

Phylogenetics (Recap)

Simon J. Greenhill



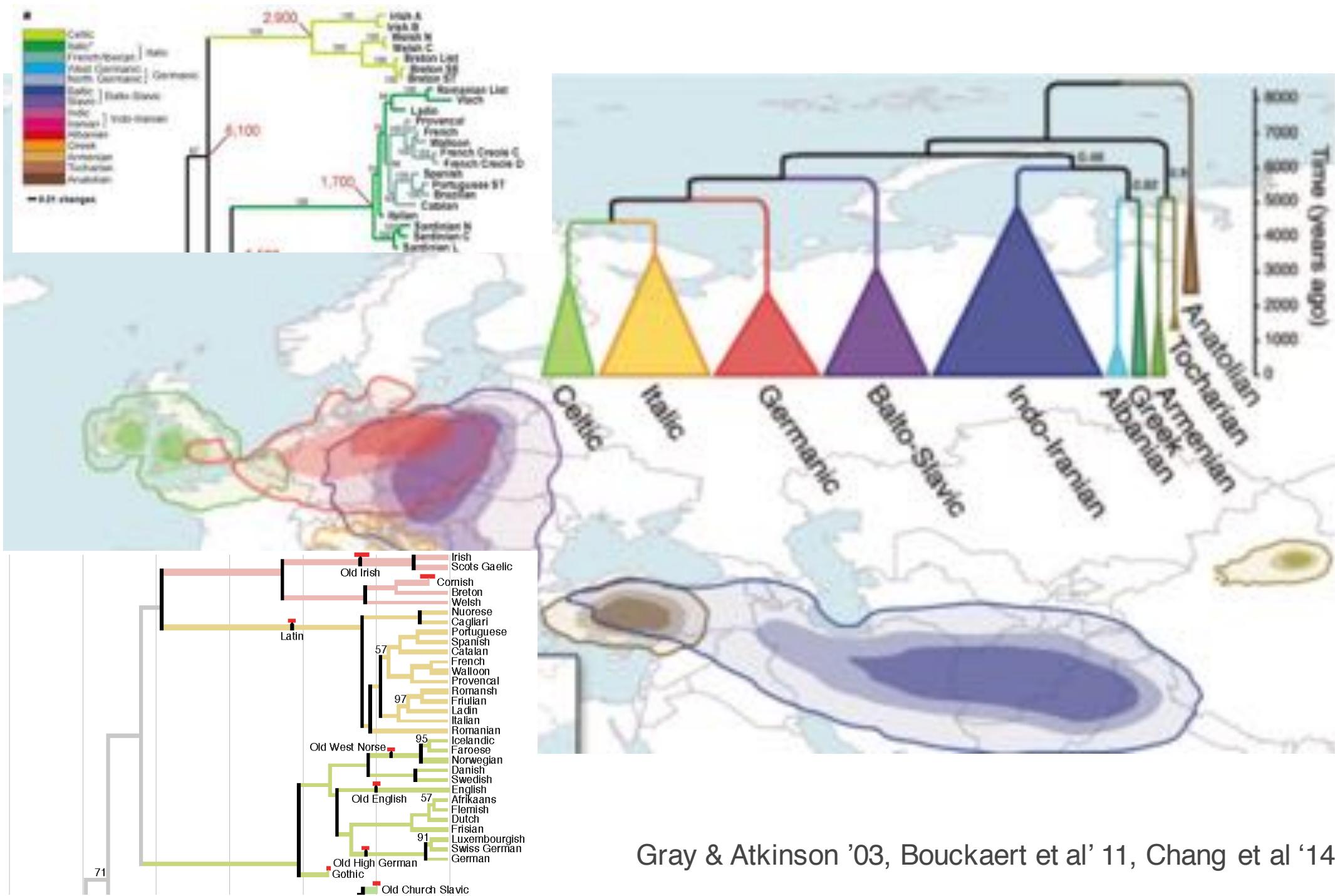
ARC CENTRE OF EXCELLENCE FOR
THE DYNAMICS OF LANGUAGE

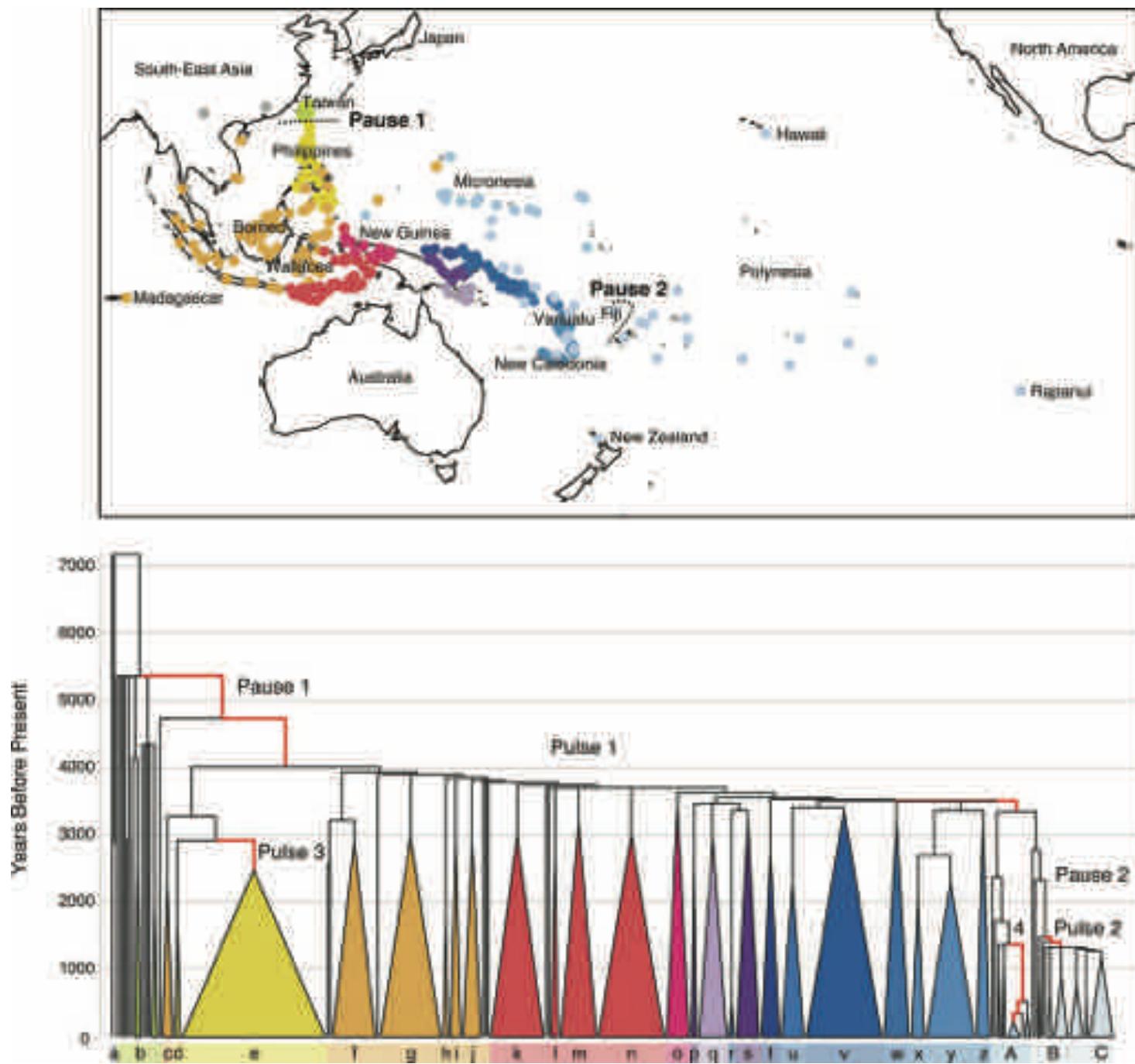


Max Planck Institute for the
Science of Human History

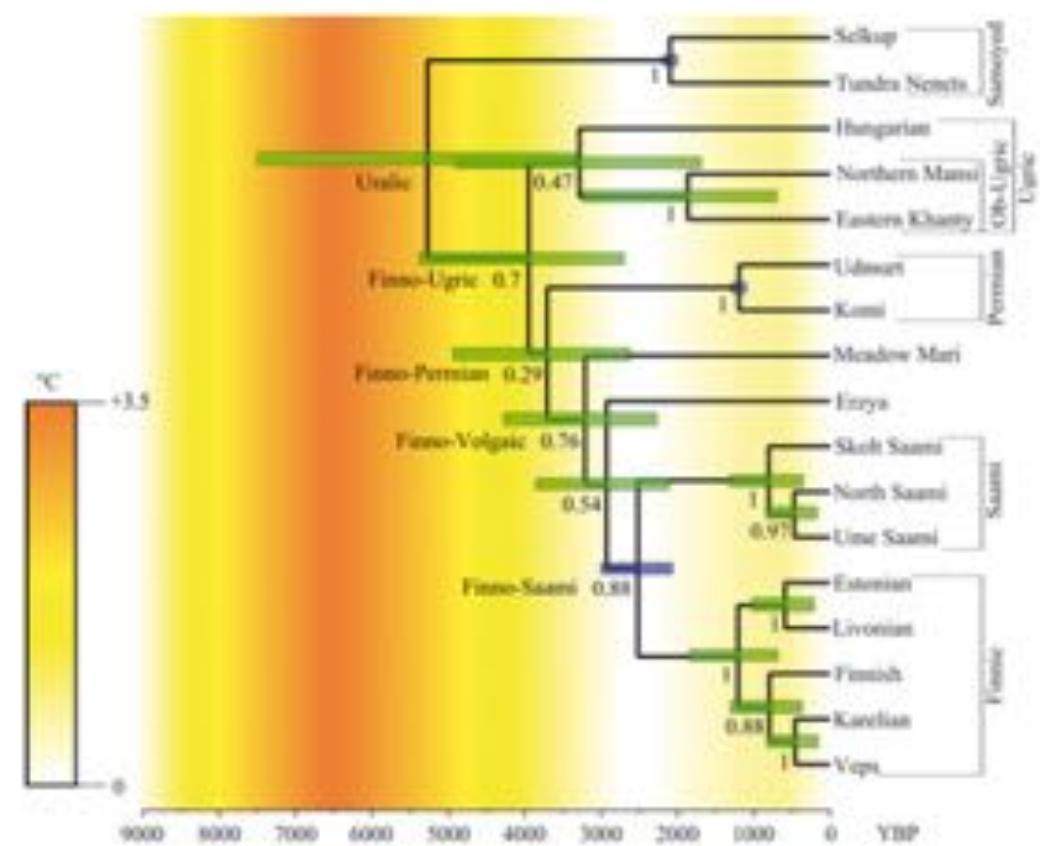
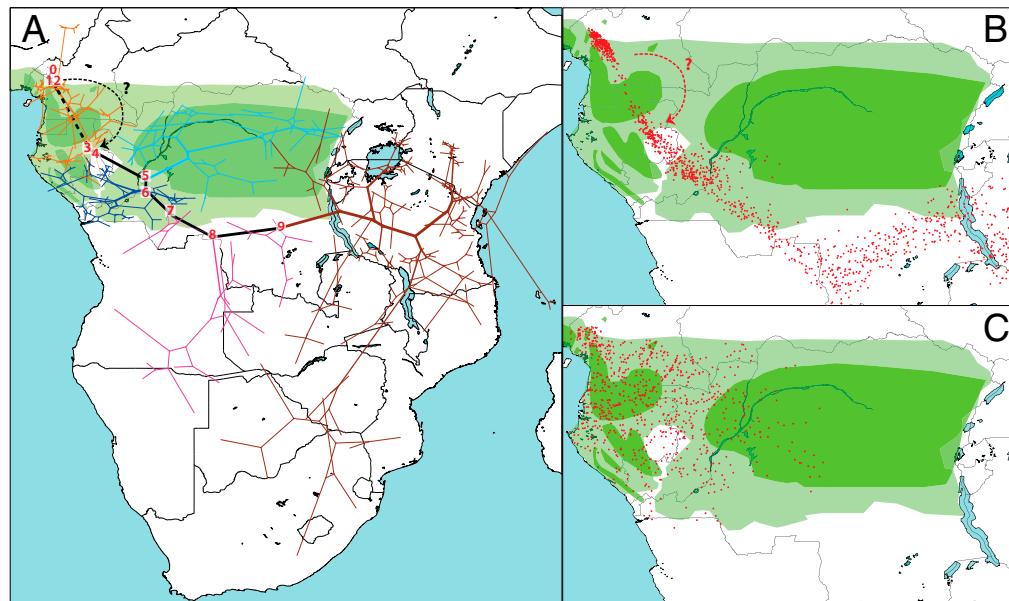
Phylogenetics

Range of methods in a robust, statistical/inferential framework for **testing** evolutionary hypotheses.

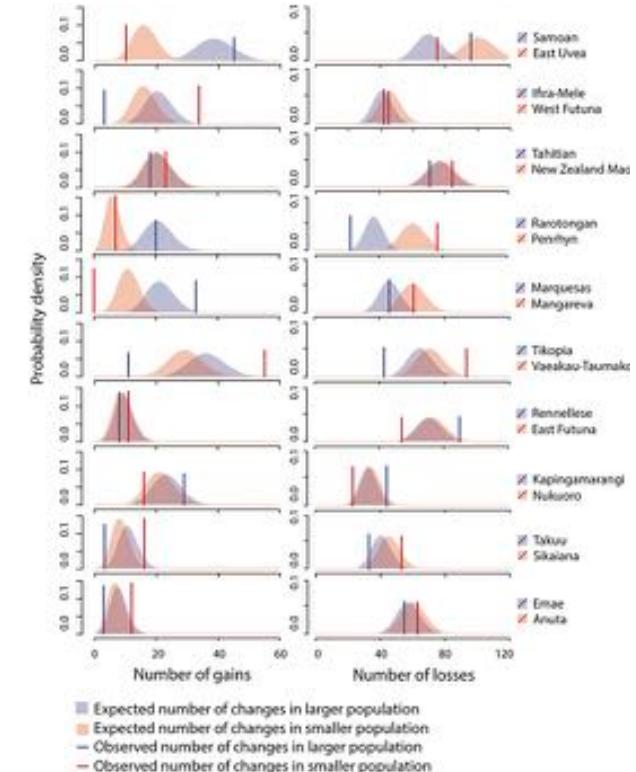
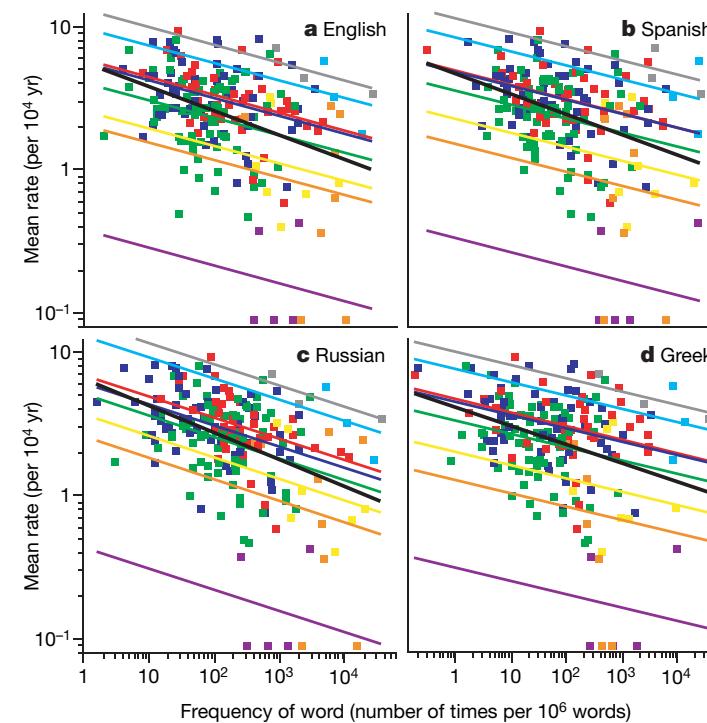
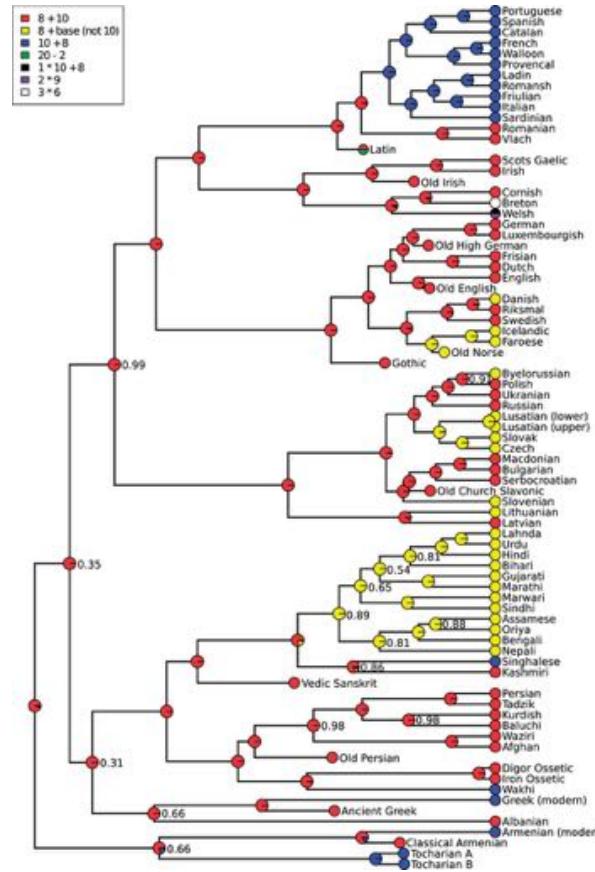
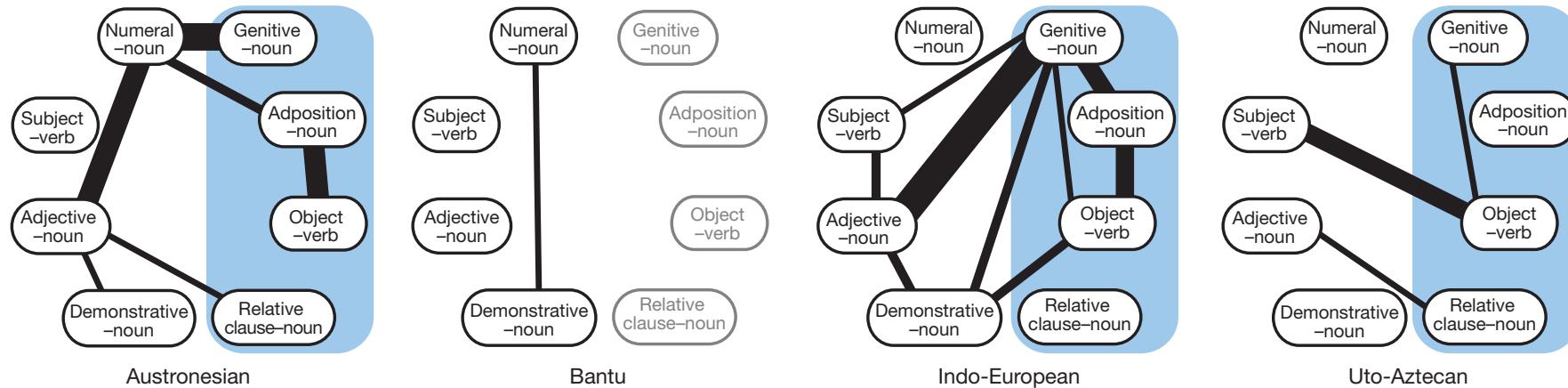




Gray et al '09



Dunn et al. '15

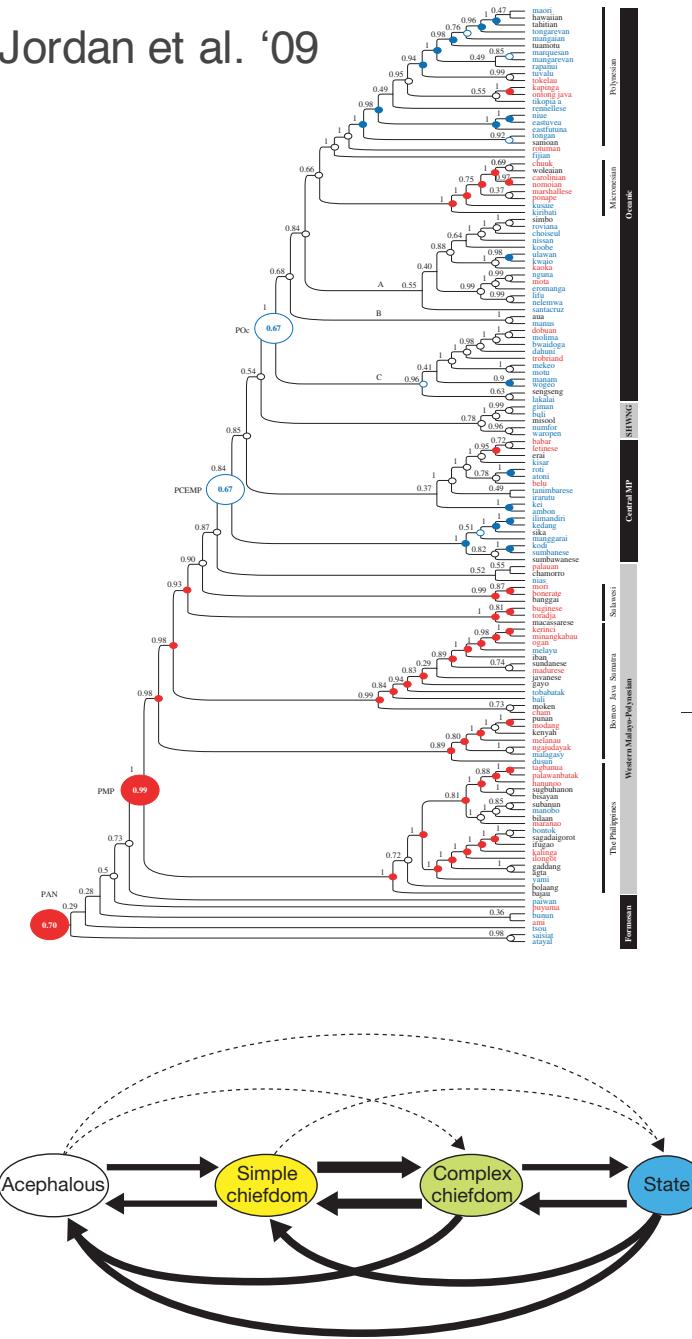


Calude and Verkerk '16

Pagel et al '07

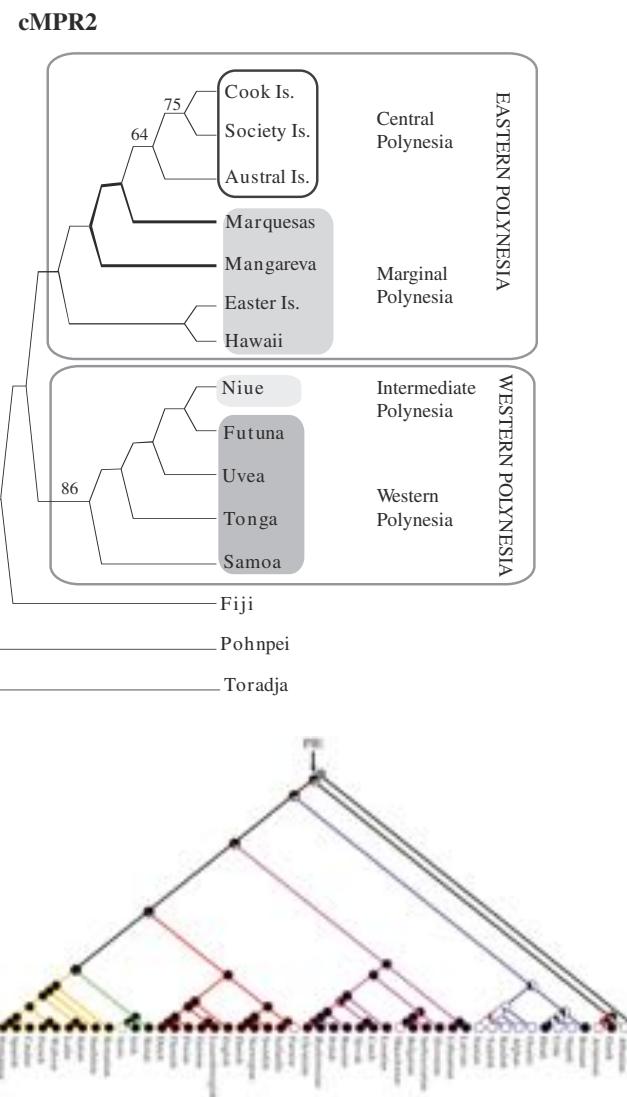
Bromham et al. '15

Jordan et al. '09



Currie et al. '10

Larsen '11



Da Silva & Tehrani '16

Matthews et al '11

Pazyryk

Yomut

Shahsevan

Qashqai

Boyer Ahmad

Bakhtiari

A

1.00
(1.00)

0.98
(0.85)

0.99
(0.91)

0.99
(0.95)

Plain-weave

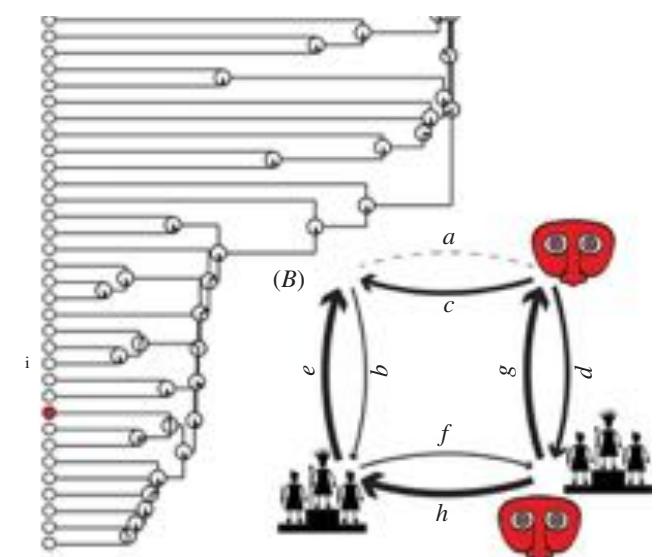
Weft-wrapping

Pile-weave technique

"Infinite knot" motif

"Rooster" motif

"Jagged border"



Watts et al. '15

Controversial

“most vibrant stream
of contemporary
linguistics”

“Computational methodologies
of this kind can only be helpful
for historical linguistics”

“languages and biomolecular
sequences evolve in very
different ways”

“more questions than
answers”

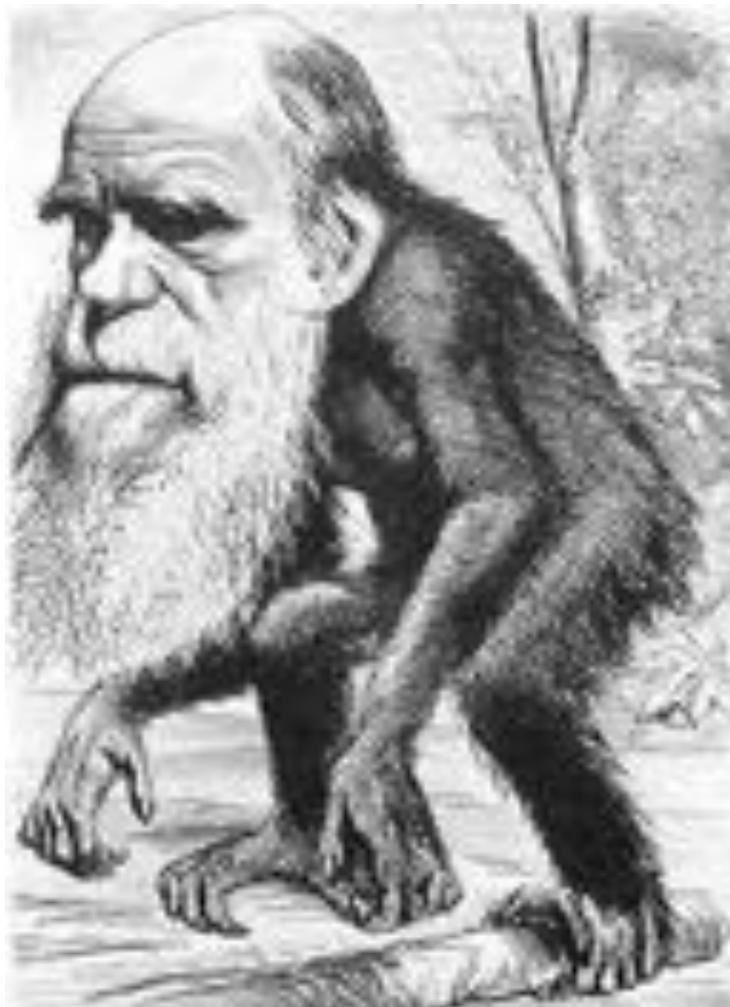
“utter bollocks”

“biggest intellectual fraud
since Chomsky”

“this isn’t history, it’s history put in nested boxes!”



What is evolution?



Variation

Heritability

Differential survival

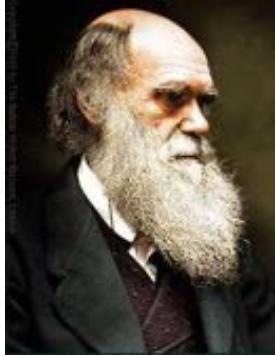
When and where did  originate?

What **differences** are there between 's?

How are  related to other 's?

What **processes** shaped ?

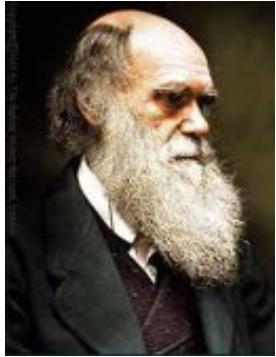
Can we infer what  were in the **past**?



Darwin (1871)

"Languages, like organic beings, can be classed in groups under groups..."

"The formation of different languages and of distinct species, and the proofs that both have been developed through a gradual process, are **curiously parallel**"



Darwin (1871)

"Languages, like organic beings, can be classed in groups under groups..."

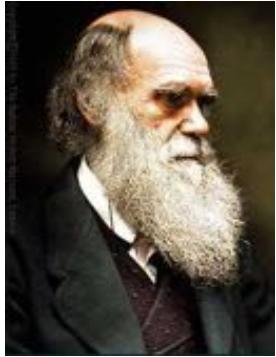
"The formation of different languages and of distinct species, and the proofs that both have been developed through a gradual process, are **curiously parallel**"



Schleicher 1863

Darwinism Tested by the Science of Language

"same process has long been generally assumed for linguistic organisms"



Darwin (1871)

"Languages, like organic beings, can be classed in groups under groups..."

"The formation of different languages and of distinct species, and the proofs that both have been developed through a gradual process, are **curiously parallel**"



Schleicher 1863

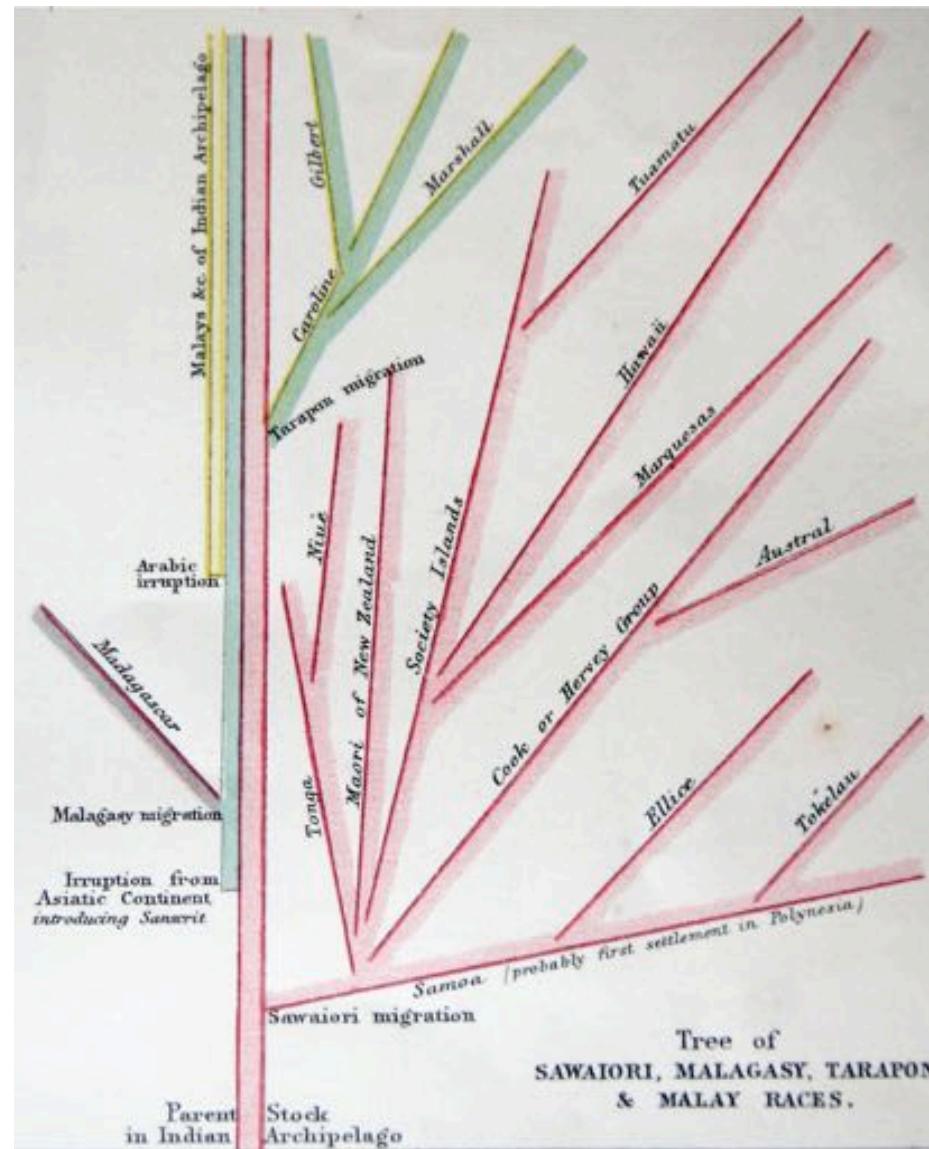
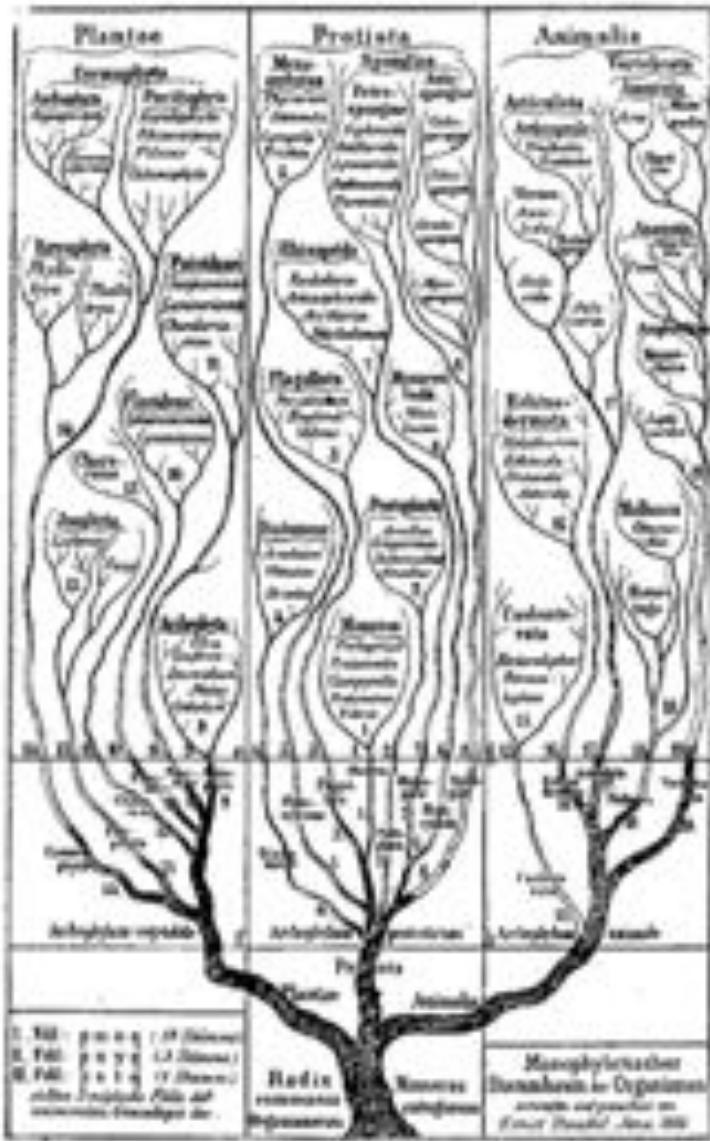
Darwinism Tested by the Science of Language

"same process has long been generally assumed for linguistic organisms"

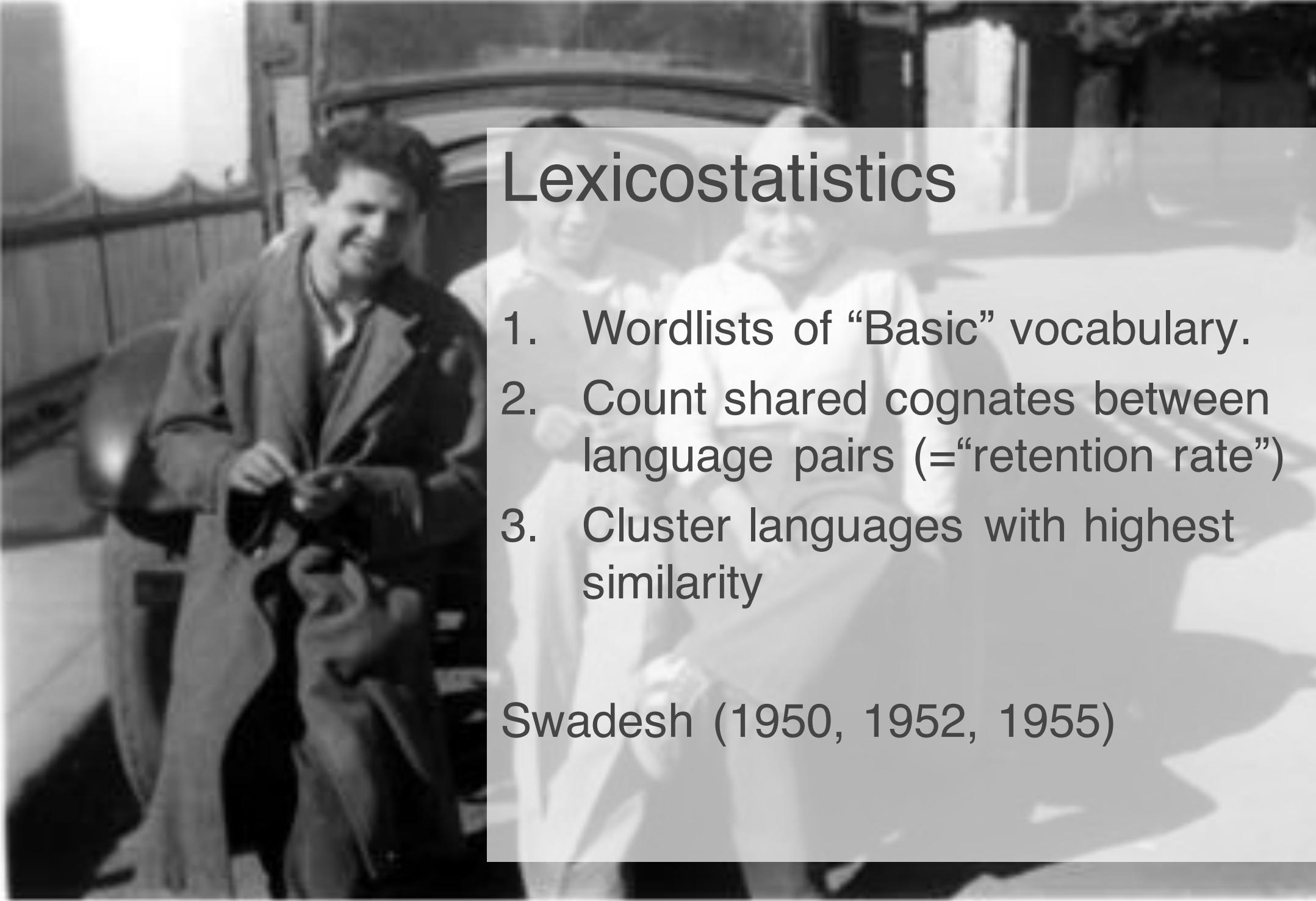


Brugmann (1884)

Importance of using "shared innovations" to identify clades and not "shared retentions"







Lexicostatistics

1. Wordlists of “Basic” vocabulary.
2. Count shared cognates between language pairs (=“retention rate”)
3. Cluster languages with highest similarity

Swadesh (1950, 1952, 1955)

	Taboo	Blood	To Suck
Fijian	tabu	drā	sucu-ma
Tahitian	tapu	toto	ngote
Maori	tapu	toto	ngote
Hawaiian	kapu	koko	omo
Marquesan	tapu	toto	omo

Identified by Systematic Sound Correspondences
 - e.g. Maori “t” = “k” in Hawaiian.

Elbert 1953

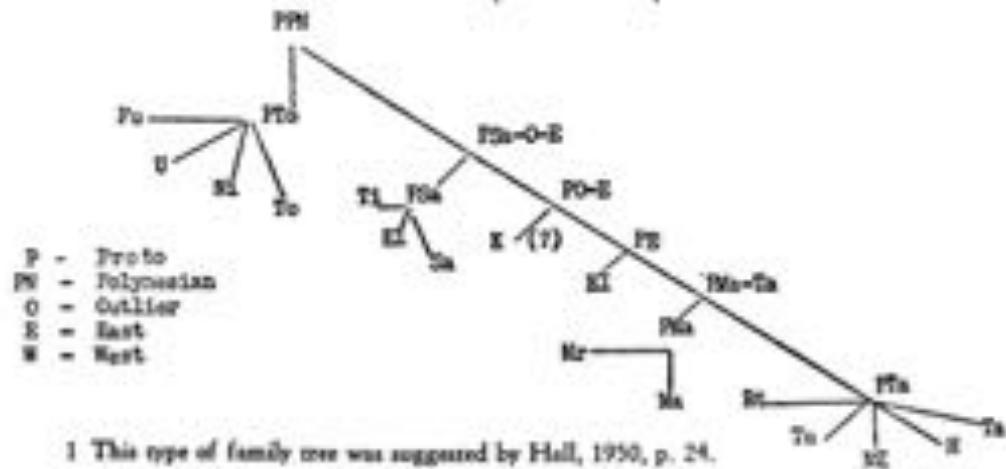
TABLE 2
Polynesian cognate percentages

<u>P</u>	<u>Si</u>	<u>Ta</u>	<u>Ti</u>	<u>S</u>	<u>Ma</u>	<u>Si¹</u>	<u>Ti¹</u>	<u>Ma¹</u>	<u>K¹</u>	<u>E</u>	<u>II</u>	<u>Mr</u>	<u>Ma</u>	<u>Si</u>	<u>Ta</u>	<u>EI</u>	<u>H</u>	<u>Ta</u>
0	62	76	03	79	76	57	62	68	56	52	50	57	56	66	62	61	58	58
	72	66	76	76	70	53	59	50	52	51	53	53	51	62	62	61	55	59
	64	68	61	63	50	59	48	58	49	45	49	55	47	56	55	58	49	51
	70	68	66	66	46	55	55	49	45	48	49	45	45	58	53	54	49	52
	82	76	66	66	68	59	59	62	67	63	71	68	71	67	66	71		
			66	62	52	59	58	61	62	63	66	66	67	68	66	66		
			64	60	62	55	53	53	55	52	67	62	57	59	60	66		
	55	60	52	55	52	55	56	56	60	58	60	60	60	59	59	51		
			57	52	58	58	59	53	55	61	60	56	56	51	51	51		
				51	54	53	53	53	51	53	56	55	56	53	53	52		
				53	49	45	45	46	53	50	52	49	49	50	51	49	48	
					47	49	45	45	54	51	51	51	49	50	50	50	48	
						64	63	64	62	63	64	64	62	64	64	64	64	64
							73	75	72	70	69	68	73	73	69	67	68	68
								79	77	77	70	67	67	79	77	73	73	73
								83	83	79	85	83	83	83	83	71	73	76
									79	77	83	79	77	71	73	76	76	76
										83	83	85	83	83	83	83	83	83
											73	75	72	70	69	67	68	68
												73	75	72	70	69	67	68
													73	75	72	70	69	68
														73	75	72	70	69
															73	75	72	70
																73	75	72
																	73	75
																		73

1 Percentage based on incomplete data.

TABLE 4

A tentative family tree for Polynesia¹



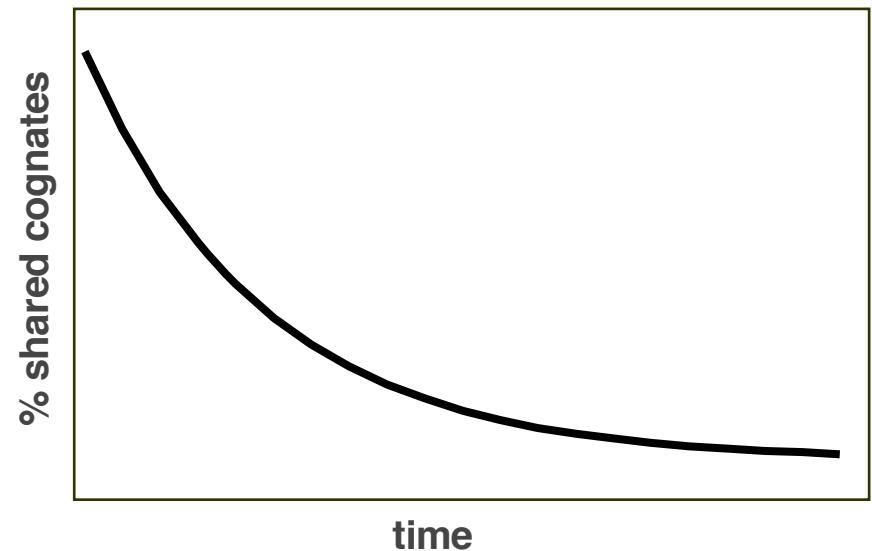


Glottochronology

Loss of cognates happens at a constant rate
(=radioactive decay)

19% loss per 1000 years (Lees 1953)

$$time = \frac{\log(\% \text{ shared cognates})}{2 \log(\text{retention rate})}$$



The Rise of Lexicostatistics...

IN THE LAST DECADE glottochronology has excited international interest and acquired a literature of its own. To the anthropologist it promises a measure of time depth for language families without documented history, and yet another linguistic example of regularity in cultural phenomena.

Hymes (1960): “Lexicostatistics so far”

“... a significant work—one which may conceivably be as revolutionary for Oceanic linguistics and culture history as was the work of Greenberg (1949–54) for the interpretation of African languages and cultures”

Murdock (1964) p.117

...and the fall of Lexicostatistics

Major Criticism: Universality of Rates

Old Norse & Icelandic?

- Glottochronology: 200 years.
- Reality: 1000 years

Bergsland & Vogt 1962: “Our findings clearly disprove the basic assumption of glottochronology ‘that fundamental vocabulary changes at a constant rate’ ”



Jungner, Hugo; Elisabeth Svärdström (1940-1971). Sveriges runinskrifter: V. Västergötlands runinskrifter. Stockholm: Kungl. Vitterhets Historie och Antikvitets Akademien. ISSN 0562-8016. p. 260

Fallout.

"a tradition of hostility towards probabilistic modelling in historical linguistics" (Sankoff '73)

"In summary, glottochronology is not accurate; all its basic assumptions have been severely criticized. It should not be accepted, it should be rejected" (Campbell '04)

"Linguists don't do dates" (McMahon & McMahon '03)





U.P.G.M.A

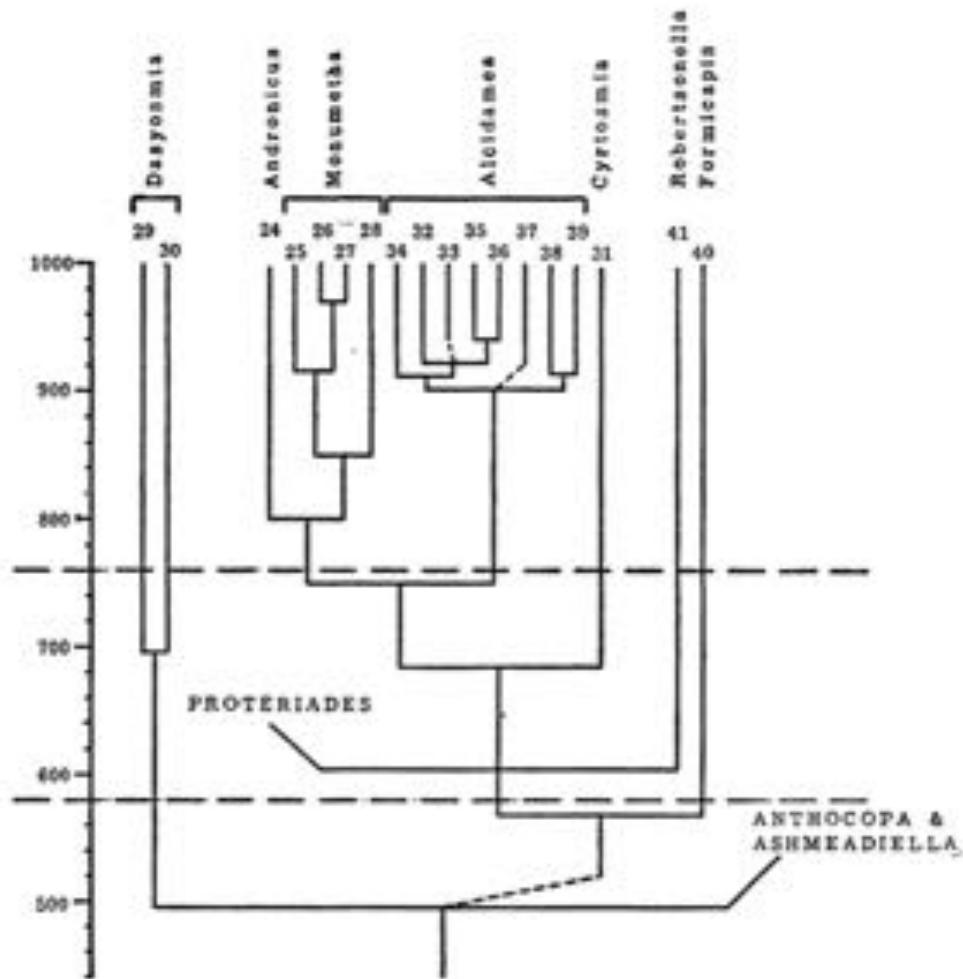


FIG. 6. Diagram of relationships for the genus *Hoplitis* obtained by the weighted variable group method.

A QUANTITATIVE APPROACH TO A PROBLEM
IN CLASSIFICATION¹

CHARLES D. MICHENER AND ROBERT R. SOKAL²

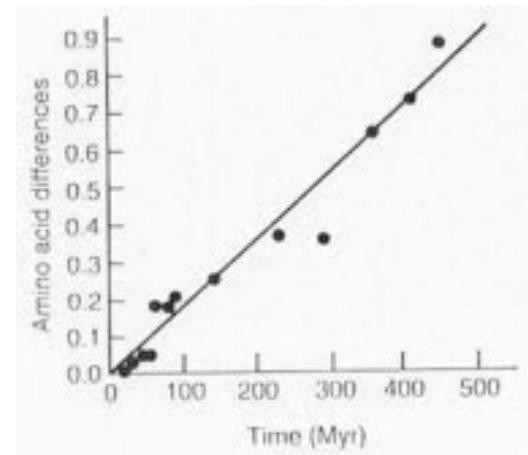
Department of Entomology, University of Kansas, Lawrence



Molecular Clock

Zuckerkandl & Pauling 1962

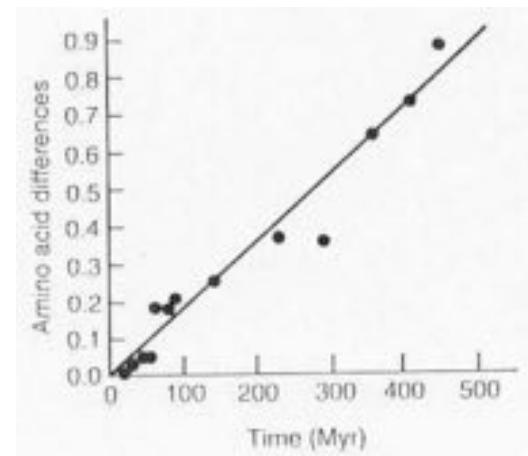
Number of AA differences were proportional to species divergence times.



Molecular Clock

Zuckerkandl & Pauling 1962

Number of AA differences were proportional to species divergence times.



Kimura 1968

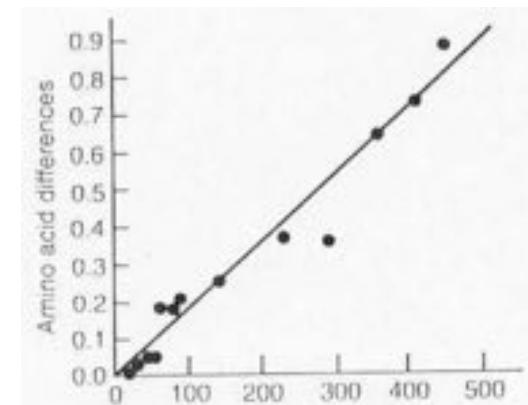
The average time taken for one base pair replacement within a genome is therefore

$$28 \times 10^6 \text{ yr} \div \left(\frac{4 \times 10^9}{300} \right) \div 1.2 \div 1.8 \text{ yr}$$

Molecular Clock

Zuckerkandl & Pauling 1962

Number of AA differences were proportional to species divergence times.



Kimura 1968

The average time taken for one base pair replacement within a genome is therefore

$$28 \times 10^6 \text{ yr} \div \left(\frac{4 \times 10^9}{300} \right) \div 1.2 \div 1.8 \text{ yr}$$

Sarich & Wilson 1967

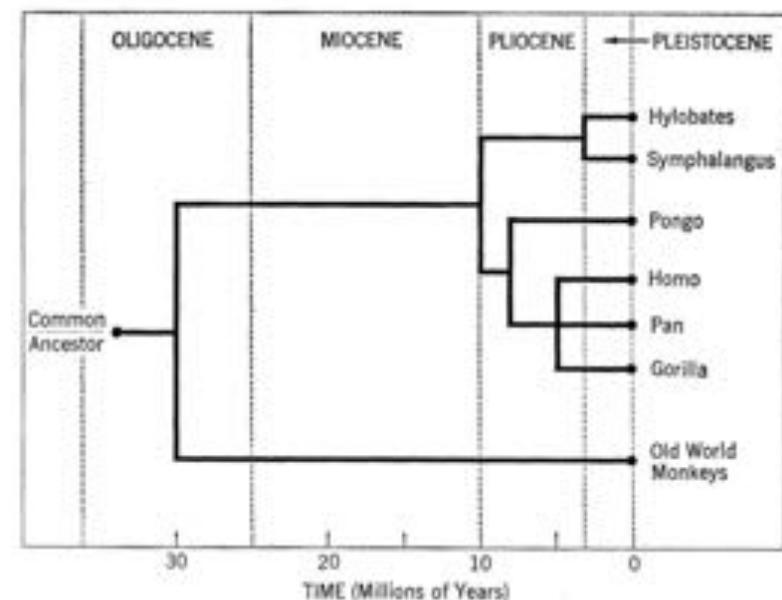


Fig. 1. Times of divergence between the various hominoids, as estimated from immunological data. The time of divergence of hominoids and Old World monkeys is assumed to be 30 million years.

Problems?

Kirsch 1969

SEROLOGICAL DATA AND PHYLOGENETIC INFERENCE:
THE PROBLEM OF RATES OF CHANGE

JOHN A. W. KIRSCH

Felsenstein 1978

CASES IN WHICH PARSIMONY OR COMPATIBILITY
METHODS WILL BE POSITIVELY MISLEADING¹

JOSEPH FELSENSTEIN

Britten 1986

Rates of DNA Sequence Evolution Differ
Between Taxonomic Groups

ROY J. BRITTEN

The Cladistics Wars



SCIENCE
as a
PROCESS



An Evolutionary Account
of the Social and Conceptual
Development of Science

DAVID L. HULL

Solutions.

Cavalli-Sforza &
Edwards 1967

Yang 1993

Sanderson 1997

Drummond et al. 2006

Phylogenetic Analysis Models and Estimation Procedures

L. L. CAVALLI-SFORZA AND A. W. F. EDWARDS*

Maximum-Likelihood Estimation of Phylogeny from DNA Sequences When Substitution Rates Differ over Sites¹

Ziheng Yang

A Nonparametric Approach to Estimating Divergence Times in the Absence of Rate Constancy

Michael J. Sanderson

Relaxed Phylogenetics and Dating with Confidence

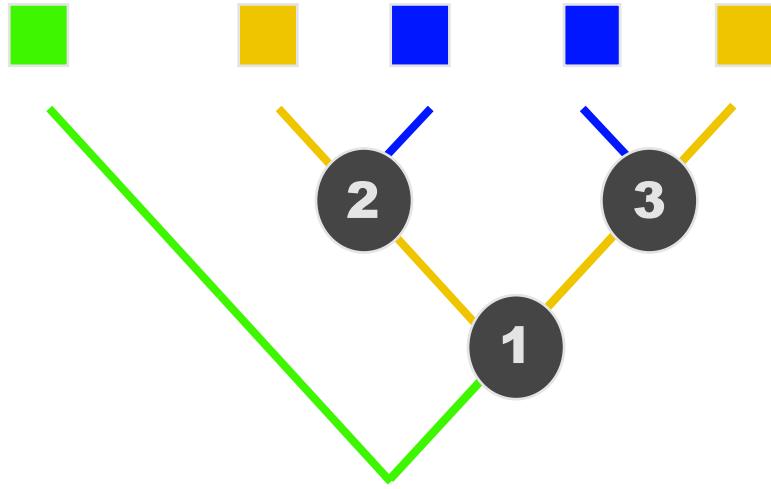
Alexei J. Drummond[✉], Simon Y. W. Ho, Matthew J. Phillips, Andrew Rambaut[✉]

How do we build trees?

1. Maximum Parsimony.
2. Maximum Likelihood.
3. **Bayesian Phylogenetic Analyses.**

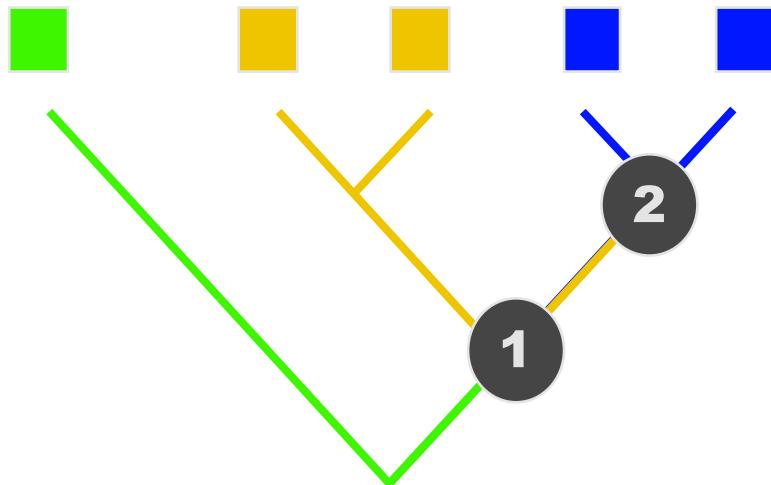


Maximum Parsimony



Aim to find “most parsimonious” tree

Smallest amount of evolution

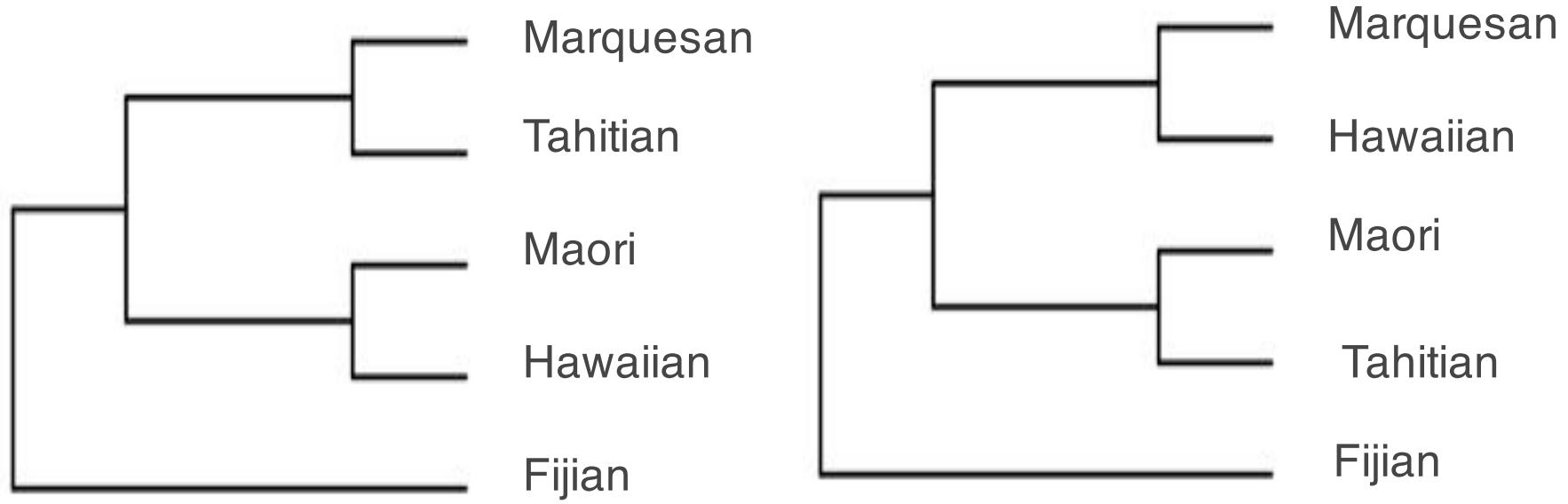


Unlikely that things should arise more than once
(i.e. cognates should innovate once)

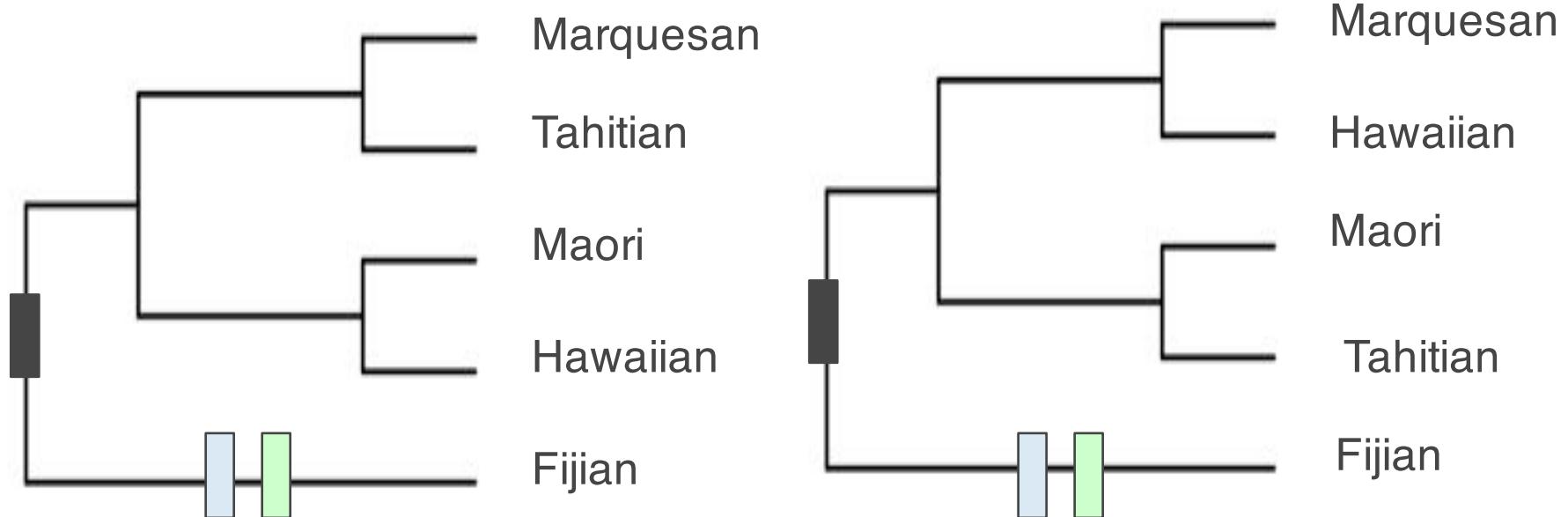
	Taboo	Blood	To Suck
Fijian	tabu	drā	sucu-ma
Tahitian	tapu	toto	ngote
Maori	tapu	toto	ngote
Hawaiian	kapu	koko	omo
Marquesan	tapu	toto	omo

	Taboo	Blood	To Suck
Fijian	tabu	drā	sucu-ma
Tahitian	tapu	toto	ngote
Maori	tapu	toto	ngote
Hawaiian	kapu	koko	omo
Marquesan	tapu	toto	omo

Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



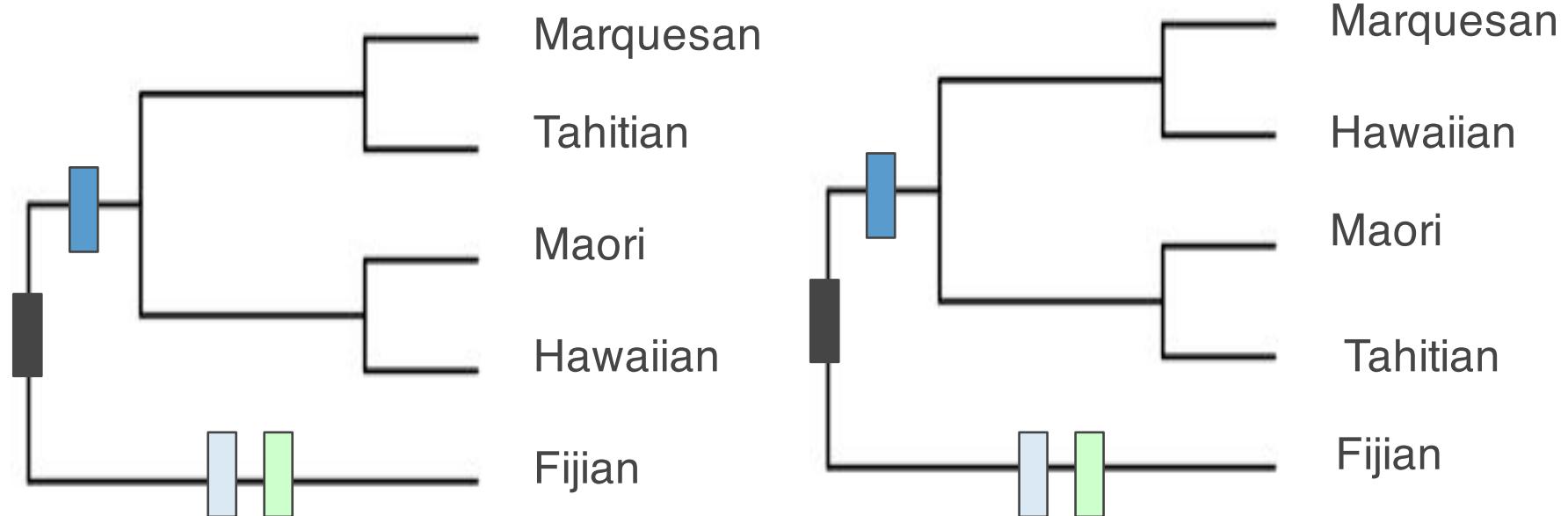
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=3

Length=3

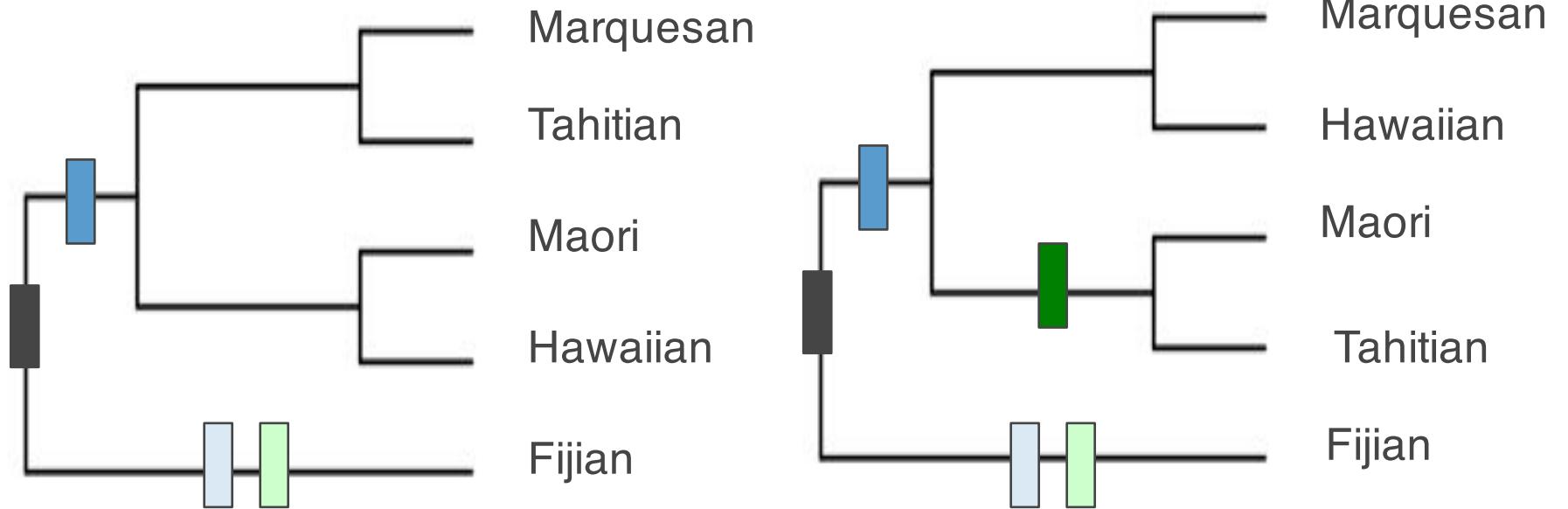
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=4

Length=4

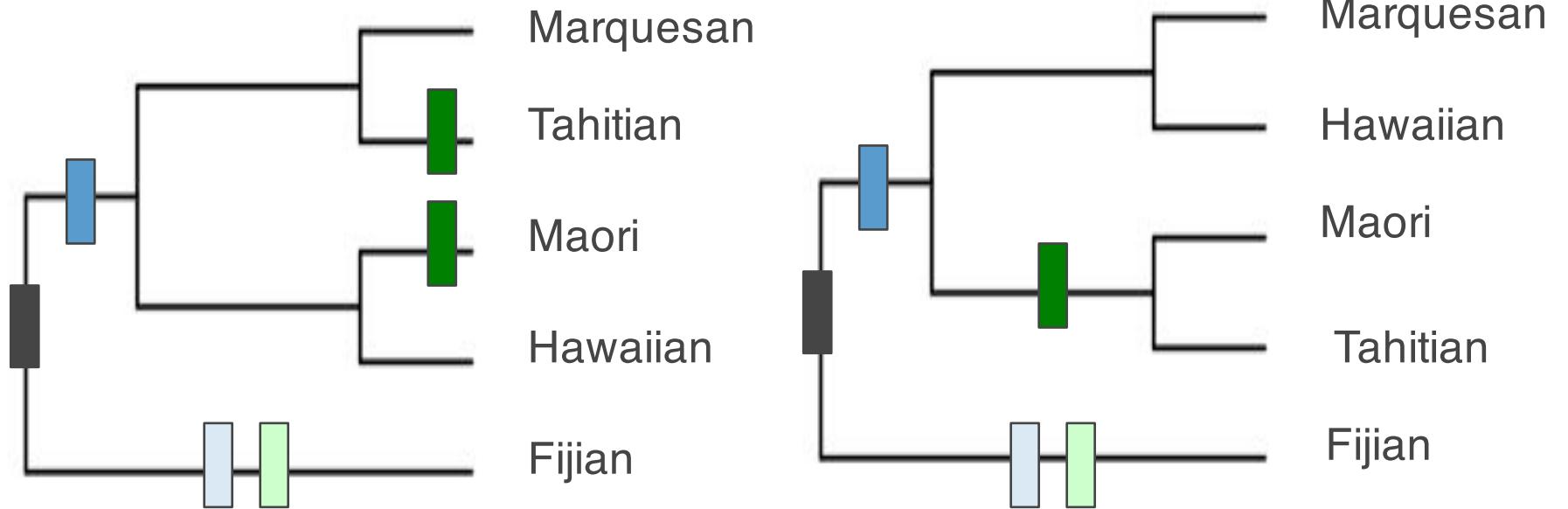
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=4

Length=5

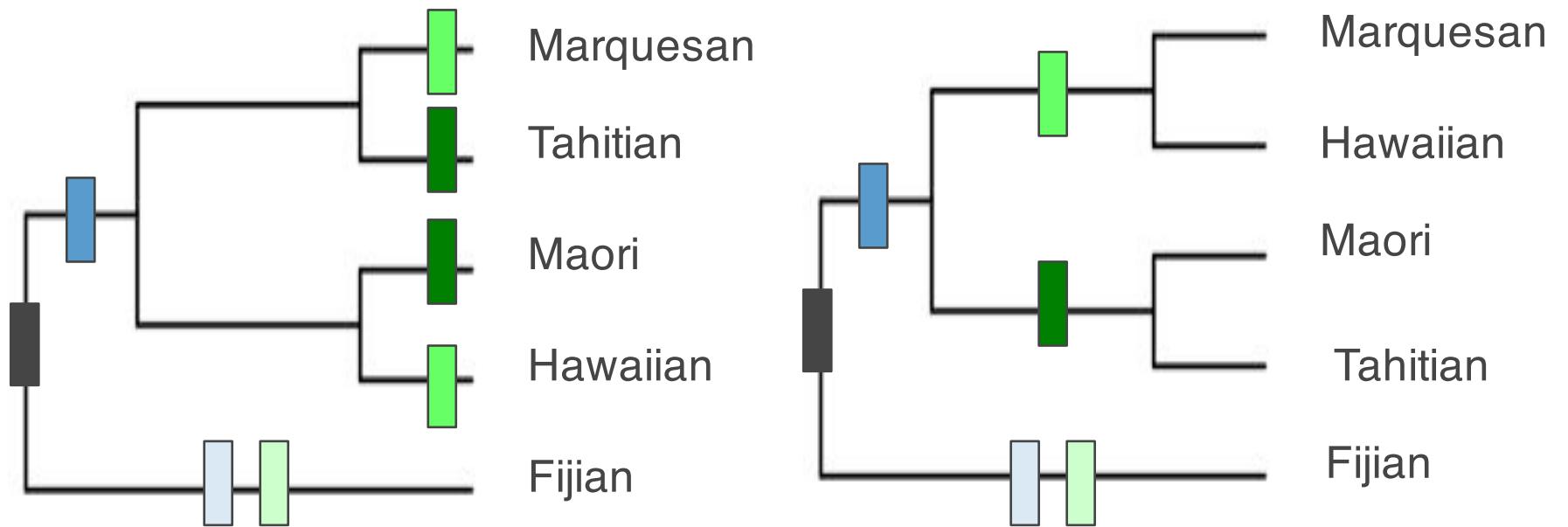
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=6

Length=5

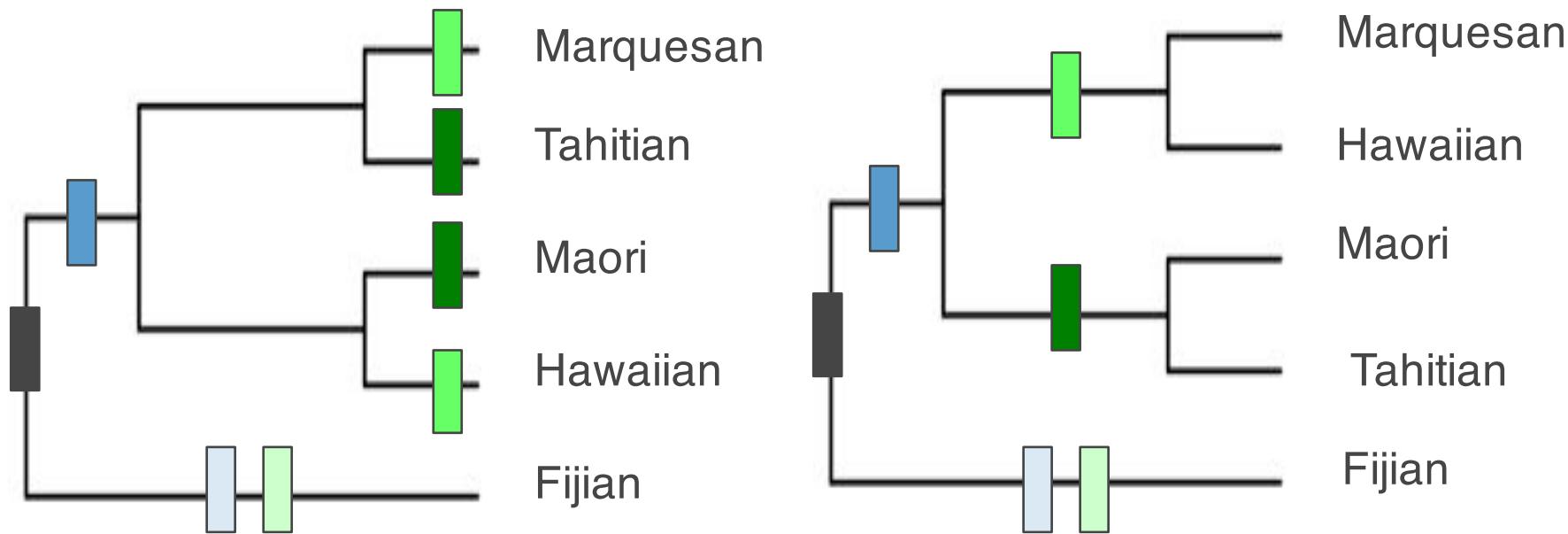
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=8

Length=6

Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



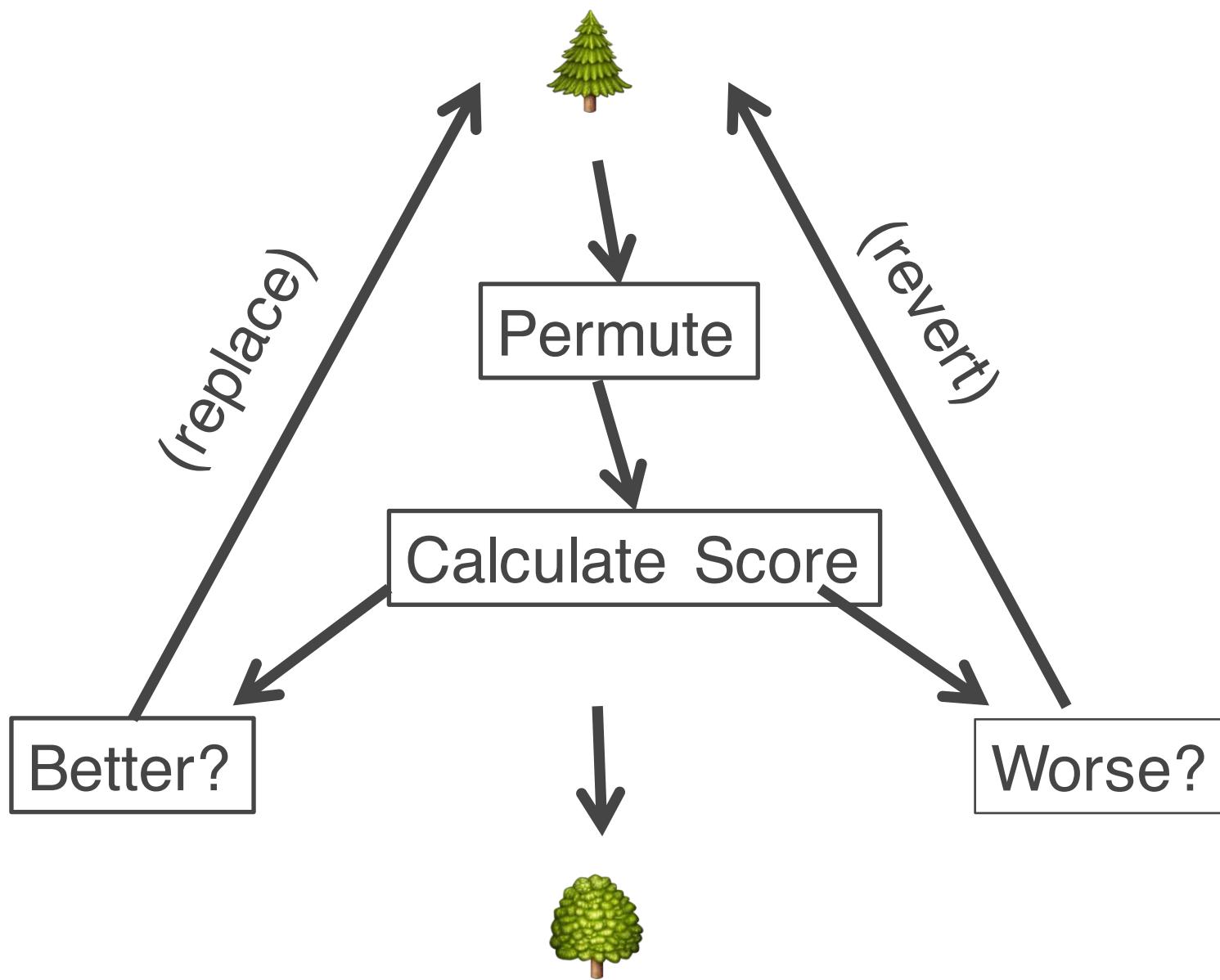
Length=8

Length=6



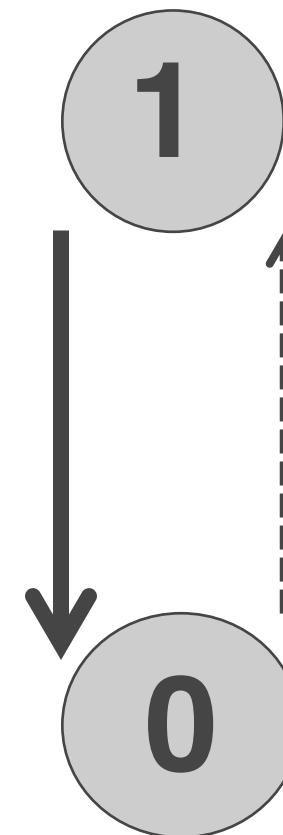
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1

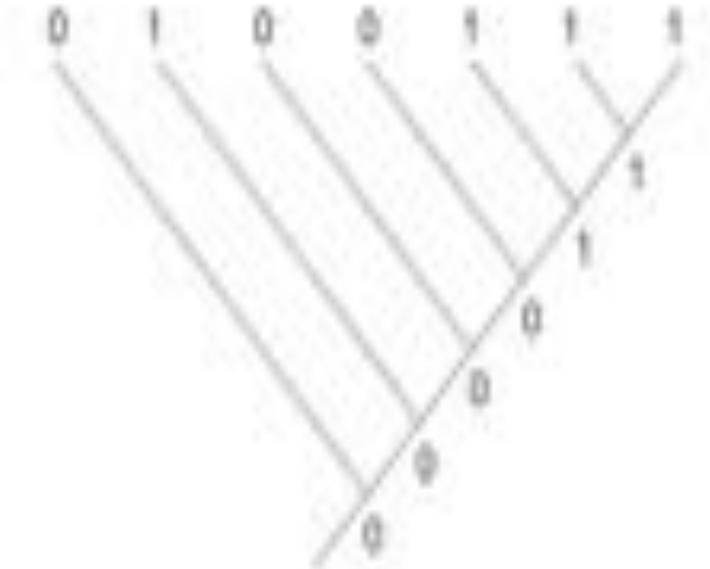
Algorithm



Maximum Likelihood

- Builds on Max. Parsimony
- Stochastic **model** of change.
- **Likelihood** = fit of data to tree under a model.
 - Very small number = $\log(L_h)$
 - Closer to zero = better fit.





$$L(a) = P(0 \rightarrow 0b1) \times P(0 \rightarrow 0b2) \times P(1 \rightarrow 1b3) \times P(1 \rightarrow 1b4) \times \\ P(0 \rightarrow 0b5) \times P(0 \rightarrow 0b6) \times P(0 \rightarrow 0b7) \times P(0 \rightarrow 1b8) \times P(1 \rightarrow 1b9) \\ \times P(1 \rightarrow 1b10) \times P(1 \rightarrow 1b11) \times P(1 \rightarrow 1b12)$$

$L_h =$

P of being in state 0, and staying state 0 on branch 1.

\times

P of being in state 0, and staying state 0 on branch 2.

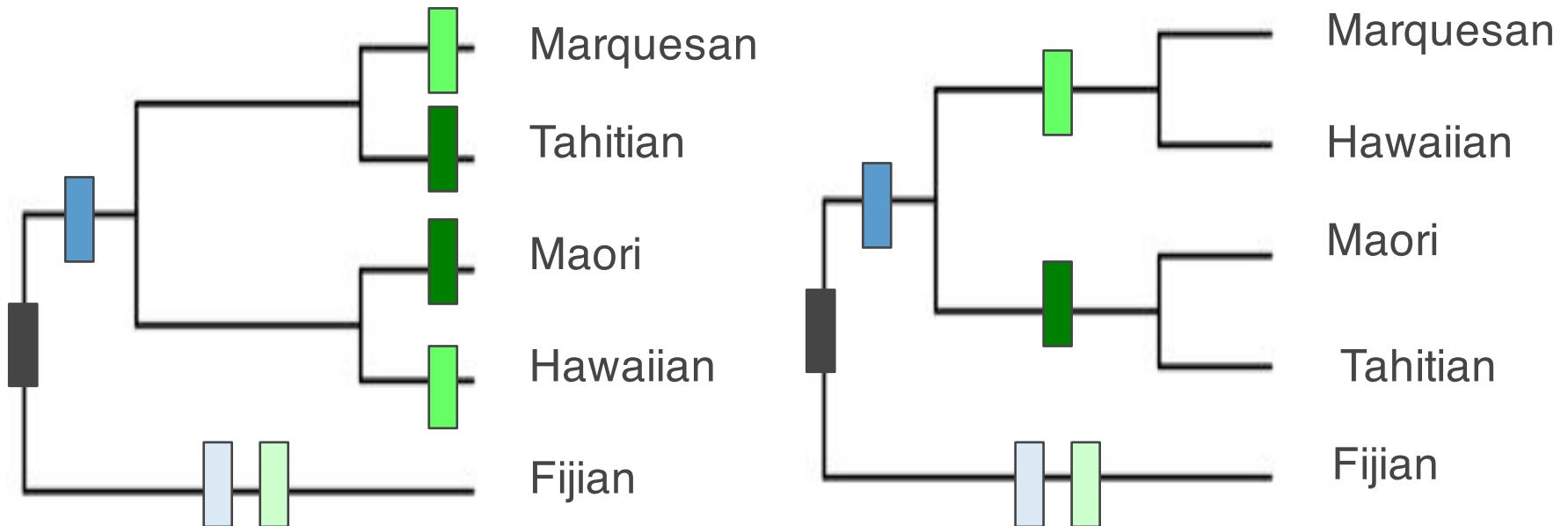
\times

P of being in state 0, and staying state 0 on branch 2.

.... etc

Site Likelihood(a) = $(\text{branch 1}) \times \dots \times (\text{branch n})$

Site Likelihood(a) = $P(\text{reconstruction 1}) \times \dots \times P(\text{reconstruction n})$

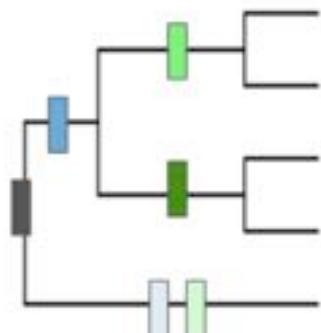
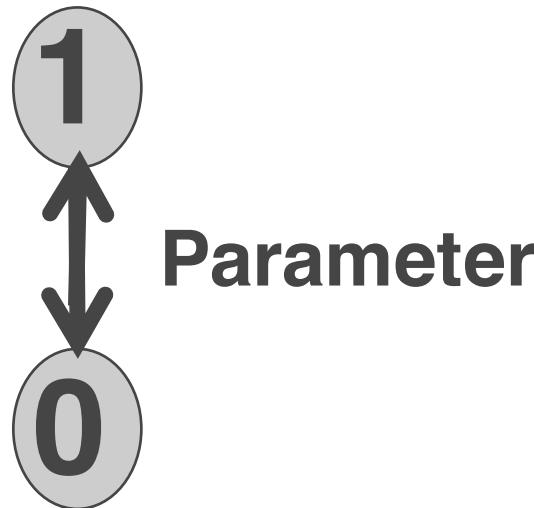


$$\ln(L) = -14.804$$

$$\ln(L) = -12.007 \quad \leftarrow$$

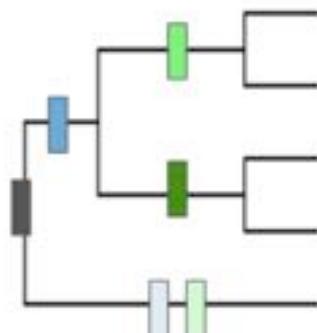
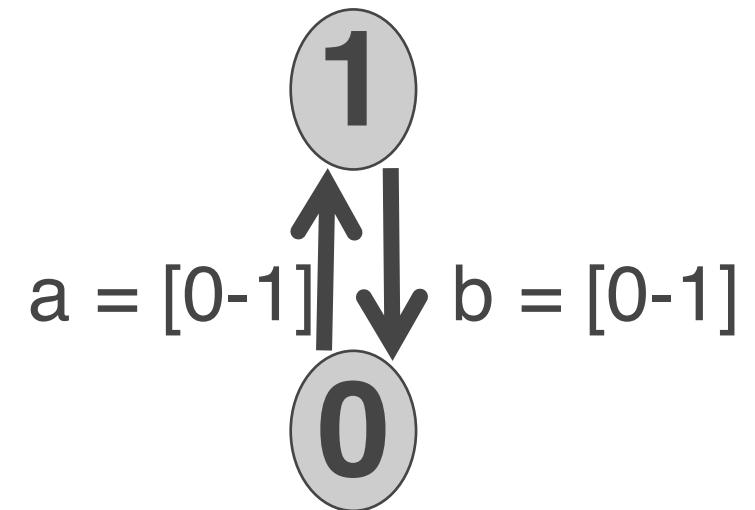
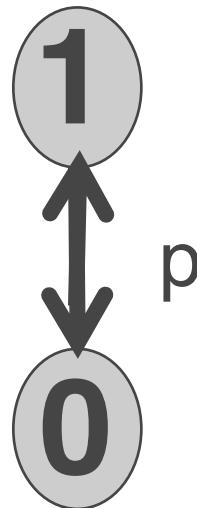
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1

Models



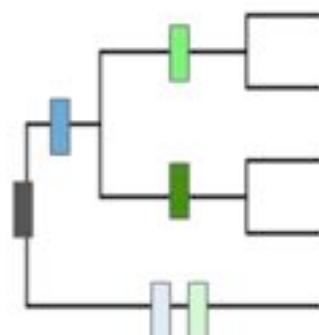
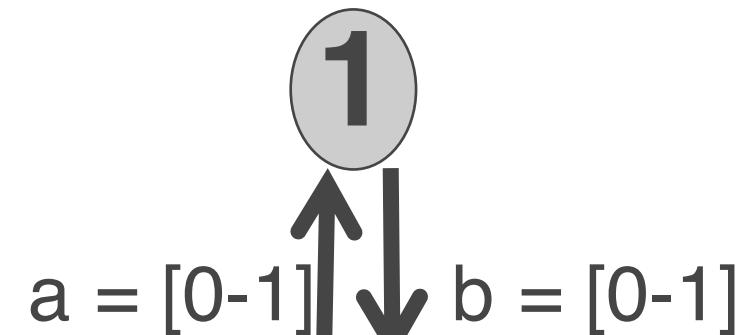
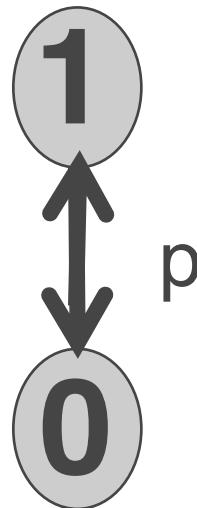
$$\ln(L) = -12.007$$

Models



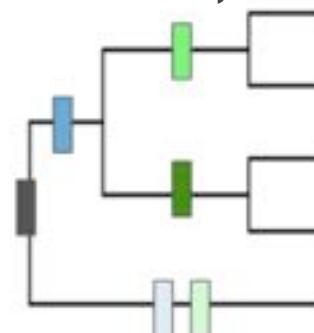
$$\ln(L) = -12.007$$

Models



$\ln(L) = -12.007$

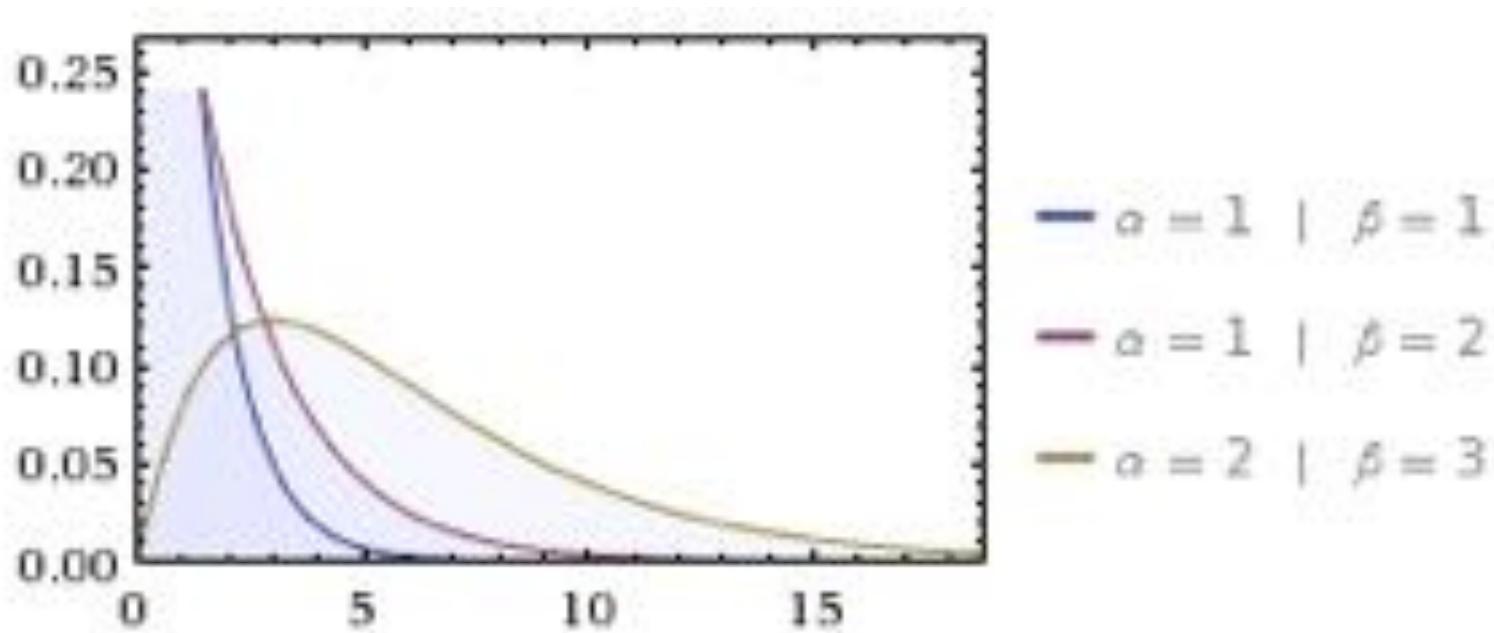
$a = 0.92, b = 0.08$

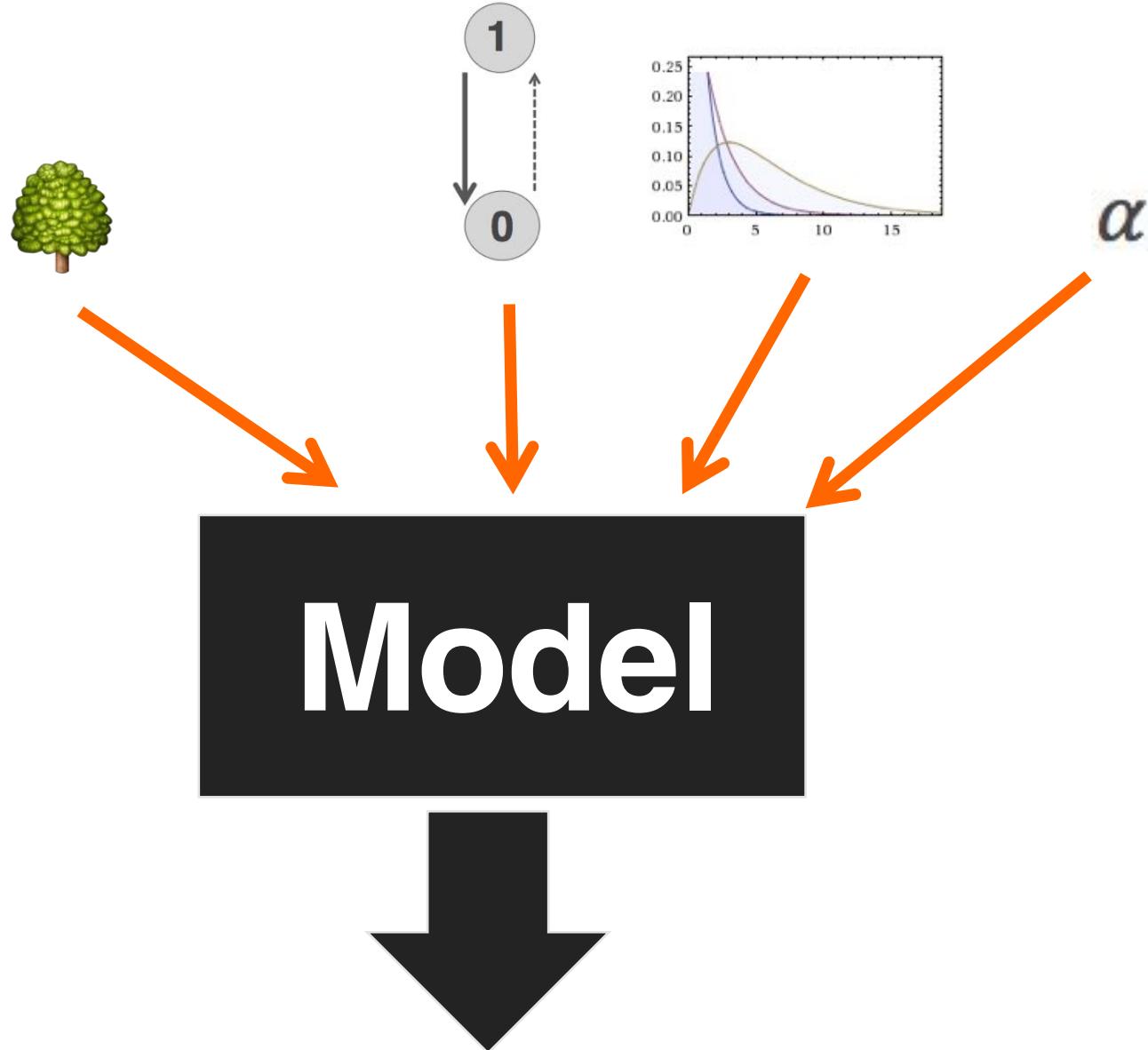


$\ln(L) = -9.072$

Rate Variation - Gamma

- Gamma Distribution
- One parameter, α , controls the shape
- Estimate the best value for α using Lh.





Likelihood score



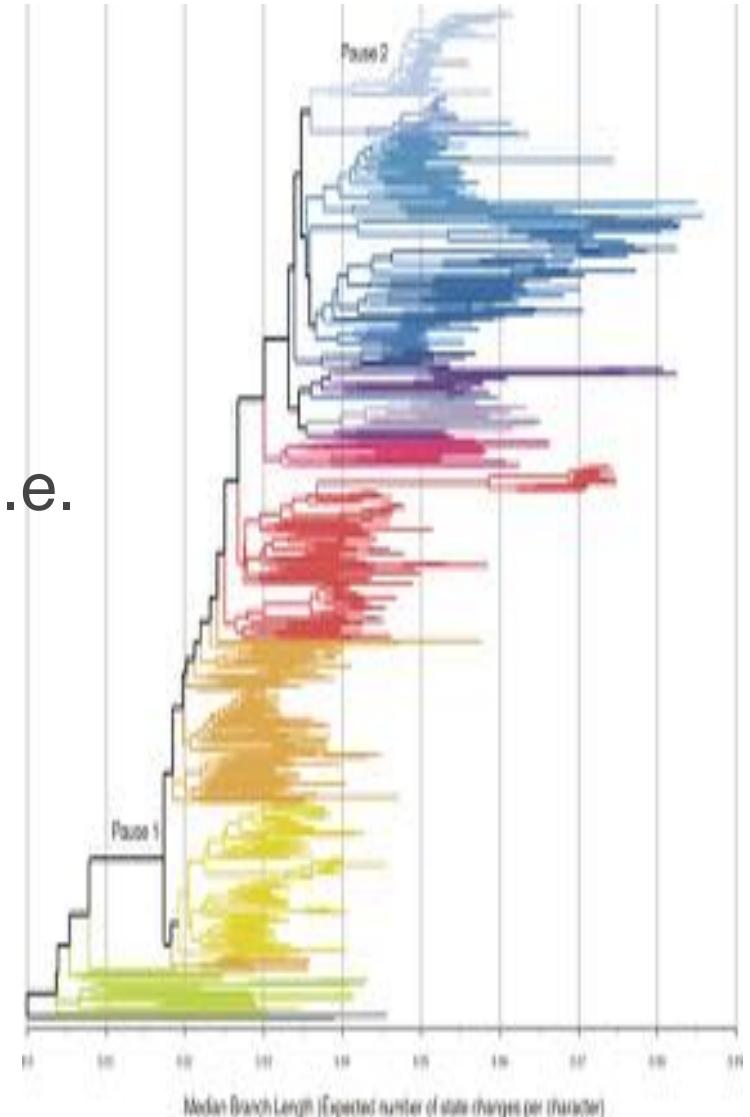
Linguists don't do dates.

Phylogenetic Dating

Estimate how long branches are.

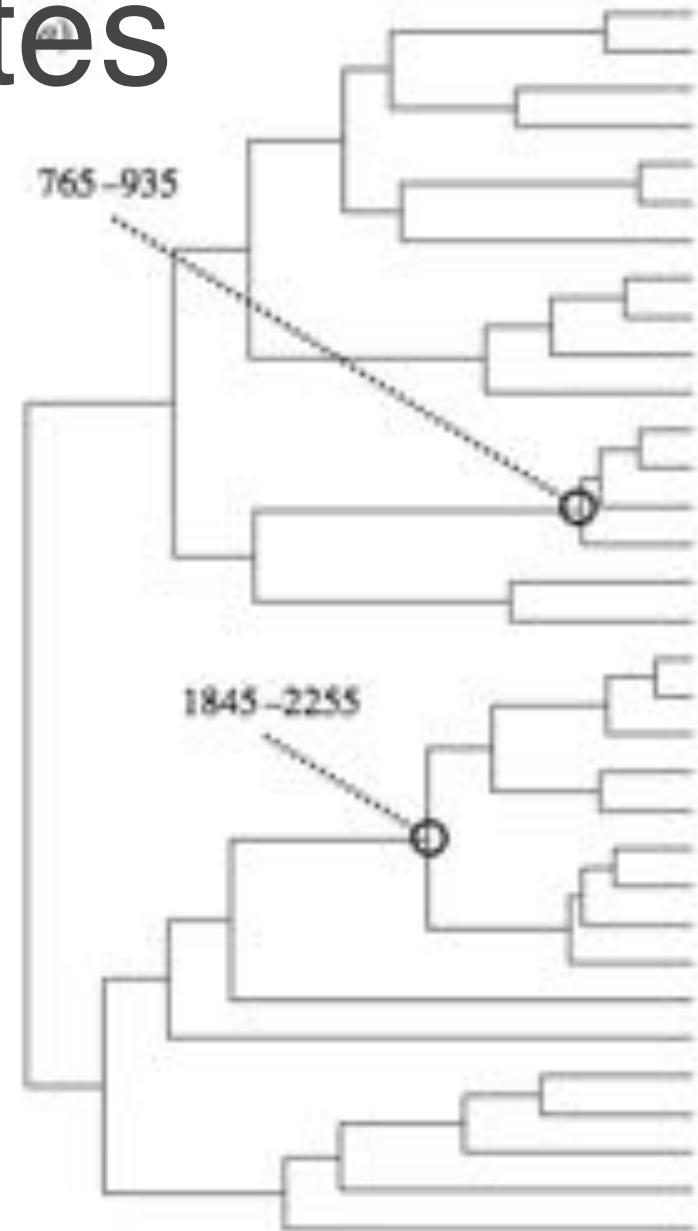
(number of changes per cognate set)

Awesome: Not a global retention rate (i.e. glottochronology) but a **per-language** estimate of the amount of change.



Convert Rates to Dates

- Use (pre)historical information to calibrate nodes
 - e.g. Archaeology suggests initial settlement was..
 - e.g. Historical evidence says that X and Y were separate by...
- Smooth rates over these calibrations



Strict Clock

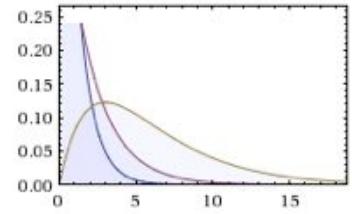
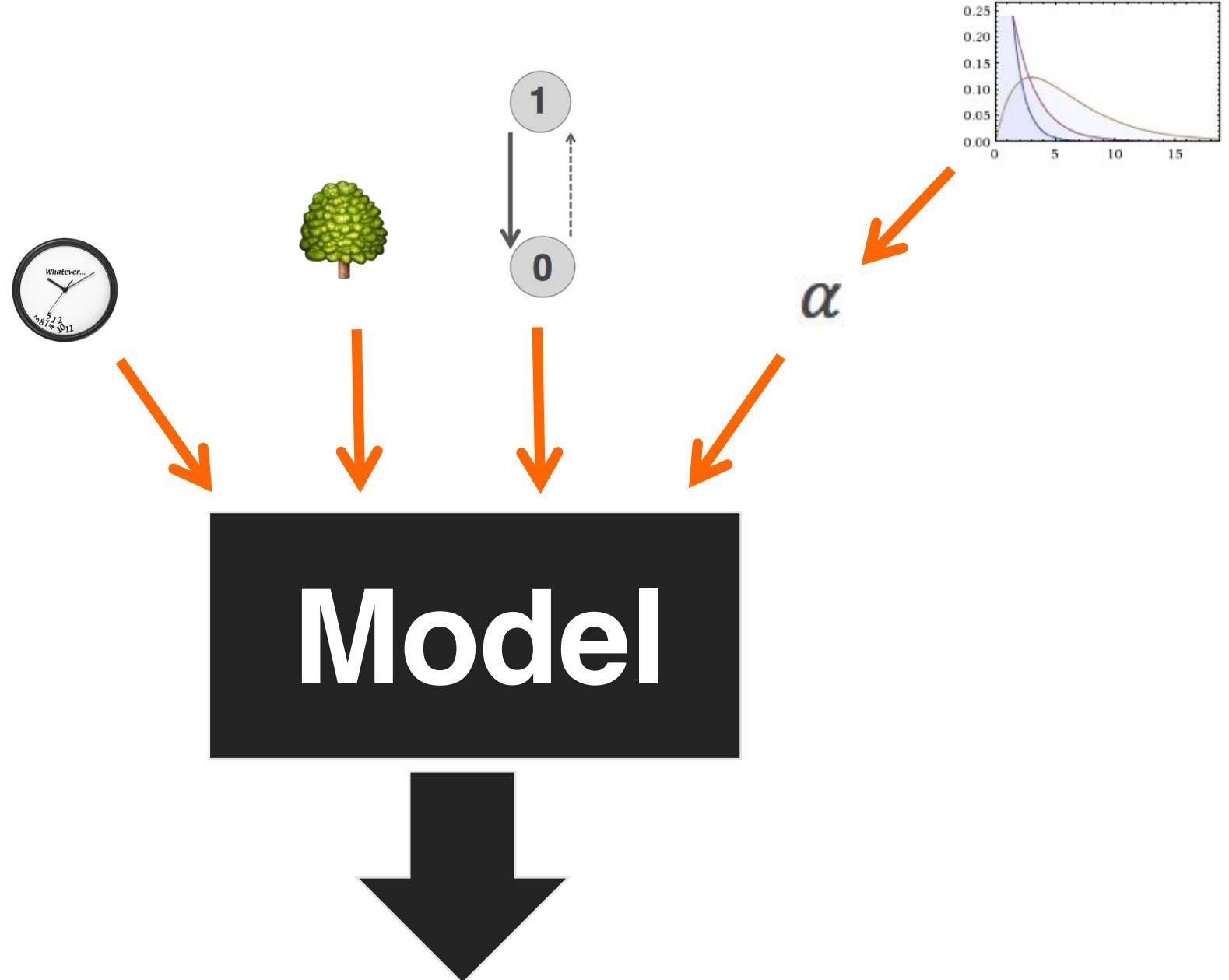


- One rate for all languages.
- No variation
- = glottochronology

Relaxed Clock



- Allows rate to vary across branches.
- Rates are drawn from a parametric distribution with parameters estimated from the data



Problem: uncertainty

- ML and MP give **a point estimate**.
- But a single tree is not enough.
 - Reality is complicated.
 - Need to estimate uncertainty around that estimate.
 - “confidence intervals” = how confident can I be about my estimate.
- MP/ML Solution is “bootstrapping” = painfully labor intensive.

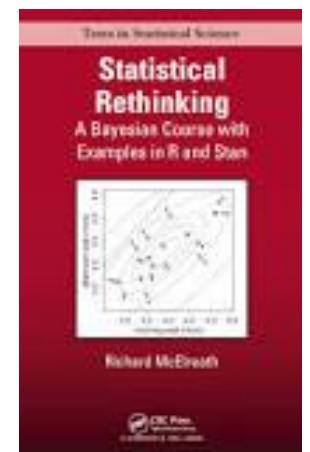


Bayesian methods

1. Explicitly account for uncertainty
2. Directly incorporate different information (*priors*)
3. It's 2016. Everyone is Bayesian.

<http://www.stat.wisc.edu/~ane/bot940/bayes.pdf>

McElreath “Statistical Rethinking”
(<http://xcelab.net/rm/statistical-rethinking/>)

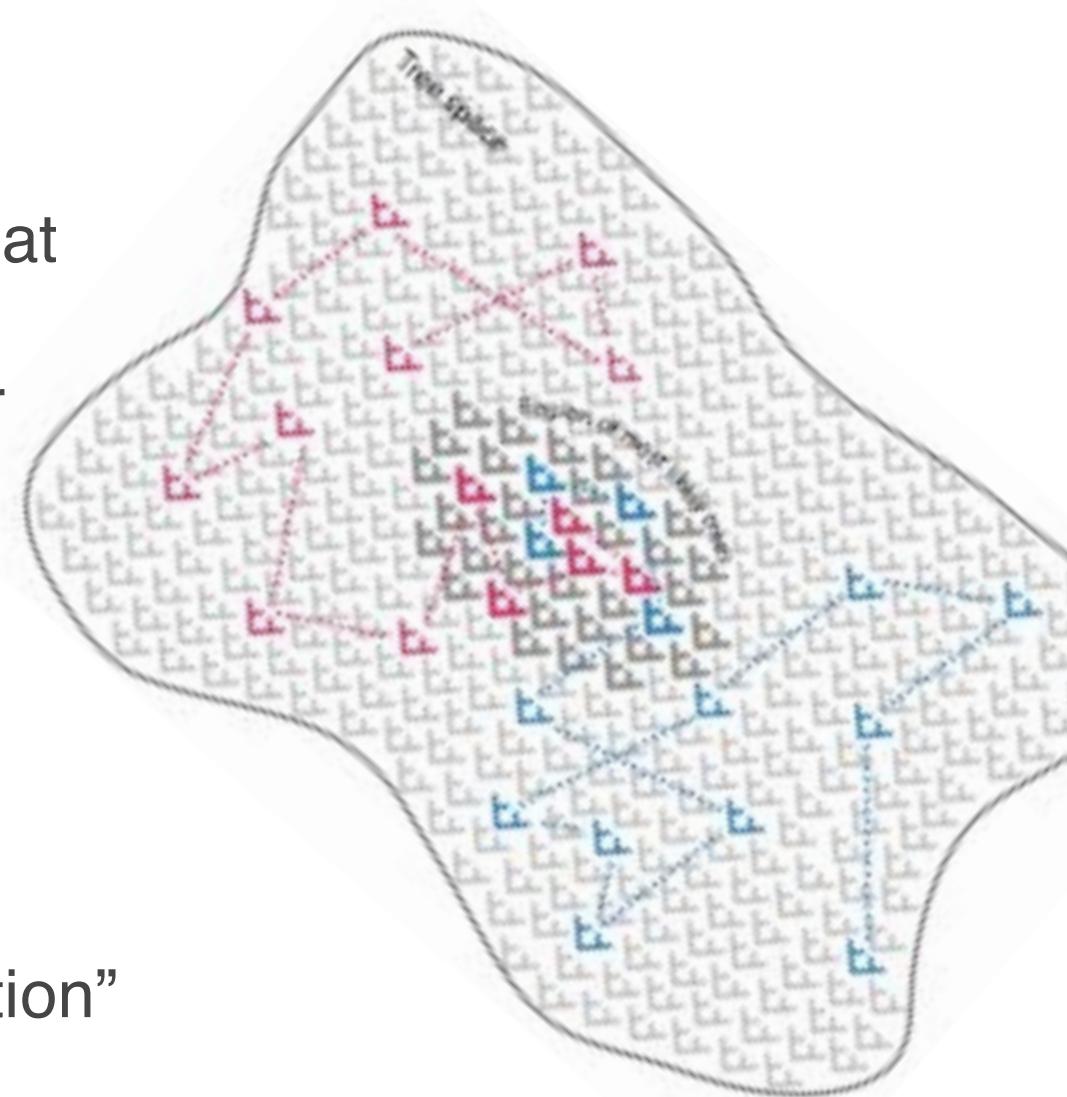


Bayesian Phylogenetics

1. Data & Model & Tree
2. Sample a Tree
3. Calculate the **Likelihood** of that tree
4. Modify the tree or a parameter
5. Repeat (MCMC walk through treespace)

Samples best fitting trees from all possible trees.

- a distribution of trees
- = “Posterior Probability Distribution”



Posterior Probability Distribution: Uto-Aztec Languages

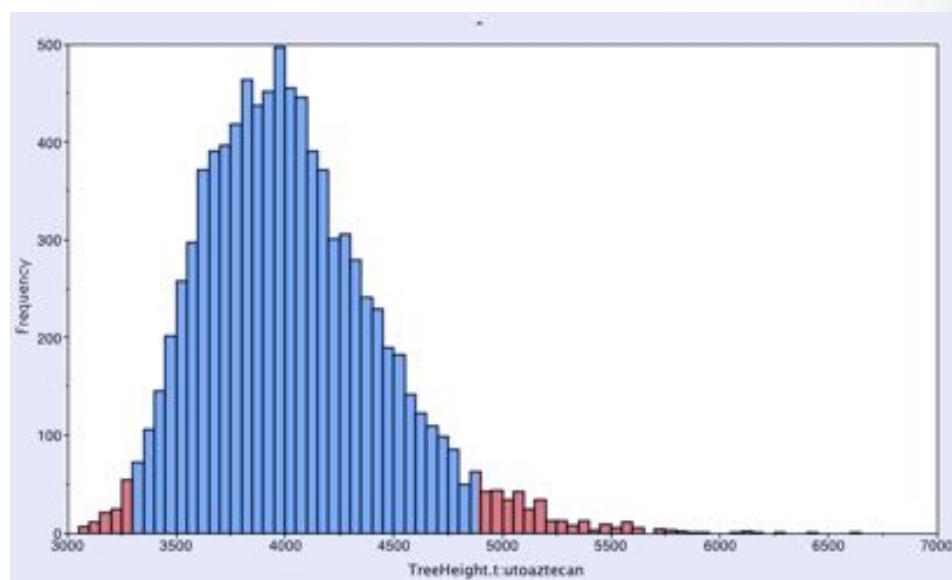
Sampled 10,000 trees.

Draw each one (Densitree, Bouckaert '10)

Some well supported regions

Some less...

Some conflict

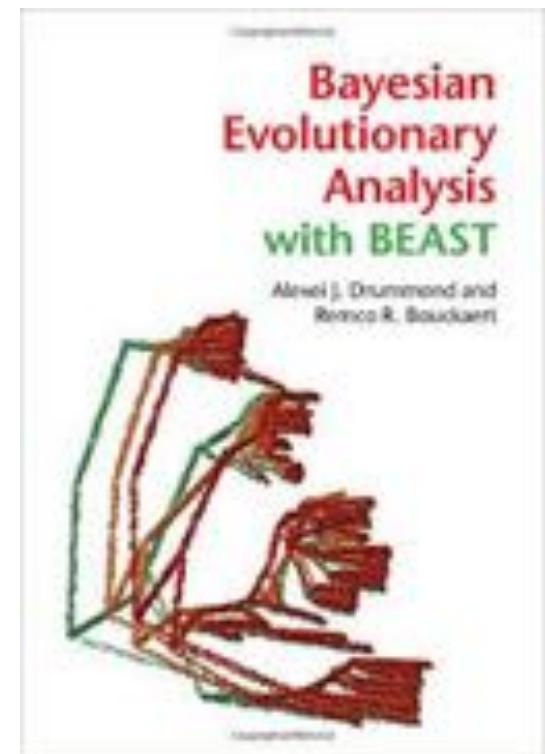


95% HPD: 3302 – 4884 years



Summary

- Many conceptual parallels across disciplines.
- Grew out of many of the same concerns found in Linguistics/Anthropology (innovations vs retentions, rate variation, conflicting signal, etc).
- Wide range of techniques available, but Bayesian methods are clear winners to date.
- Practical: My BEAST tutorial in QMSS materials, BEAST Tutorial online, BEAST book.



Questions?

