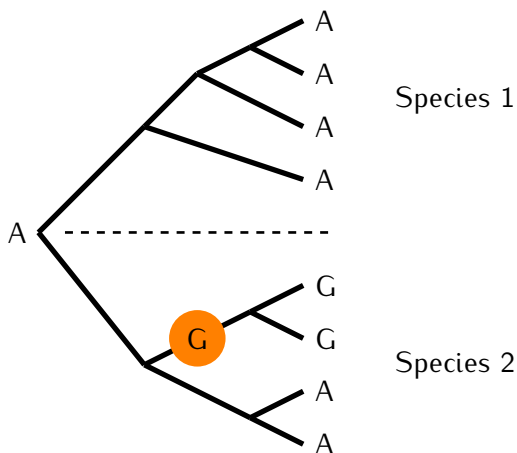


- Mutations on *interspecies branches* lead to fixed differences between species

Within vs. between species



Within vs. between species



- Mutations on *interspecies branches* lead to fixed differences between species
- Mutations on *intraspecies branches* lead to polymorphism in one species

If all mutations are neutral...

- ☞ The ratio of polymorphic sites vs. fixed differences sites is constant along the genome!

If all mutations are neutral...

☞ The ratio of polymorphic sites vs. fixed differences sites is constant along the genome!

If mutation rate varies between sites but is constant over time in the two species, two predictions:

- 1 the ratio of polymorphism vs. divergence is constant between genes

If all mutations are neutral...

☞ The ratio of polymorphic sites vs. fixed differences sites is constant along the genome!

If mutation rate varies between sites but is constant over time in the two species, two predictions:

- 1 the ratio of polymorphism vs. divergence is constant between genes
- 2 the ratio of non-synonymous to synonymous polymorphism equals the ratio of non-synonymous to synonymous divergence

Polymorphism and divergence are two facets of the same process

The HKA test

Hudson, Kreitman and Aguadé (1987)

- Compare at least 2 loci in 2 species, with polymorphism data in at least 1 species

The HKA test

Hudson, Kreitman and Aguadé (1987)

- Compare at least 2 loci in 2 species, with polymorphism data in at least 1 species
- If mutation rate is constant in time:
 - Regions with high mutation rate display high levels of polymorphism and divergence
 - Regions with low mutation rate display low levels of polymorphism and divergence

The HKA test

Hudson, Kreitman and Aguadé (1987)

- Compare at least 2 loci in 2 species, with polymorphism data in at least 1 species
- If mutation rate is constant in time:
 - Regions with high mutation rate display high levels of polymorphism and divergence
 - Regions with low mutation rate display low levels of polymorphism and divergence
- 'Goodness-of-fit' test to assess how consistent distinct regions are with a constant mutation rate

The HKA test

Hudson, Kreitman and Aguadé (1987)

- Compare at least 2 loci in 2 species, with polymorphism data in at least 1 species
- If mutation rate is constant in time:
 - Regions with high mutation rate display high levels of polymorphism and divergence
 - Regions with low mutation rate display low levels of polymorphism and divergence
- 'Goodness-of-fit' test to assess how consistent distinct regions are with a constant mutation rate
- Assumes free recombination between regions and no recombination within regions

The MK test

McDonald and Kreitman (1991)

- One coding gene in at least 2 species, with polymorphism data for at least 1 species

McDonald and Kreitman (1991)

- | | | D simulans | | | | | D yakuba | | | | | | | | | | | |
|------|------|-------------------------|---|-----|---|---|-------------|---|---|---|---|-------------------------|---|---|---|-------|-------|--|
| | Con. | a b o d e f g h i j k l | | | | | a b o d e f | | | | | a b o d e f g h i j k l | | | | | | |
| 781 | G | T | T | T | T | T | T | T | T | T | T | - | - | - | - | Repl. | Fixed | |
| 789 | T | - | - | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 808 | A | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 816 | G | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 834 | T | T | T | T | - | - | T | T | T | T | T | - | - | - | - | Syn. | Fixed | |
| 859 | C | - | - | - | - | - | C | C | - | - | - | - | - | - | - | Syn. | Fixed | |
| 867 | C | - | - | - | - | - | - | - | - | - | - | G | G | G | G | Syn. | Fixed | |
| 879 | C | T | T | T | T | T | T | T | T | T | T | - | - | - | - | Syn. | Fixed | |
| 900 | G | - | - | - | - | - | - | - | - | - | - | A | - | - | - | Syn. | Fixed | |
| 974 | G | - | - | - | - | - | - | - | - | - | - | T | T | T | T | Syn. | Fixed | |
| 983 | T | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 1019 | C | - | - | - | - | - | - | - | - | - | - | - | A | - | - | Syn. | Fixed | |
| 1031 | A | - | - | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1034 | T | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 1043 | C | - | - | - | - | - | - | - | - | - | - | - | A | - | - | Syn. | Fixed | |
| 1060 | C | T | T | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1069 | C | - | - | - | - | - | A | A | A | A | A | A | A | A | A | Syn. | Fixed | |
| 1101 | G | - | - | - | - | - | - | - | - | - | - | A | A | A | A | Syn. | Fixed | |
| 1127 | T | - | - | - | - | - | - | - | - | - | - | T | C | C | C | Syn. | Fixed | |
| 1131 | T | - | - | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1160 | T | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 1175 | T | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 1178 | C | - | - | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1184 | C | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 1190 | C | - | - | - | - | - | - | - | - | - | - | - | A | - | - | Syn. | Fixed | |
| 1196 | C | - | - | - | - | - | - | - | - | - | - | T | T | T | T | Syn. | Fixed | |
| 1199 | C | - | T | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1202 | T | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 1203 | C | - | - | - | - | - | - | - | - | - | - | - | T | - | - | Syn. | Fixed | |
| 1229 | T | - | G | C | C | C | C | C | C | C | C | - | - | - | - | Syn. | Fixed | |
| 1232 | T | - | - | - | - | - | - | - | - | - | - | A | A | A | A | Syn. | Fixed | |
| 1235 | C | - | - | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1244 | C | - | - | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1265 | C | - | - | - | - | - | - | - | - | - | - | G | G | G | G | Syn. | Fixed | |
| 1271 | A | - | - | - | - | - | - | - | - | - | - | T | - | T | - | Syn. | Fixed | |
| 1277 | T | - | - | - | - | - | - | - | - | - | - | C | C | C | C | Syn. | Fixed | |
| 1283 | C | A | A | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1298 | C | - | - | - | - | - | - | - | - | - | - | T | T | T | T | Syn. | Fixed | |
| 1304 | C | - | - | - | - | - | - | - | - | - | - | - | - | - | - | Syn. | Fixed | |
| 1316 | C | - | - | -</ | | | | | | | | | | | | | | |

The MK test

McDonald and Kreitman (1991)

- One coding gene in at least 2 species, with polymorphism data for at least 1 species
- Count synonymous and non-synonymous polymorphisms and fixed differences
- Build a contingency table and perform a G-test

	Fixed	Polym.
Non-syn.	7	2
Synon.	17	42

		D. melanogaster											D. simulans						D. yakuba														
	Cos.	a	b	c	d	e	f	g	h	i	j	k	l	a	b	c	d	e	f	a	b	c	d	e	f	g	h	i	j	k	l		
781	G	T	T	T	T	T	T	T	T	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Repl.	Fixed
789	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	C	C	C	C	C	C	C	C	C	C	Syn.	Fixed
833	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	G	G	G	G	G	G	G	G	G	G	G	Repl.	Fixed
845	T	T	T	T	-	-	-	-	-	-	-	-	-	T	T	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
859	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
867	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	G	G	G	G	G	G	G	G	G	G	G	2 Poly.	
870	C	T	T	T	T	T	T	T	T	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Fixed
950	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
974	G	-	-	-	-	-	-	-	-	-	-	-	-	T	-	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
983	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	C	C	C	C	C	C	C	C	C	C	Syn.	Fixed
1019	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1031	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1034	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1043	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1069	C	T	T	T	T	T	T	T	T	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1089	C	-	-	-	-	-	-	-	-	-	-	-	-	A	A	A	A	A	A	-	-	-	-	-	-	-	-	-	-	-	-	Repl.	Fixed
1101	G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	A	A	A	A	A	A	A	A	A	A	A	A	Repl.	Fixed
1127	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Fixed
1131	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1160	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	C	C	C	C	C	C	C	C	C	C	Syn.	Fixed
1175	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	C	C	C	C	C	C	C	C	C	C	Syn.	Fixed
1178	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1184	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	G	G	G	G	G	G	G	G	G	G	G	Syn.	Fixed
1190	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1196	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	T	T	T	T	-	-	-	-	-	-	-	-	Syn.	Poly.
1202	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	C	C	C	C	C	C	C	C	C	C	Syn.	Fixed
1203	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1229	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1232	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	A	A	A	A	A	A	A	A	A	A	A	A	Syn.	Fixed
1235	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1246	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1265	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	G	G	G	G	G	G	G	G	G	G	G	Syn.	Fixed
1271	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1277	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	C	C	C	C	C	C	C	C	C	C	Syn.	Fixed
1283	C	A	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1289	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Fixed
1304	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1316	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1405	C	A	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1431	T	C	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1443	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1459	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1480	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Repl.	Poly.
1504	C	T	T	T	T	T	T	T	T	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Fixed
1518	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1524	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	G	G	G	G	G	G	G	G	G	G	G	G	Syn.	Fixed
1527	C	T	T	T	T	T	T	T	T	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1530	G	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1545	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1548	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1551	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1558	C	A	A	A	A	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Repl.	Poly.
1559	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1560	G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1573	G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	C	C	C	C	C	C	C	C	C	C	C	Repl.	Fixed
1581	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	T	T	T	T	T	T	T	T	T	T	T	T	Syn.	Fixed
1584	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1590	C	T	T	T	T	T	T	T	T	T	T	T	T	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1596	G	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1611	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	T	T	T	T	T	T	T	T	T	T	T	T	Syn.	Fixed
1614	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1630	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.
1637	A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Fixed
1657	C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	Syn.	Poly.

Inter-specific codon models

Yang (1998)

- Consider at least 2 species, with one sequence per species

Inter-specific codon models

Yang (1998)

- Consider at least 2 species, with one sequence per species
- Assumes a known phylogeny (at least topology)

Inter-specific codon models

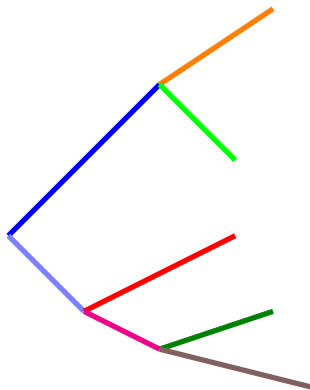
Yang (1998)

- Consider at least 2 species, with one sequence per species
- Assumes a known phylogeny (at least topology)
- Non-homogeneous model: distinct branches in the tree are allowed to have evolved with distinct $\omega = dN/dS$:

Inter-specific codon models

Yang (1998)

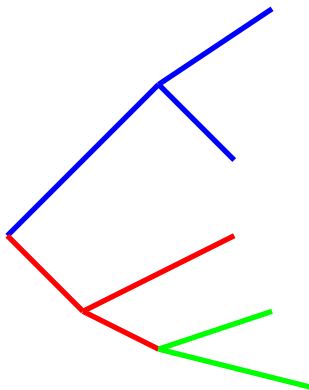
- Consider at least 2 species, with one sequence per species
- Assumes a known phylogeny (at least topology)
- Non-homogeneous model: distinct branches in the tree are allowed to have evolved with distinct $\omega = dN/dS$:
 - One per branch \Rightarrow *branch* model



Inter-specific codon models

Yang (1998)

- Consider at least 2 species, with one sequence per species
- Assumes a known phylogeny (at least topology)
- Non-homogeneous model: distinct branches in the tree are allowed to have evolved with distinct $\omega = dN/dS$:
 - One per branch \Rightarrow *branch* model
 - Several clades \Rightarrow *clade* model



Inter-specific codon models

Yang (1998)

- Consider at least 2 species, with one sequence per species
- Assumes a known phylogeny (at least topology)
- Non-homogeneous model: distinct branches in the tree are allowed to have evolved with distinct $\omega = dN/dS$:
 - One per branch \Rightarrow *branch* model
 - Several clades \Rightarrow *clade* model
- Other parameters are constant throughout the tree

Finding the best model

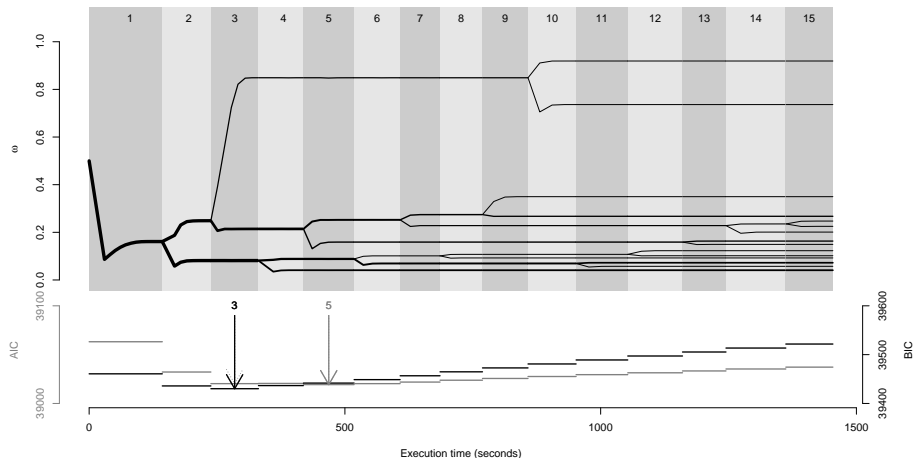
Dutheil et al. (2012)

- The branch model suffers from overparametrization issues
- The clade model needs an *a priori* knowledge

Finding the best model

Dutheil et al. (2012)

- The branch model suffers from overparametrization issues
- The clade model needs an *a priori* knowledge



Combining site and branch heterogeneity

Yang and Nielsen (2002), Zhang, Nielsen and Yang (2005)

- Consider a dataset with several species, with one species per branch and known phylogeny

Combining site and branch heterogeneity

Yang and Nielsen (2002), Zhang, Nielsen and Yang (2005)

- Consider a dataset with several species, with one species per branch and known phylogeny
- Consider two models, with and without selection. Branches where positive selection might have occurred are known *a priori*

Combining site and branch heterogeneity

Yang and Nielsen (2002), Zhang, Nielsen and Yang (2005)

- Consider a dataset with several species, with one species per branch and known phylogeny
- Consider two models, with and without selection. Branches where positive selection might have occurred are known *a priori*
- Branches evolving under positive selection are called *foreground* branches, others *background* branches

Combining site and branch heterogeneity

Yang and Nielsen (2002), Zhang, Nielsen and Yang (2005)

- Consider a dataset with several species, with one species per branch and known phylogeny
- Consider two models, with and without selection. Branches where positive selection might have occurred are known *a priori*
- Branches evolving under positive selection are called *foreground* branches, others *background* branches
- Background branches evolve under the M1a model, foreground branches under the M2a model

Combining site and branch heterogeneity

Yang and Nielsen (2002), Zhang, Nielsen and Yang (2005)

- Consider a dataset with several species, with one species per branch and known phylogeny
- Consider two models, with and without selection. Branches where positive selection might have occurred are known *a priori*
- Branches evolving under positive selection are called *foreground* branches, others *background* branches
- Background branches evolve under the M1a model, foreground branches under the M2a model
- Likelihood ratio test to compare with a homogeneous M1a model.