

Phylogenetics

Simon J. Greenhill



ARC CENTRE OF EXCELLENCE FOR
THE DYNAMICS OF LANGUAGE

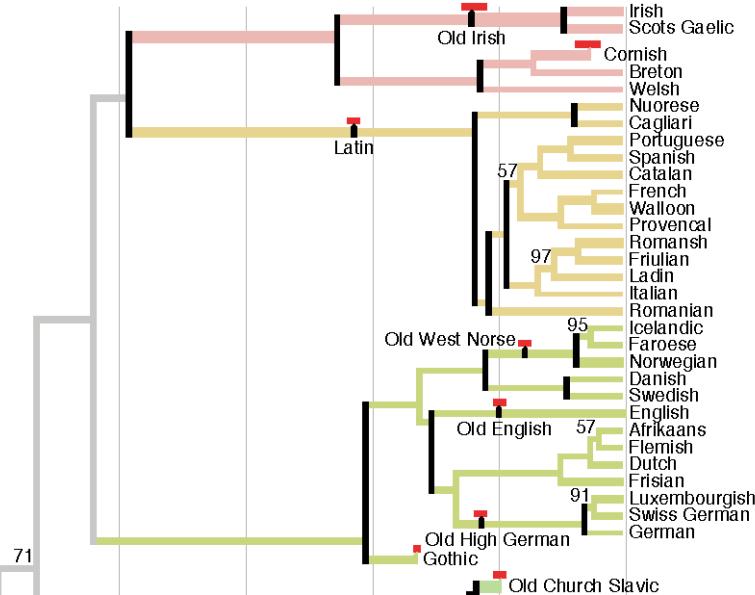
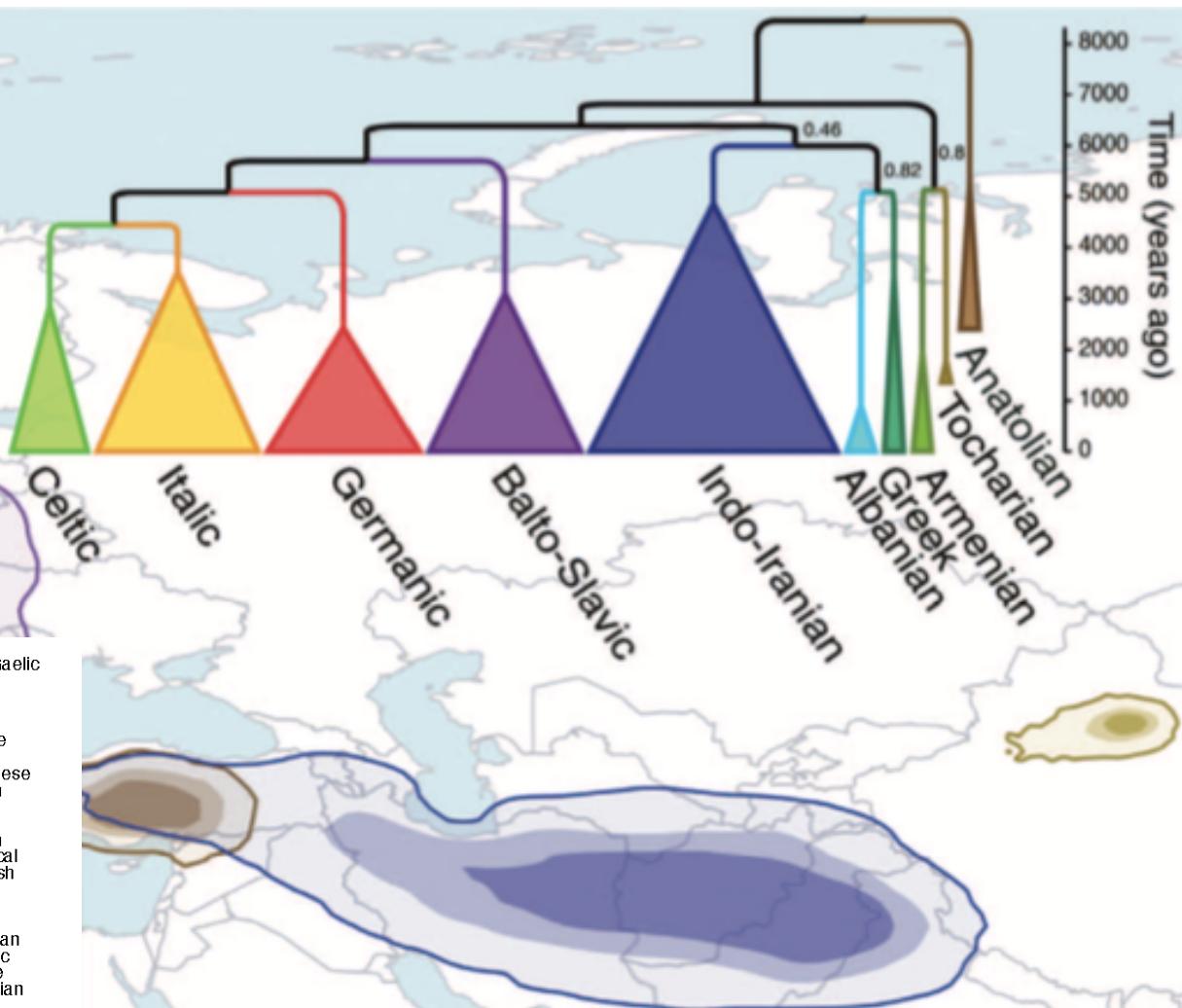
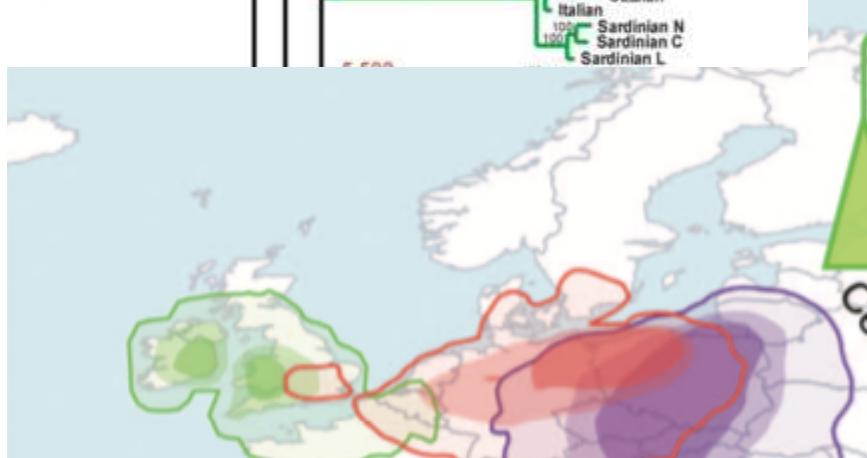
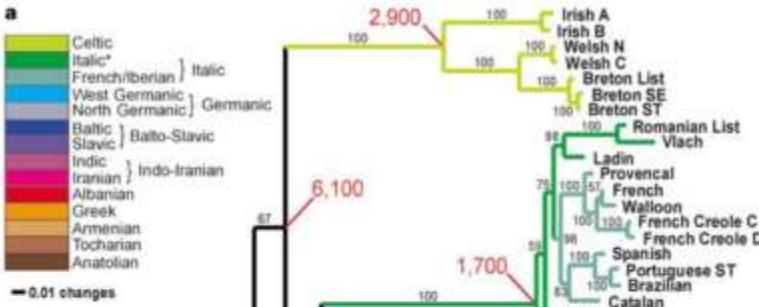


Max Planck Institute for the
Science of Human History

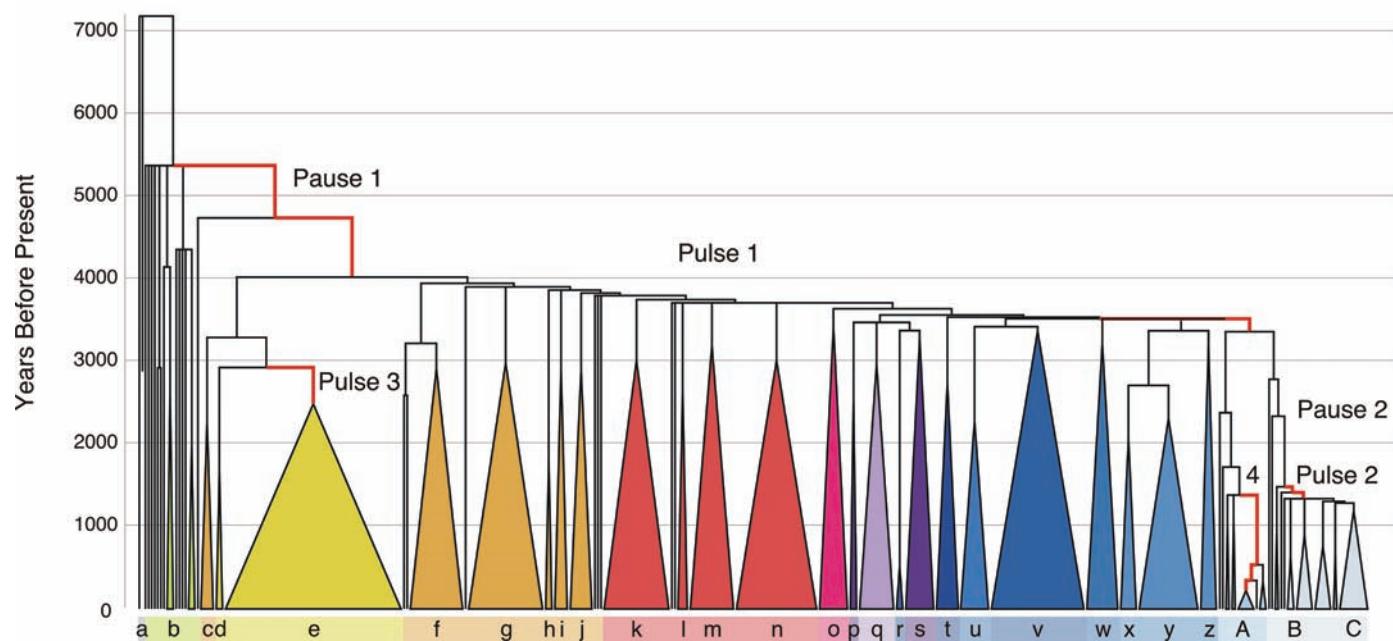
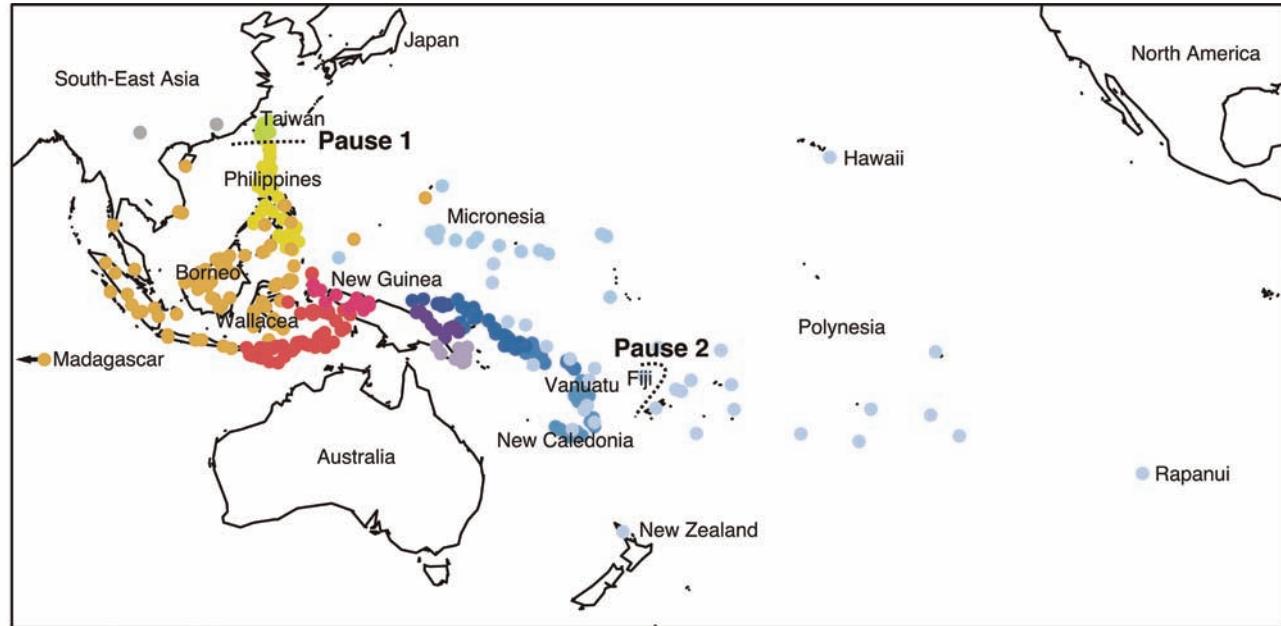
Phylogenetics

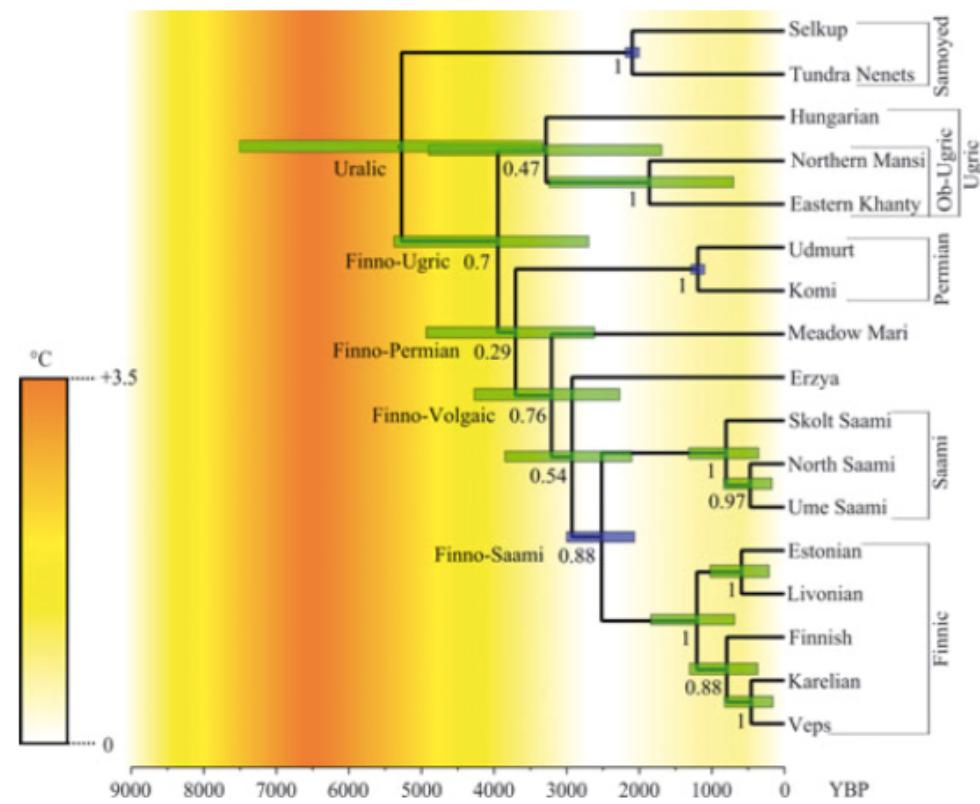
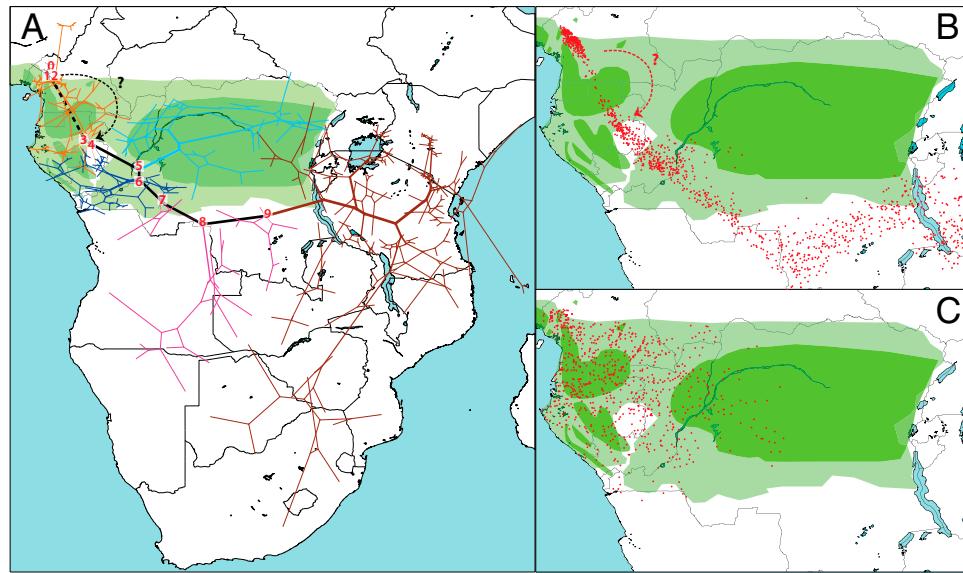
Range of methods in a robust, statistical/inferential framework for thinking about, quantifying, analysing trees.

Statistical Estimation Problem (e.g. Cavalli-Sforza & Edwards '82) supported by a theory of “tree-thinking” (O’Hara ‘88, ‘92).

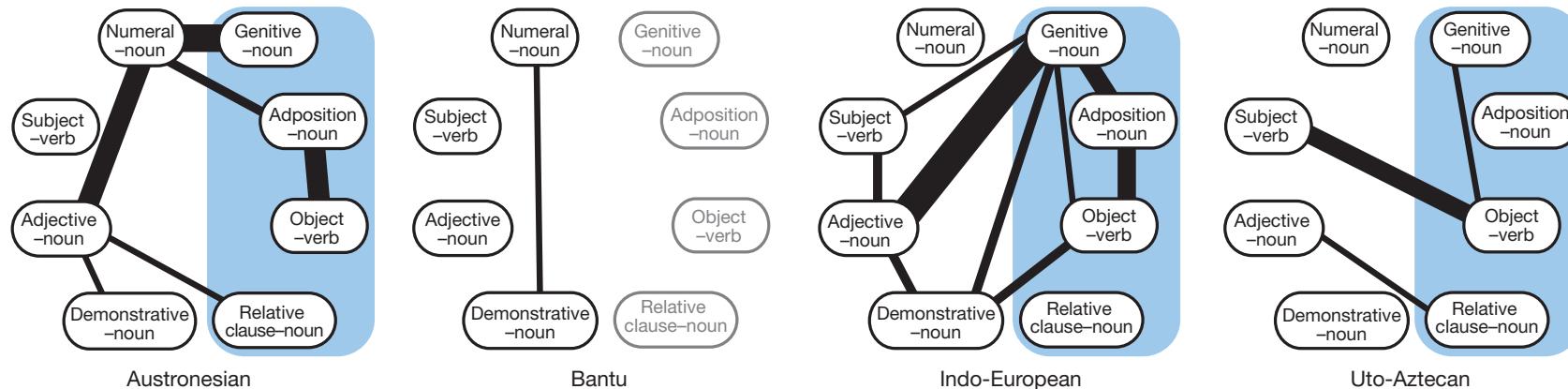


Gray & Atkinson '03, Bouckaert et al' 11, Chang et al '14.





Dunn et al. '15

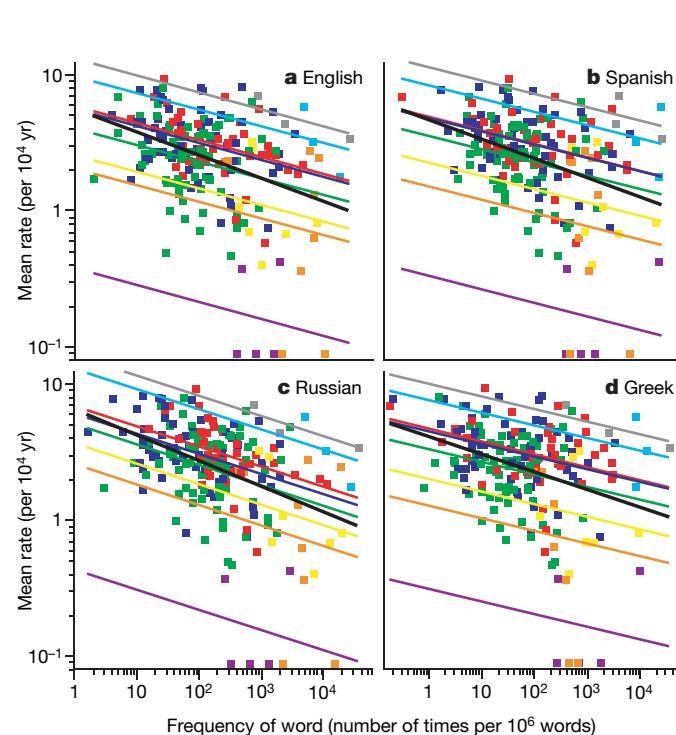
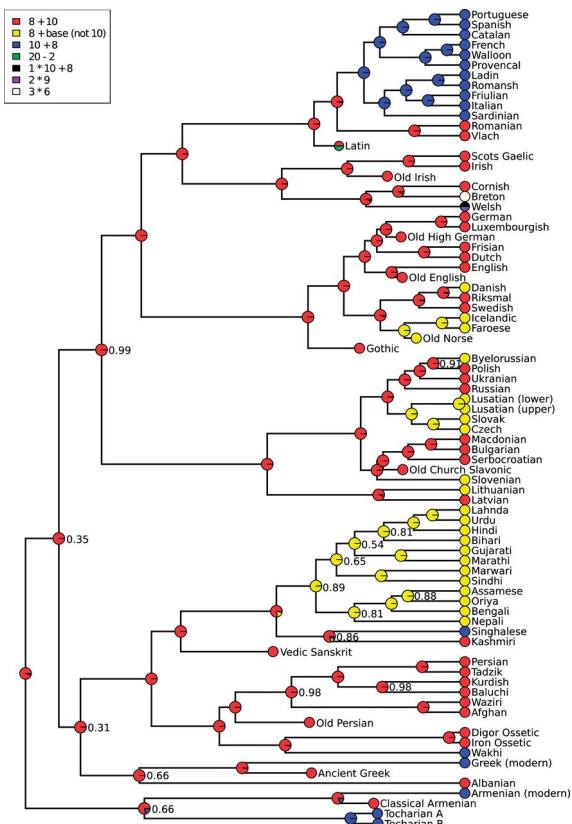


Austronesian

Bantu

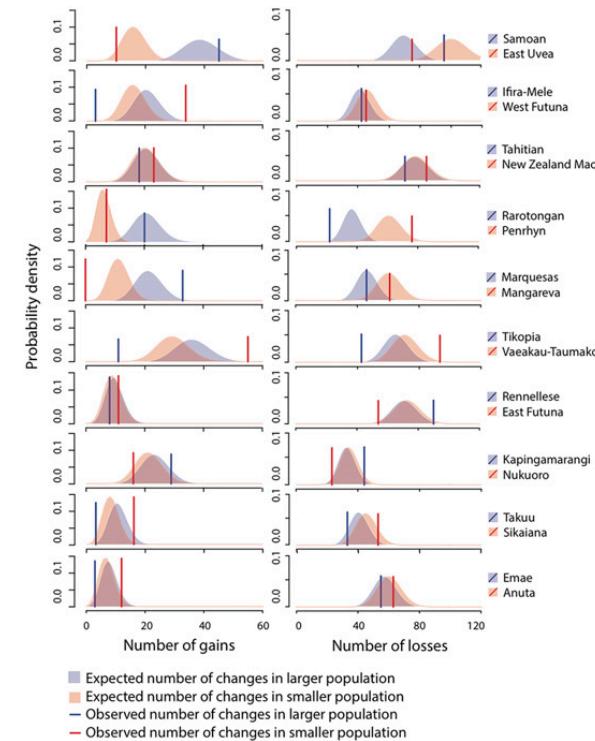
Indo-European

Uto-Aztecian



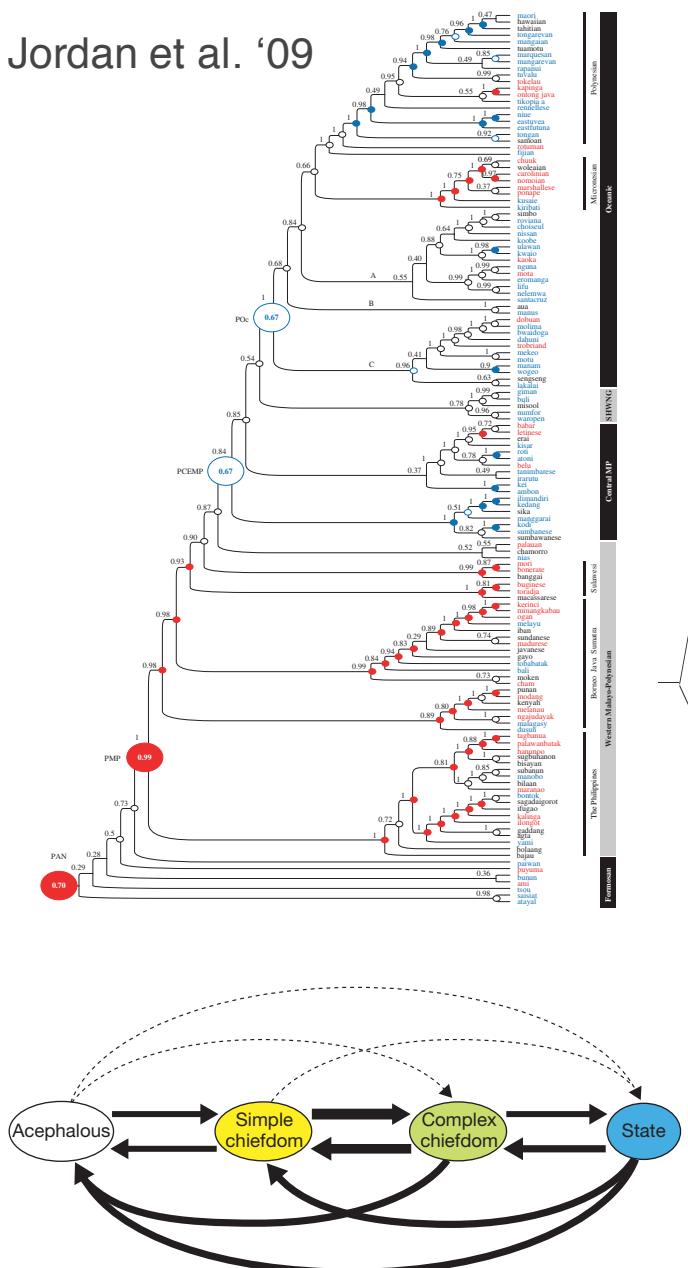
Calude and Verkerk '16

Pagel et al '07



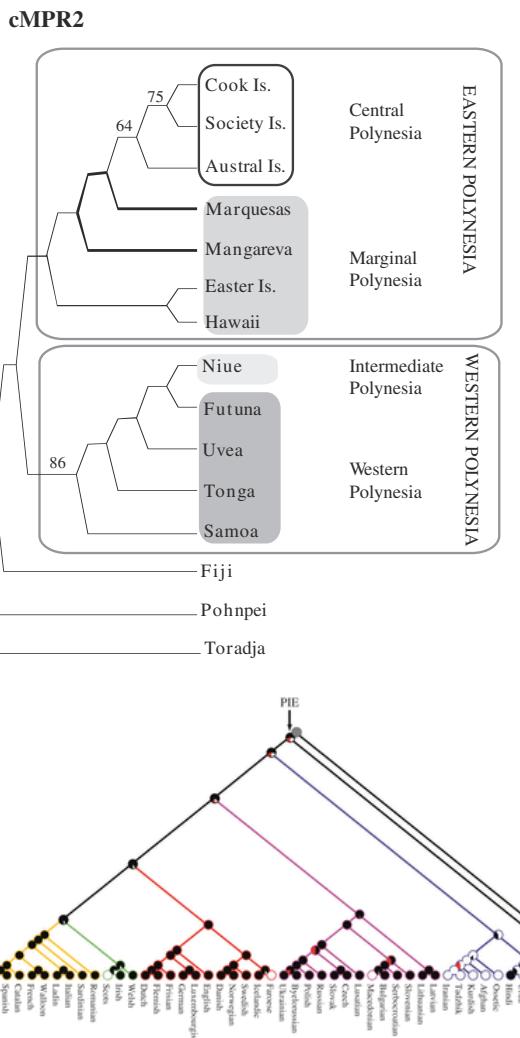
Bromham et al. '15

Jordan et al. '09

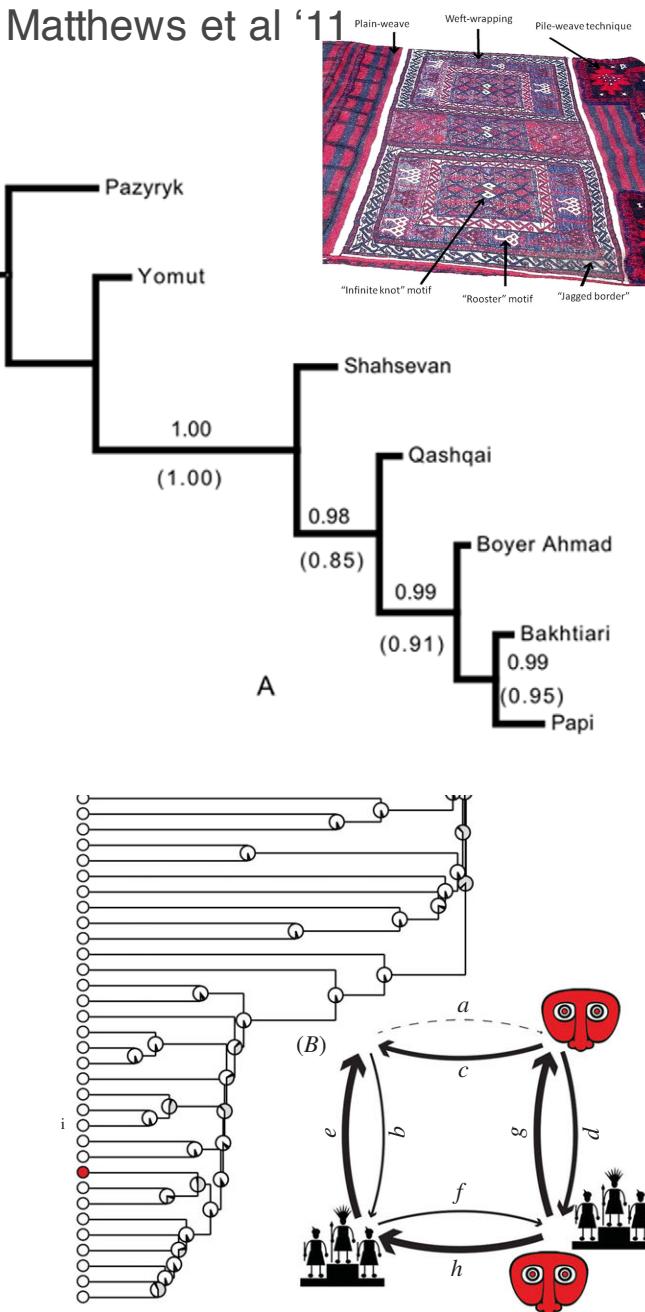


Currie et al. '10

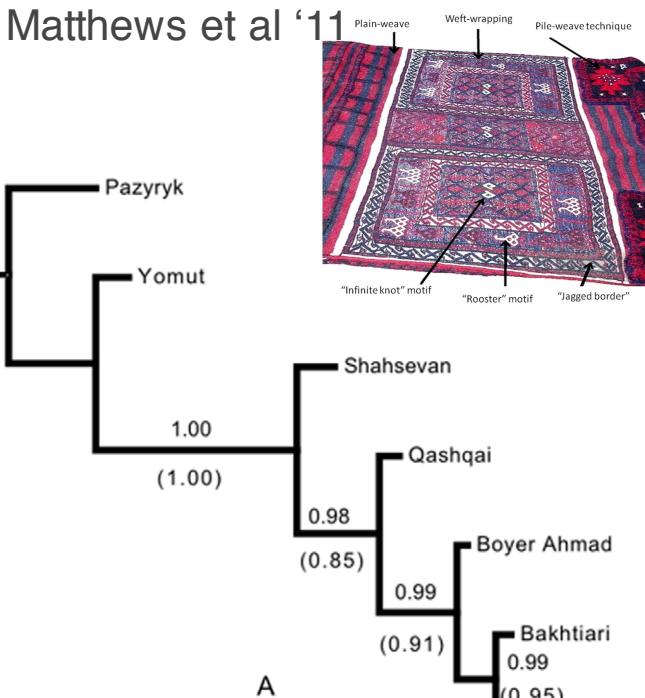
Larsen '11



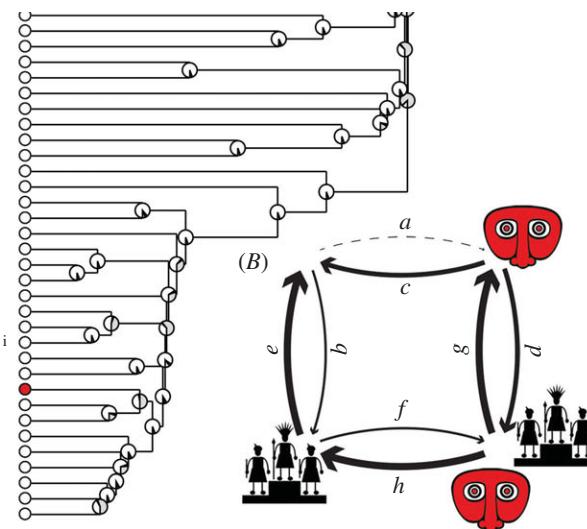
Da Silva & Tehrani '16



Watts et al. '15



A



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”

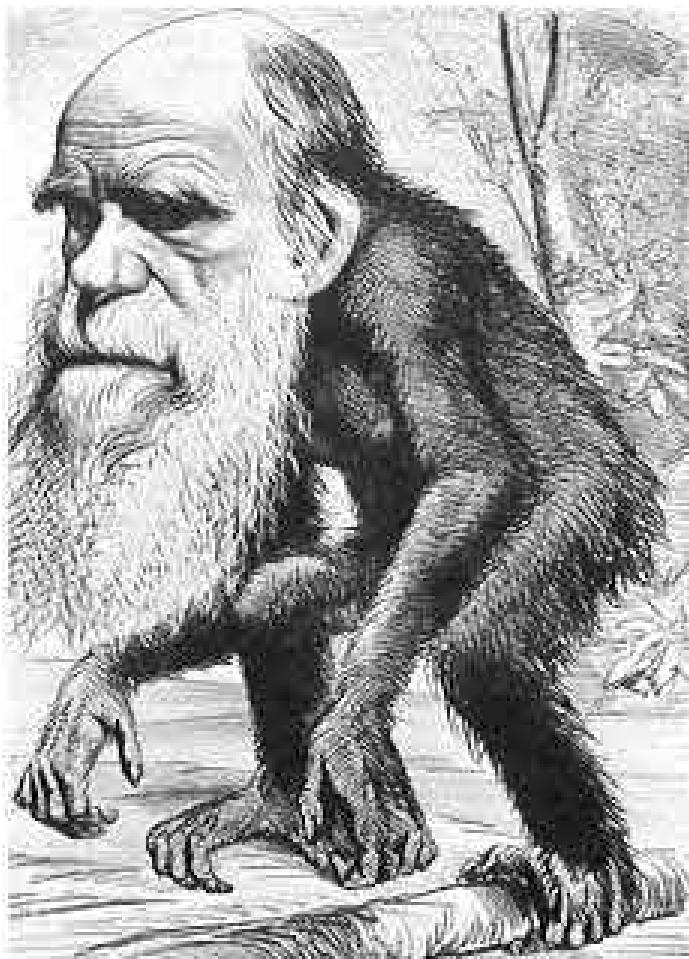
“utter bollocks”

“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?

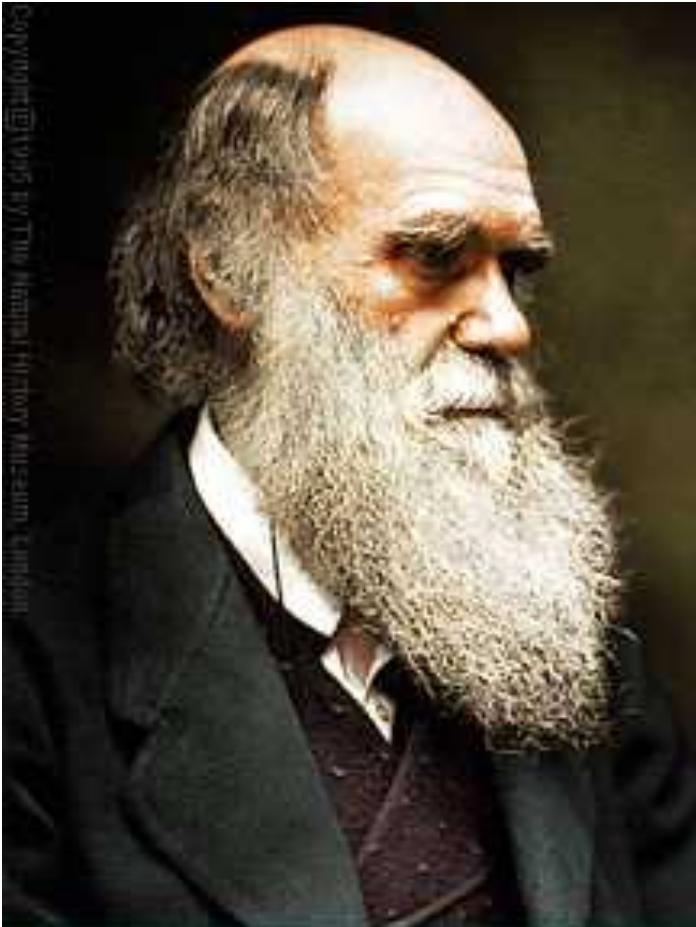
What **processes** shaped 🐨🟡💬?

Can we infer what 🐨🟡💬 were in the **past**?

Conceptual Parallels

🐨 Speciation	👩💬 Divergence
🐨 Homologies	👩💬 Cognates
🐨 Mutations	👩💬 Innovations
🐨 Random Drift	👩💬 Random Drift
🐨 Natural selection	👩💬 “Social” selection
🐨 Hybridisation	👩💬 Borrowing
🐨 Fossils	👩💬 Ancient Texts
🐨 Extinction	👩💬 Death

Descent of Man (1871)



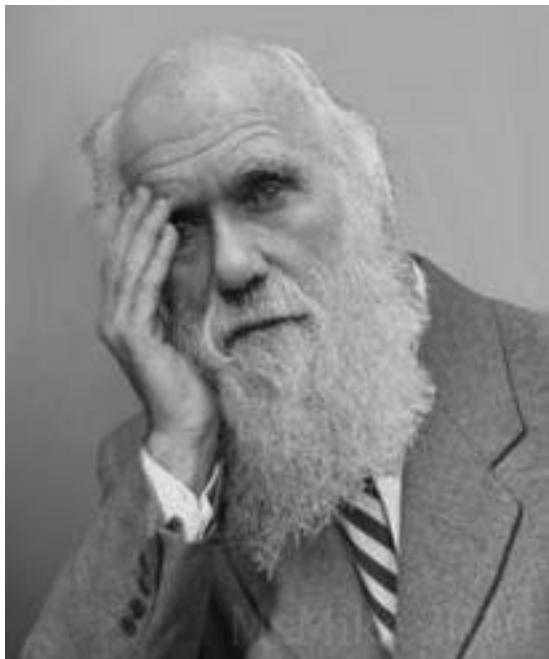
"Languages, like organic beings, can be classed in groups under groups..."

Dominant languages and dialects spread widely, and lead to the gradual extinction of other tongues. A language, like a species, when once extinct ... never reappears."

...

"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 was late



~400: BCE

Socrates & Aristotle recognised Greek had changed since Homer (730 BCE)

1786:

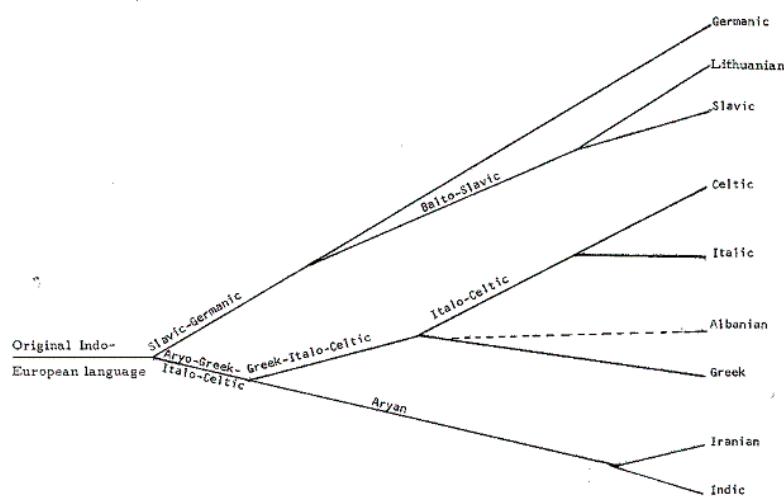
Sir William Jones: Sanskrit languages have “sprung from some common source”

1814:

Rasmus Rask: Comparative Method

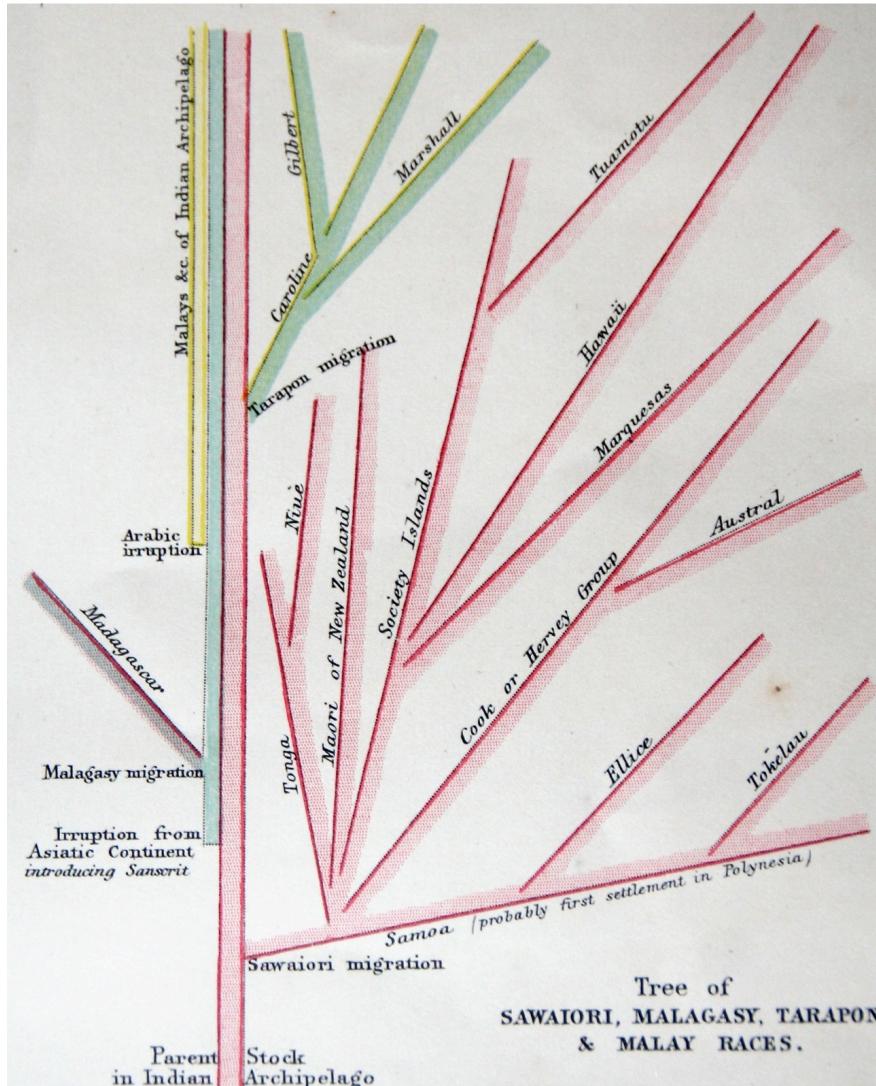
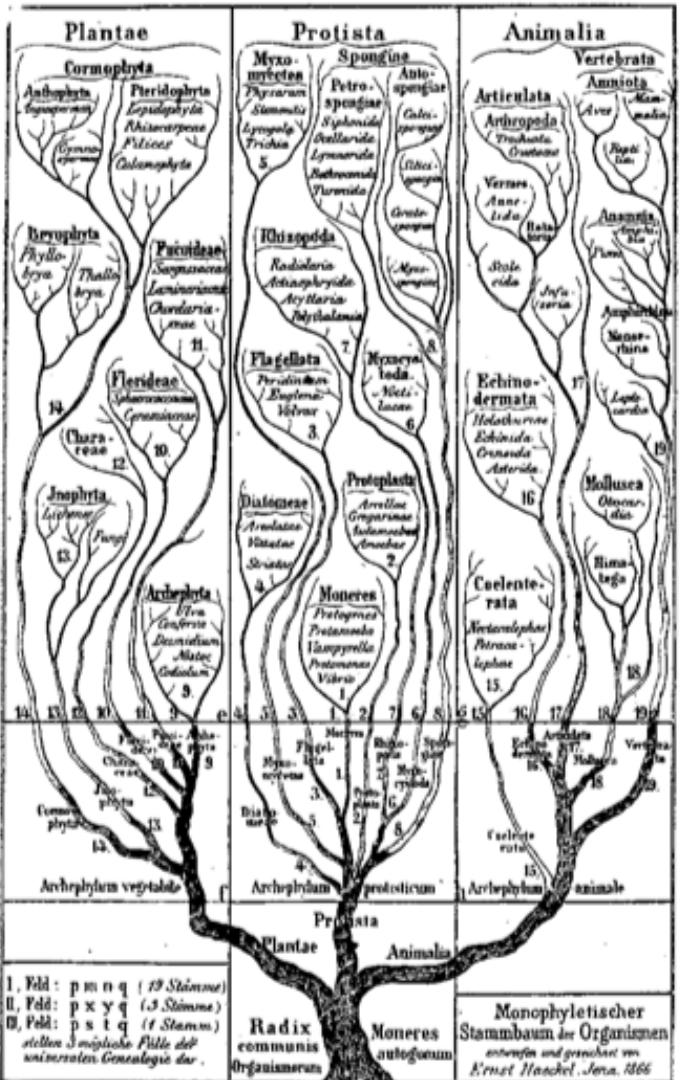
August Schleicher 1863

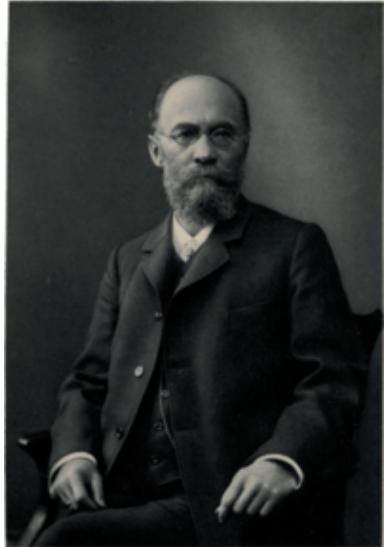
Darwinism Tested by the Science of Language



“same process has long been generally assumed for linguistic organisms”

“We set up family trees of languages known to us in precisely the same way as Darwin has attempted to do for plant and animal species.”





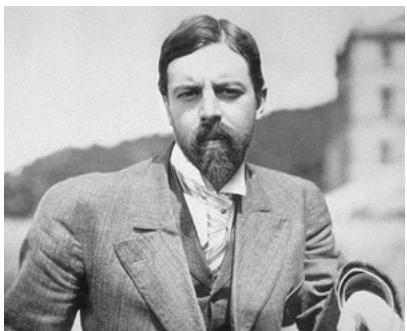
Brugmann (1884)

Importance of using “shared innovations” to identify clades and not “shared retentions”



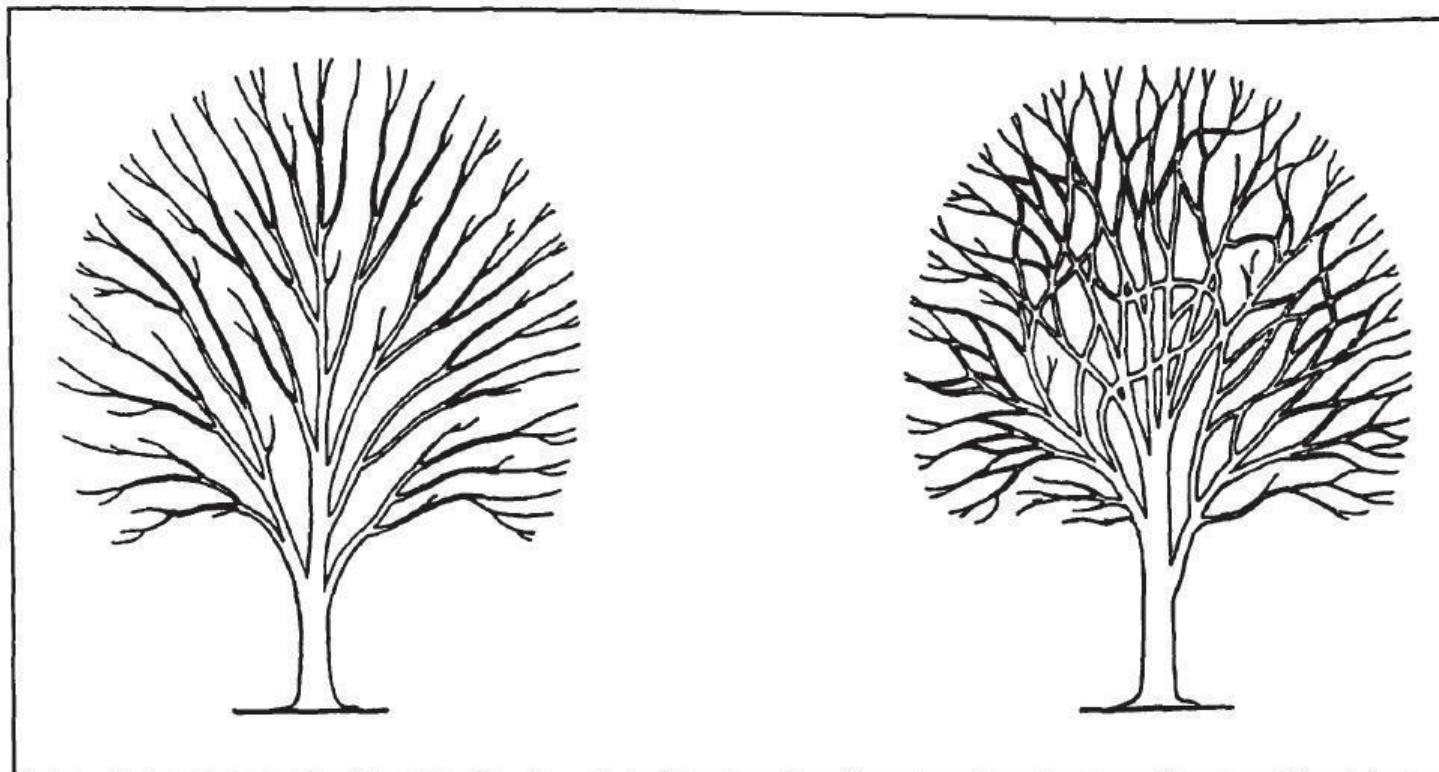
Hennig (1966)

“Synapomorphies” and
“Symplesiomorphies”



Kroeber (1948)

"The tree of life is eternally branching, and never doing anything fundamental but branching, except for the dying-away of branches. The tree of human history, on the contrary, is constantly branching and at the same time having its branches grow together again. Its plan is therefore much more complex and difficult to trace." (p.86)



Bacteria

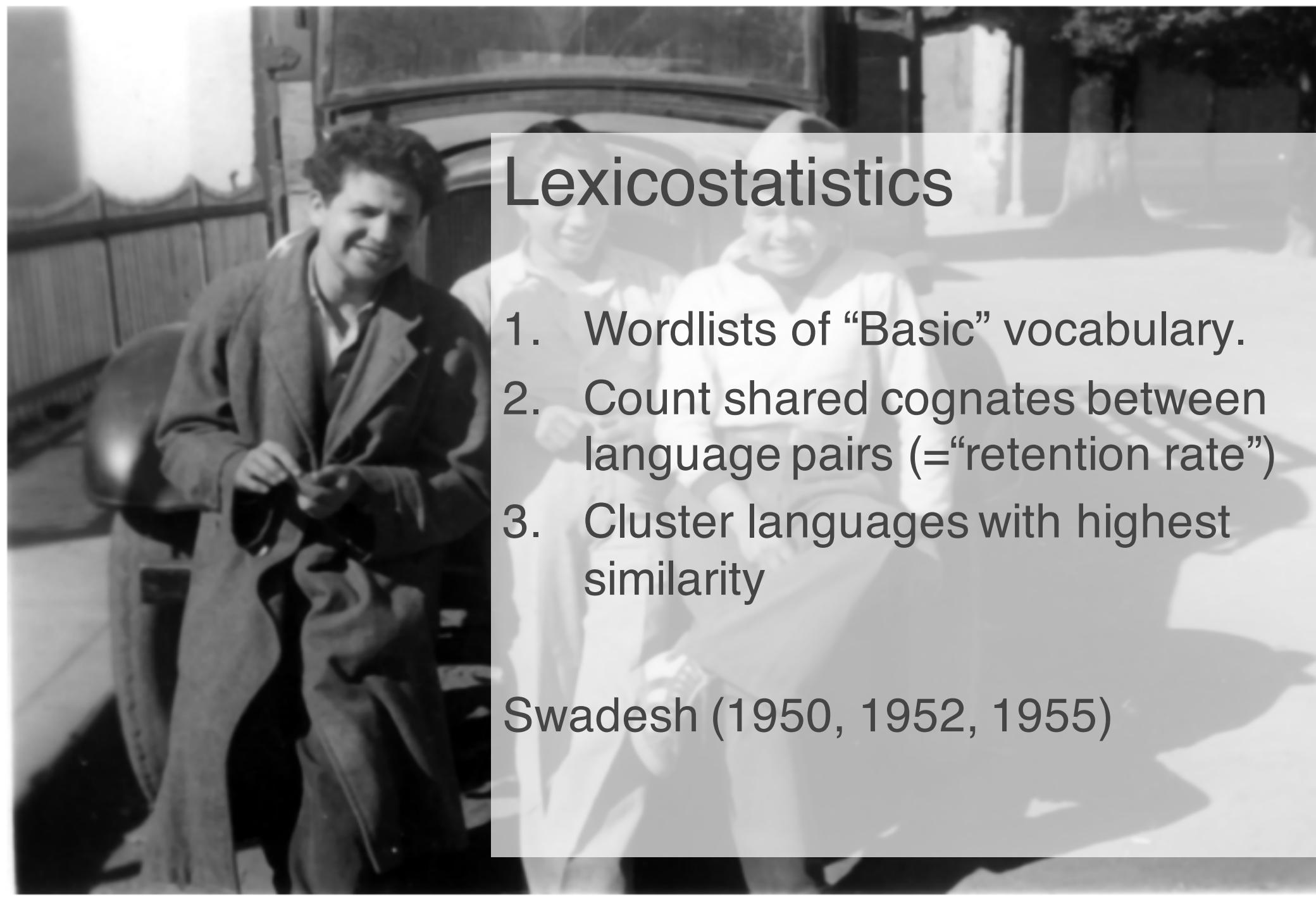


Eukarya

Archaea







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.

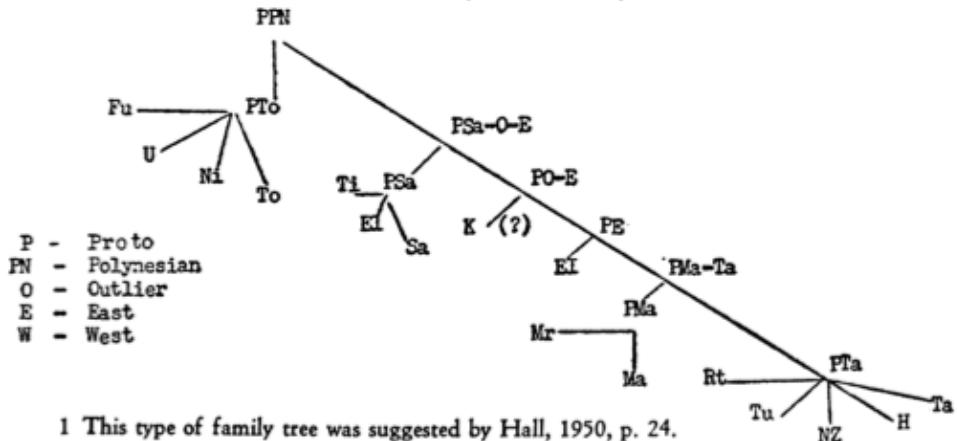
Elbert 1953

TABLE 2
Polynesian cognate percentages

U	N1	To	T1	E	Sa	Si ¹	Pi ¹	OJ ¹	Mu ¹	K	EI	Mr	Ma	Rt	Tu	NZ	H	Ta
63	62	74	83	79	74	57	62	62	54	52	54	57	56	66	62	61	58	58
	72	86	78	74	70	53	59	58	52	51	53	53	51	62	62	61	55	59
	64	68	61	63	50	59	49	54	49	49	55	47	56	55	54	49	51	
	70	64	66	46	55	55	49	49	45	48	49	45	58	53	54	49	52	
			61	76	66	66	62	59	59	62	67	63	71	68	71	67	66	T1
				78	66	61	62	59	58	61	62	63	66	66	67	68	66	E
					64	60	61	55	53	53	55	52	67	62	57	59	60	Sa
						55	60	52	55	52	55	56	60	58	60	60	59	Si
						57	52	54	54	54	59	53	55	61	60	55	54	Pi
						51	51	54	53	53	51	53	56	55	56	53	53	OJ
						53.	49	45	44	44	53	50	52	49	48	48	48	Nu
							47	49	45	54	51	51	51	49	50		K	
								64	63	64	62	63	64	62				EI
								73	75	72	70	69	68					Mr
								73	73	69	67	70	67					Ma
									83	83	79	85						Rt
									79	77	83							Tu
									71	73								NZ
										76								H

1 Percentages based on incomplete data.

TABLE 4
A tentative family tree for Polynesia¹



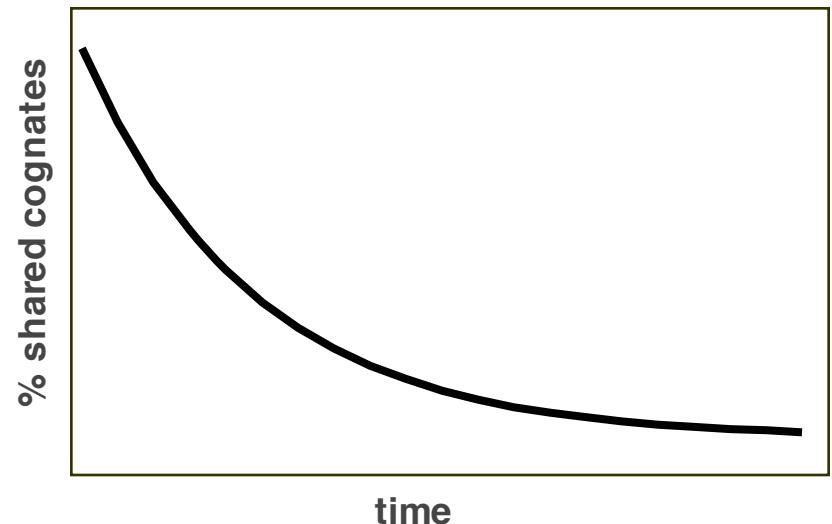


Glottochronology

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

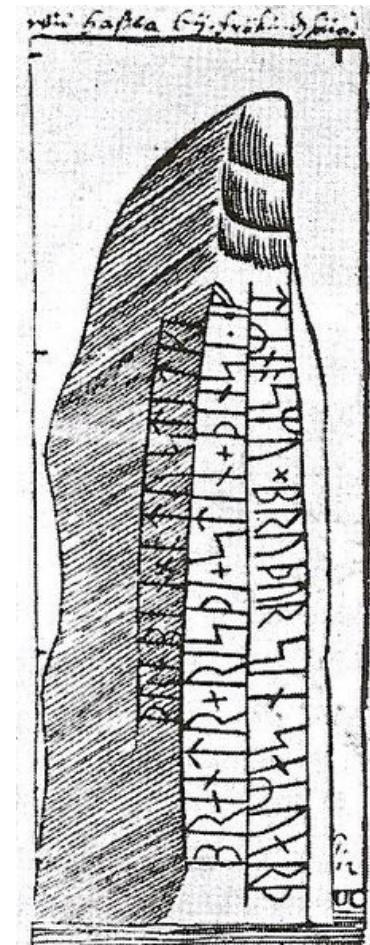
...and the fall of Lexicostatistics

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

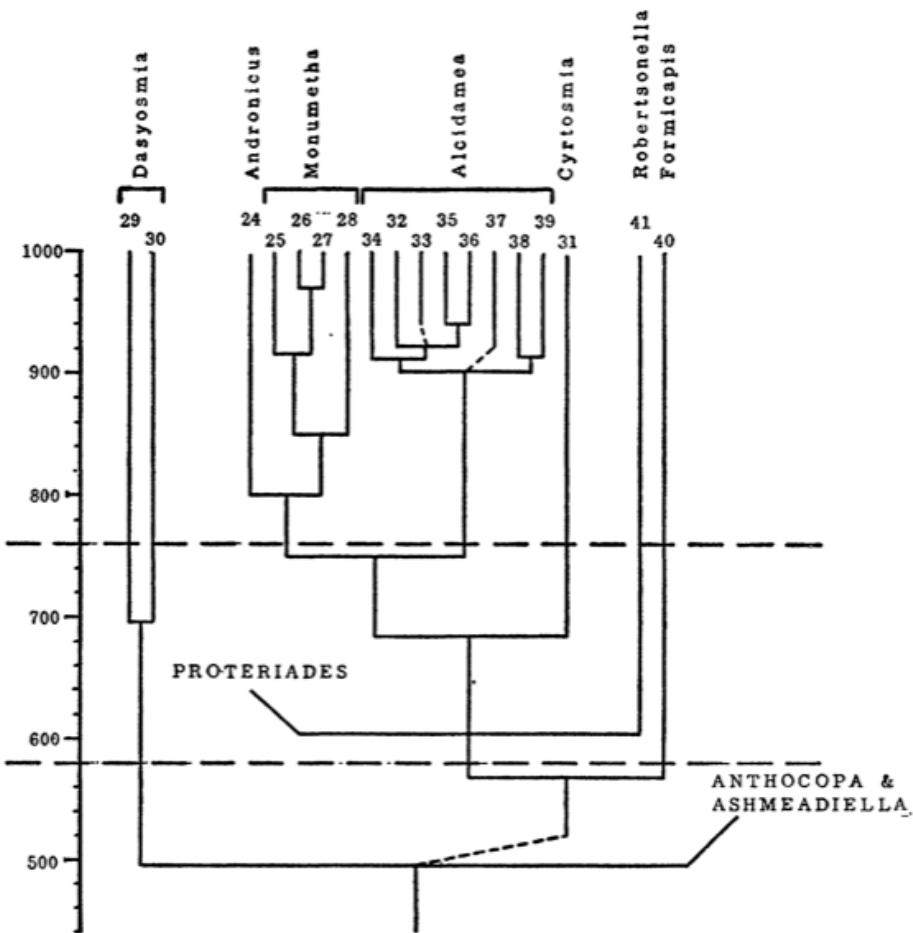


FIG. 6. Diagram of relationships for the genus *Hoplitis* obtained by the weighted variable group method.

A QUANTITATIVE APPROACH TO A PROBLEM
IN CLASSIFICATION¹

CHARLES D. MICHENER AND ROBERT R. SOKAL²

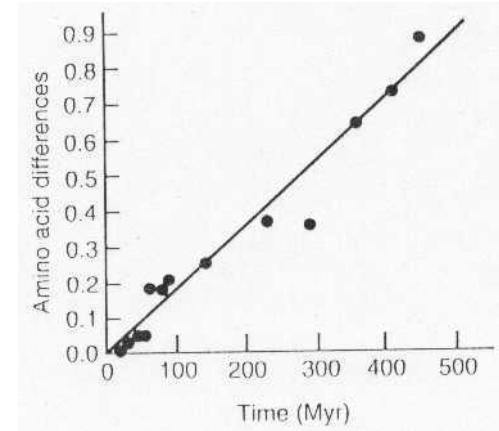
Department of Entomology, University of Kansas, Lawrence



Molecular Clock

Zuckerkandl & Pauling 1962

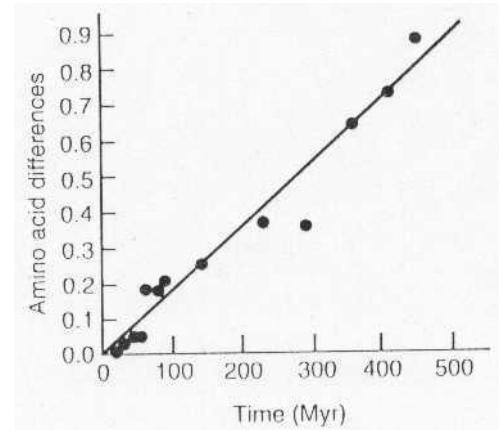
Number of AA differences were proportional to species divergence times.



Molecular Clock

Zuckerkandl & Pauling 1962

Number of AA differences were proportional to species divergence times.



Kimura 1968

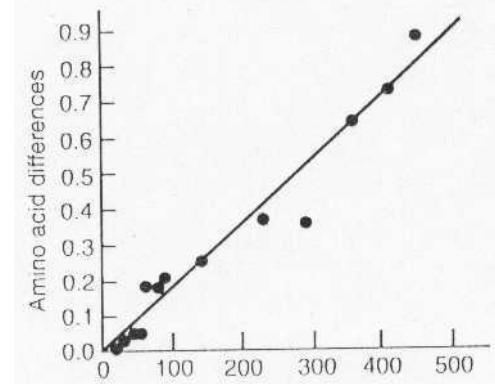
The average time taken for one base pair replacement within a genome is therefore

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

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

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

Sarich & Wilson 1967

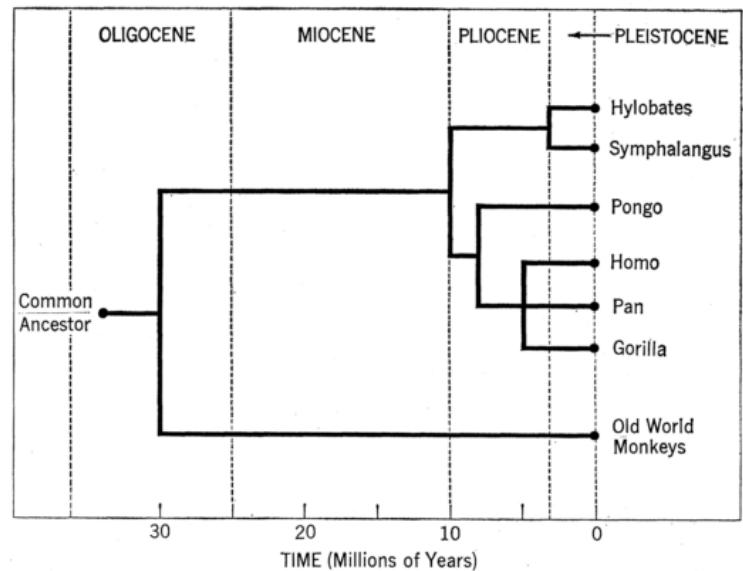


Fig. 1. Times of divergence between the various hominoids, as estimated from immunological data. The time of divergence of hominoids and Old World monkeys is assumed to be 30 million years.

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.

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^{¶¶}

Cavalli-Sforza &
Edwards 1967

Yang 1993

Sanderson 1997

Drummond et al. 2006



Joe Felsenstein

[Follow](#) ▾

Professor of Genome Sciences, and Professor of Biology, University of Washington, Seattle

Evolutionary biology, phylogenetic methods, population genetics

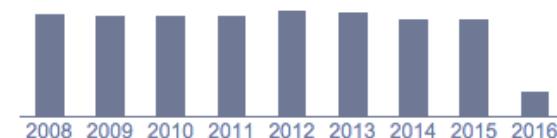
Verified email at gs.washington.edu - [Homepage](#)

Google Scholar

[Get my own profile](#)

Citation indices All Since 2011

Citations	100419	29500
h-index	68	42
i10-index	121	73



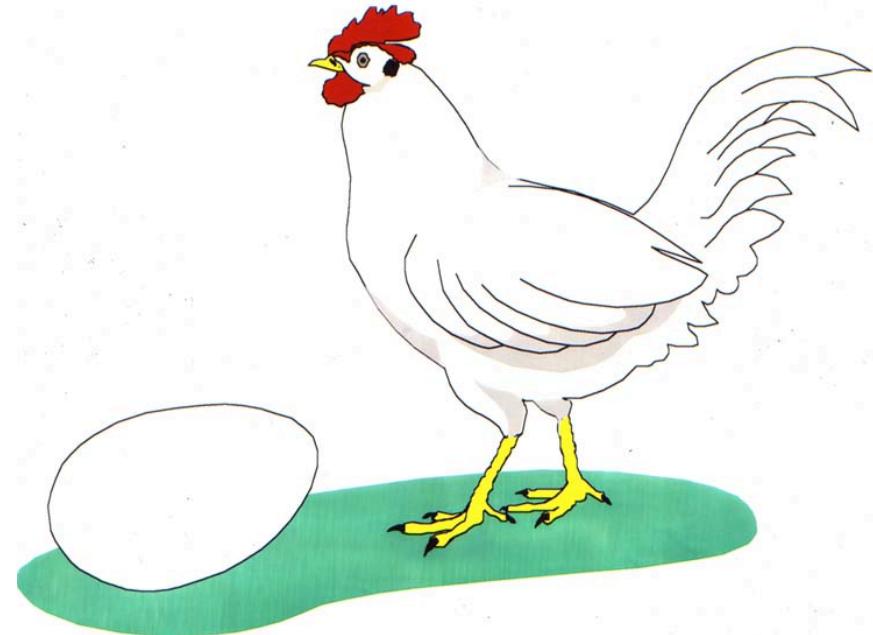
Title	1–20	Cited by	Year
Confidence limits on phylogenies: an approach using the bootstrap J Felsenstein Evolution, 783–791		30037	1985
PHYLIP-phylogeny inference package (version 3.2) D Plotree, D Plotgram cladistics 5, 163–166		23173	1989
Evolutionary trees from DNA sequences: a maximum likelihood approach J Felsenstein Journal of molecular evolution 17 (6), 368–376		8494	1981
Phylogenies and the comparative method J Felsenstein American Naturalist, 1–15		6259	1985
Inferring phylogenies J Felsenstein, J Felsenstein Sinauer Associates		3743	2004
Cases in which parsimony or compatibility methods will be positively misleading J Felsenstein Systematic Biology 27 (4), 401–410		3057	1978
Phylogenies from molecular sequences: inference and reliability J Felsenstein		2252	1988

Year

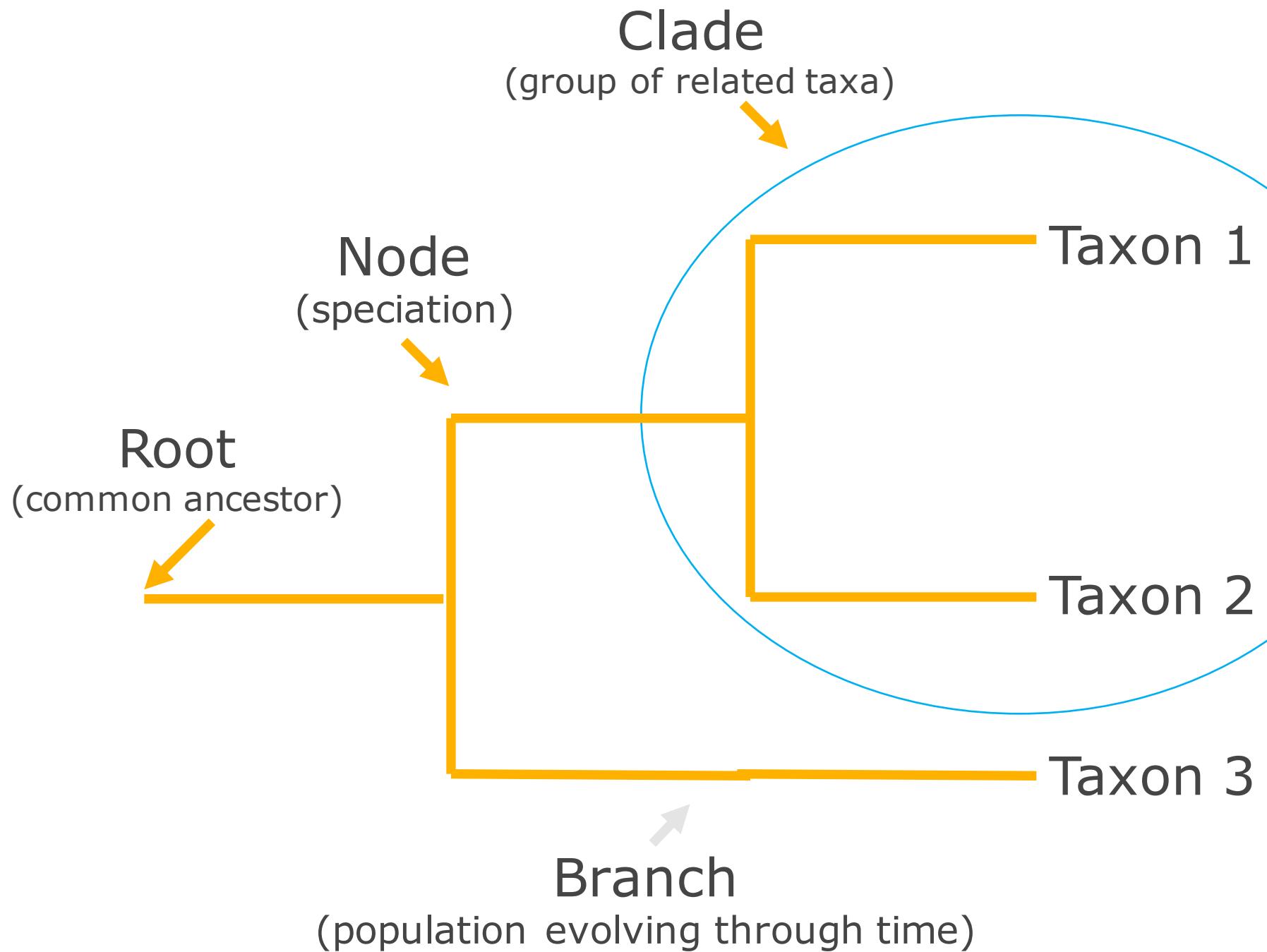
Trees are
coming

A glowing neon sign is mounted on a dark, textured brick wall. The sign features the words "Trees are" in a simple, sans-serif font, and "coming" in a flowing, cursive script. The neon light is white, casting a soft glow on the surrounding wall. The sign is held in place by four visible mounting brackets.

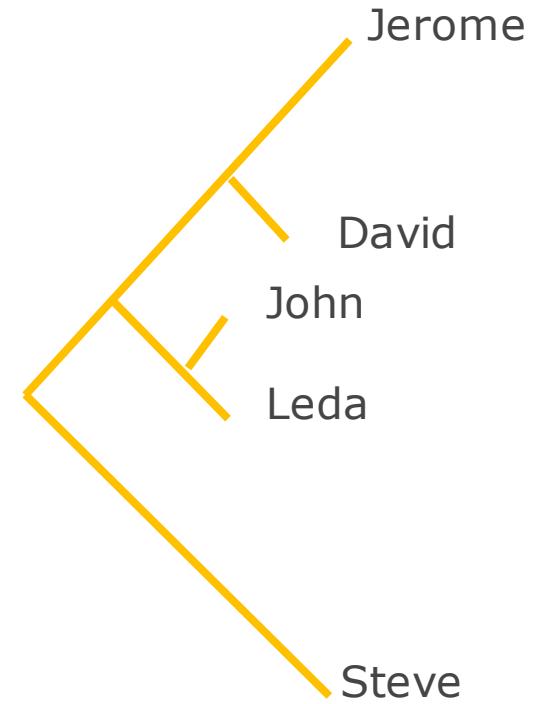
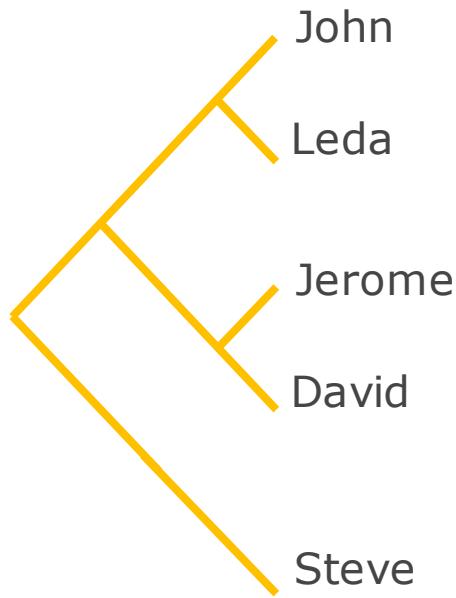
Chicken or the Egg?



Shykoff & Widmer (1998)
Trends in Ecology and Evolution, 13, 158.

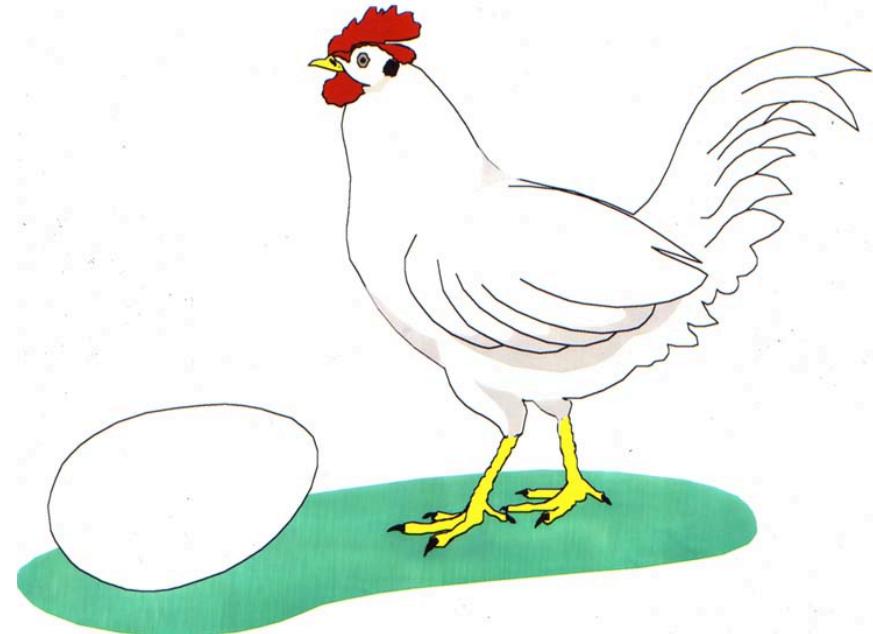
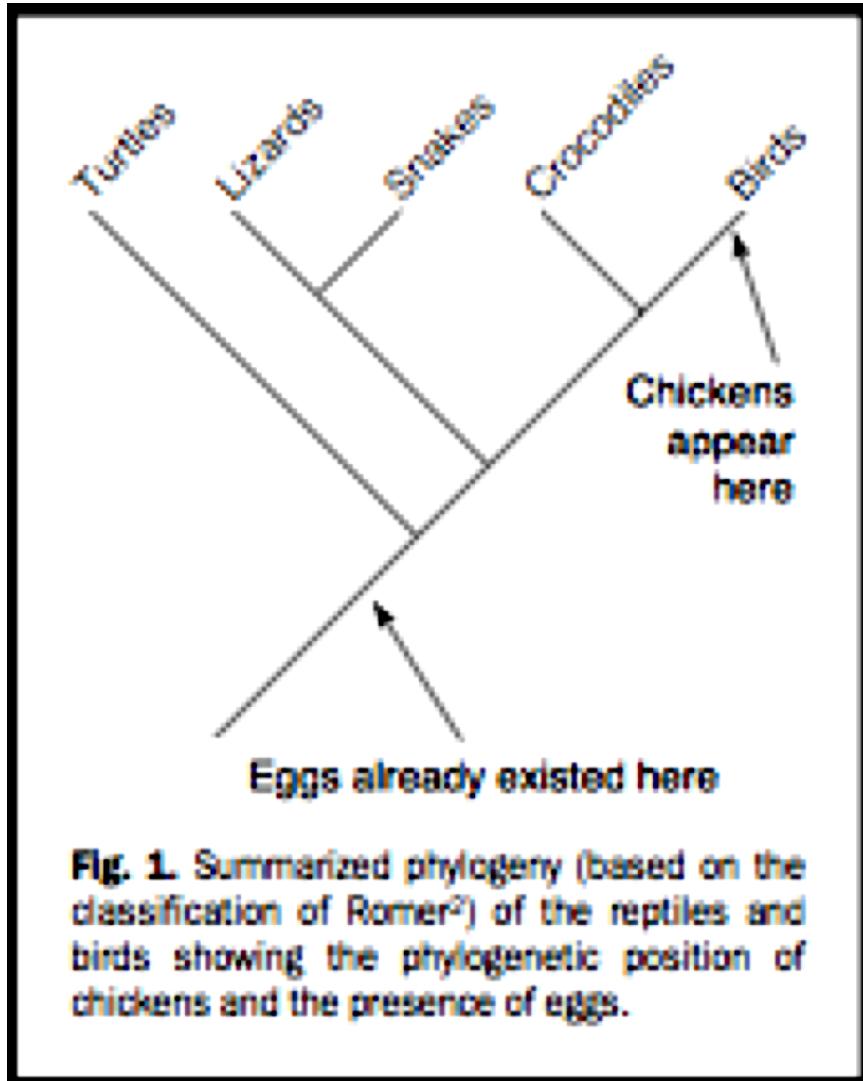


Trees rotate around nodes



ORDER is not important. GROUPINGS are.

Chicken or the Egg?



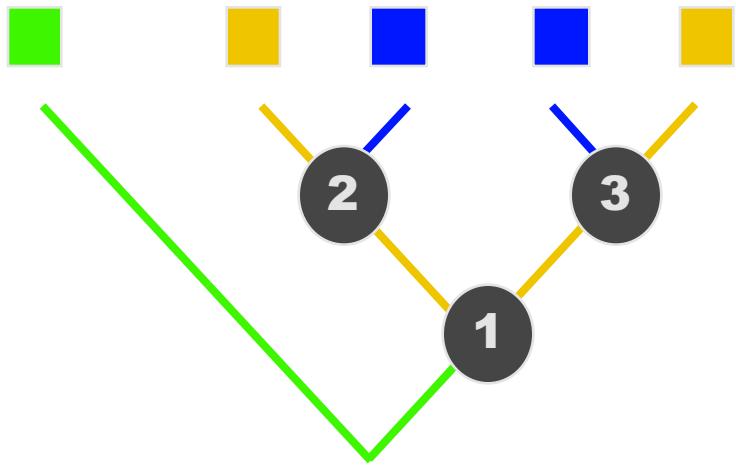
Shykoff & Widmer (1998)
Trends in Ecology and Evolution, 13, 158.

How do we build trees?

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

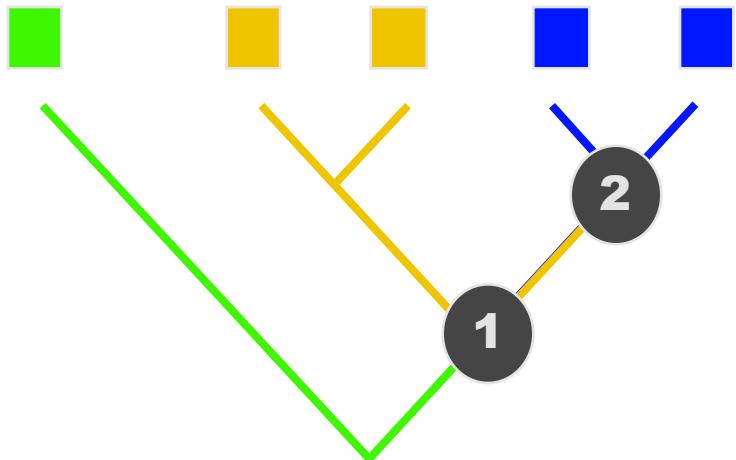


Maximum Parsimony



Aim to find “most parsimonious” tree

Smallest amount of evolution

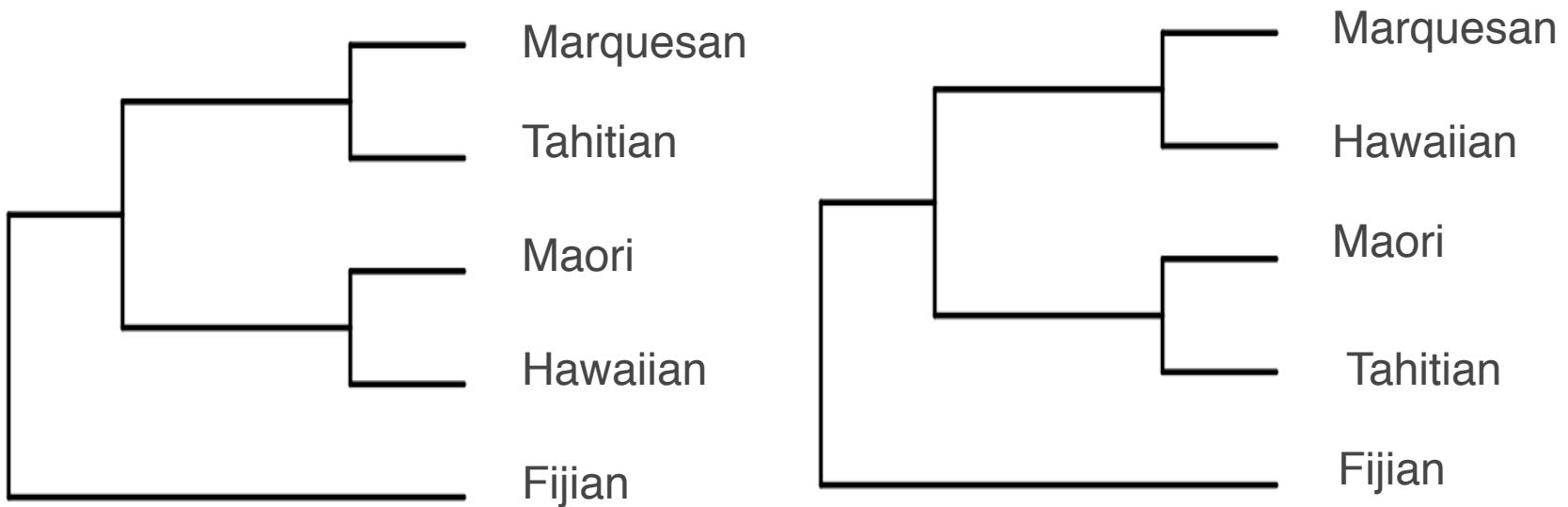


Unlikely that things should arise more than once
(i.e. cognates should innovate once)

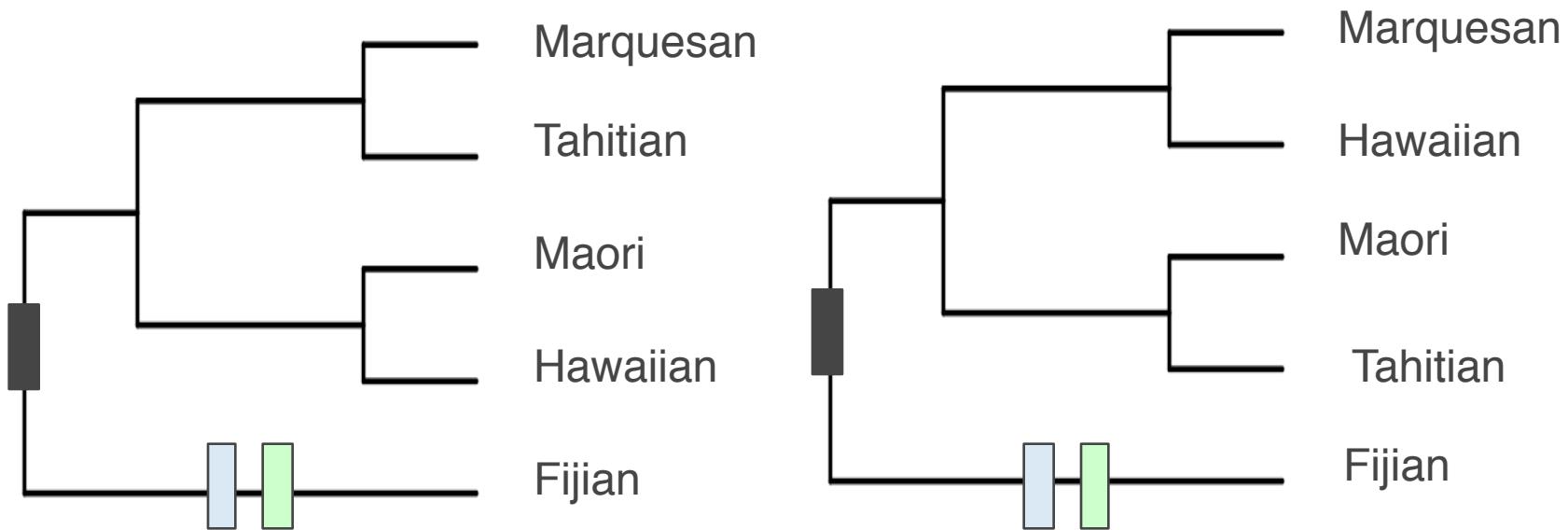
	Taboo	Blood	To Suck
Fijian	tabu	drā	sucu-ma
Tahitian	tapu	toto	ngote
Maori	tapu	toto	ngote
Hawaiian	kapu	koko	omo
Marquesan	tapu	toto	omo

	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



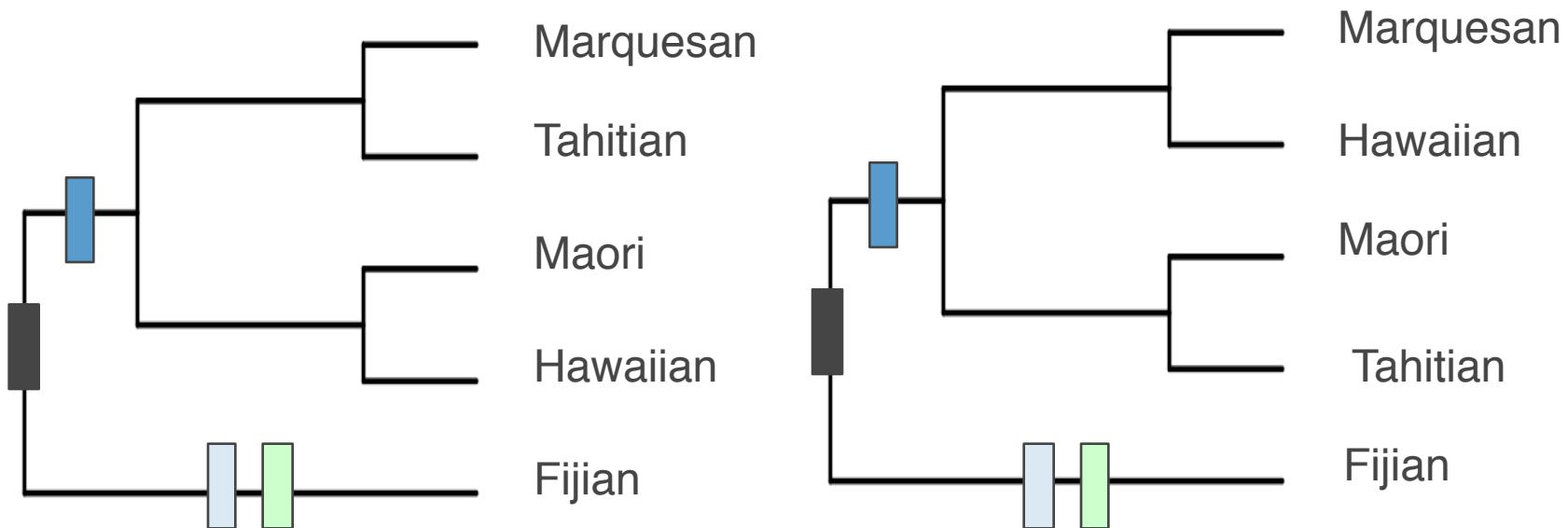
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=3

Length=3

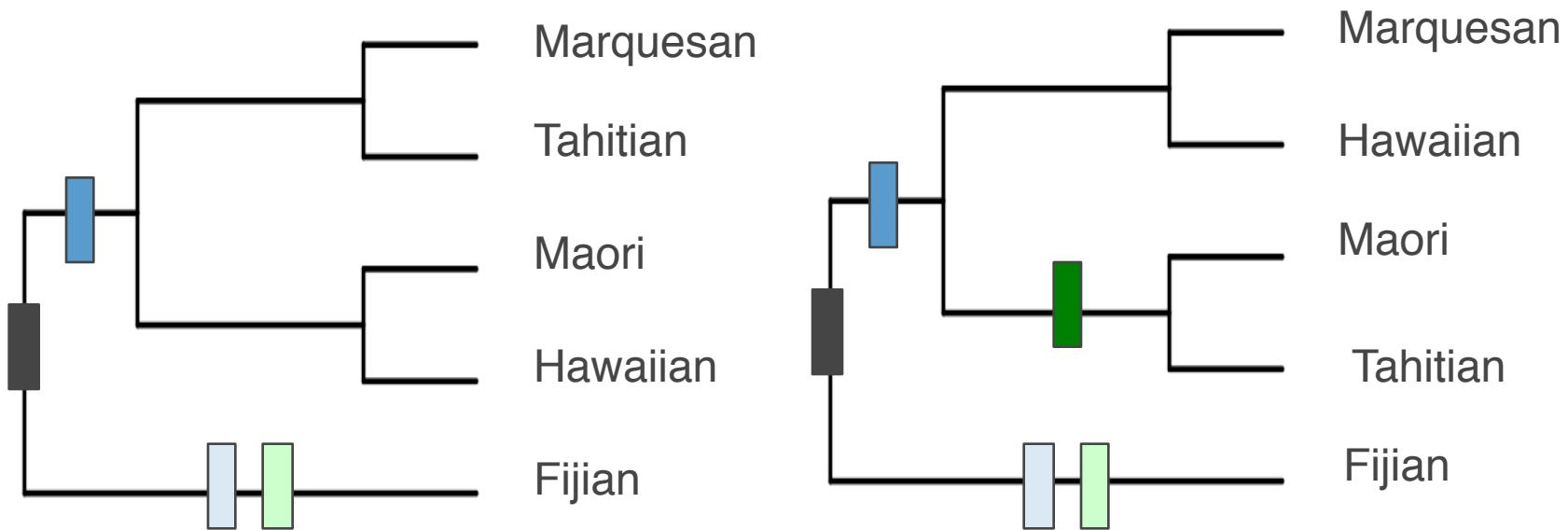
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=4

Length=4

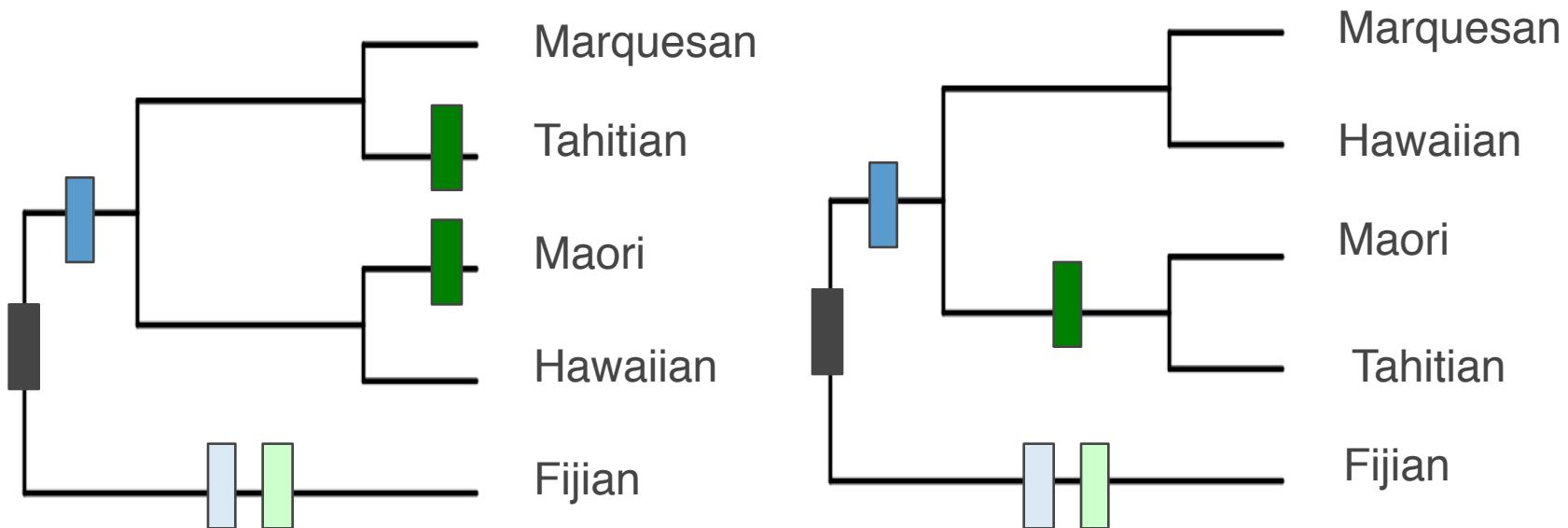
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



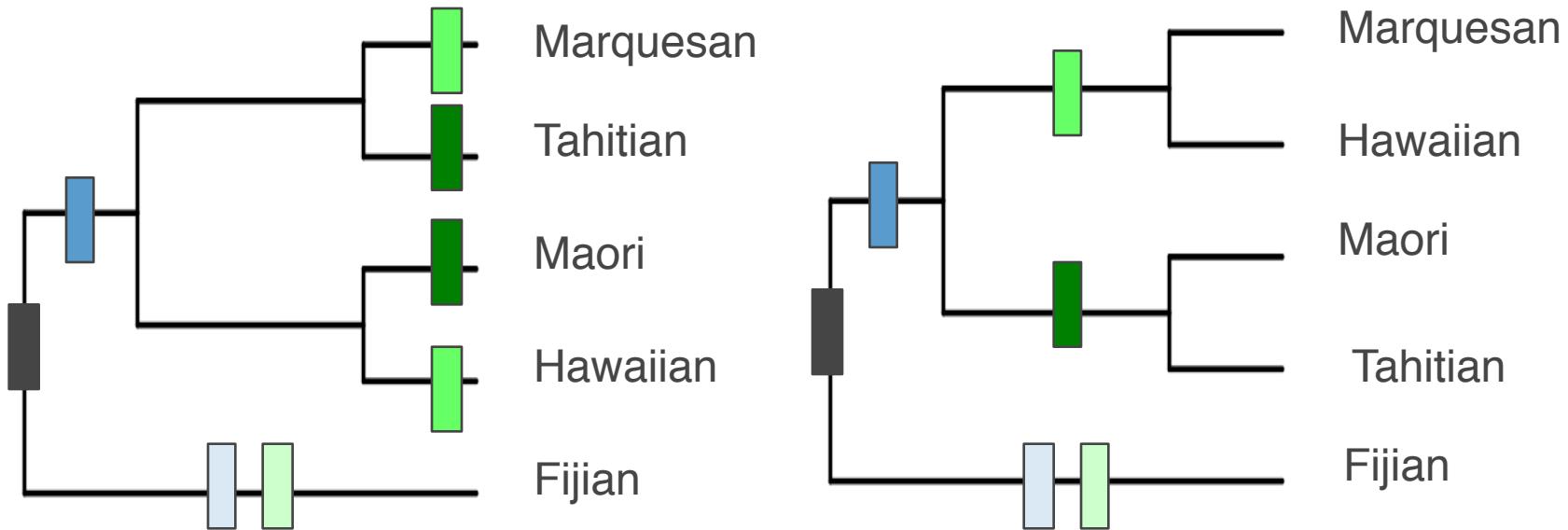
Length=4

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



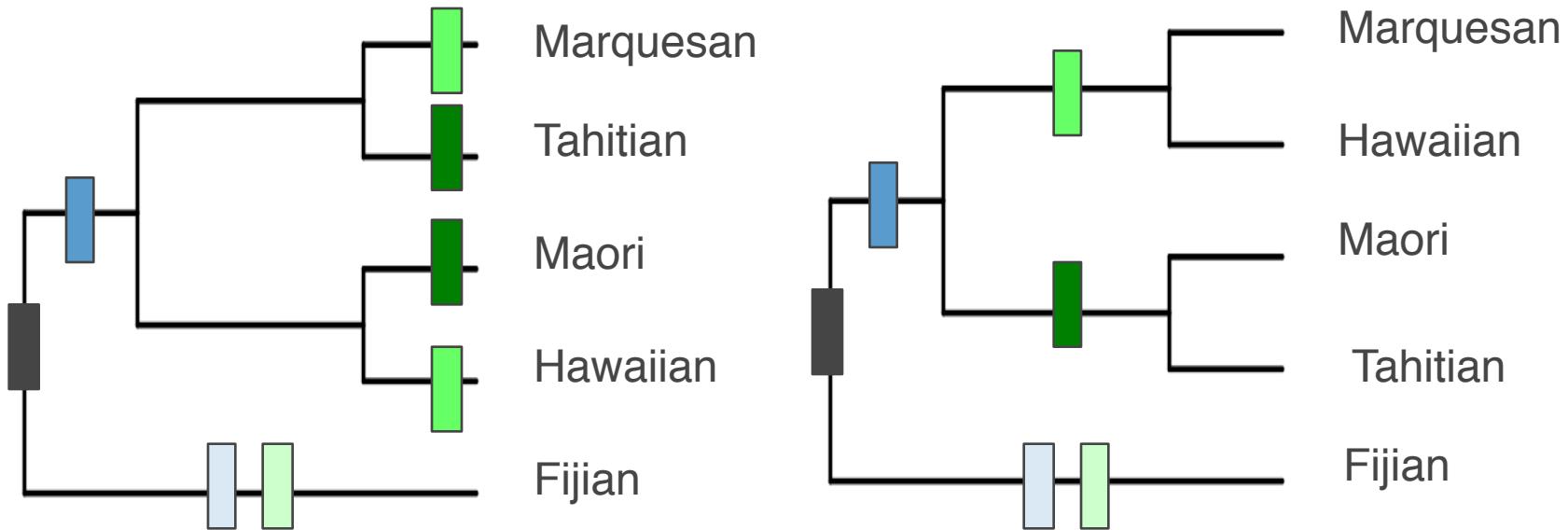
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



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



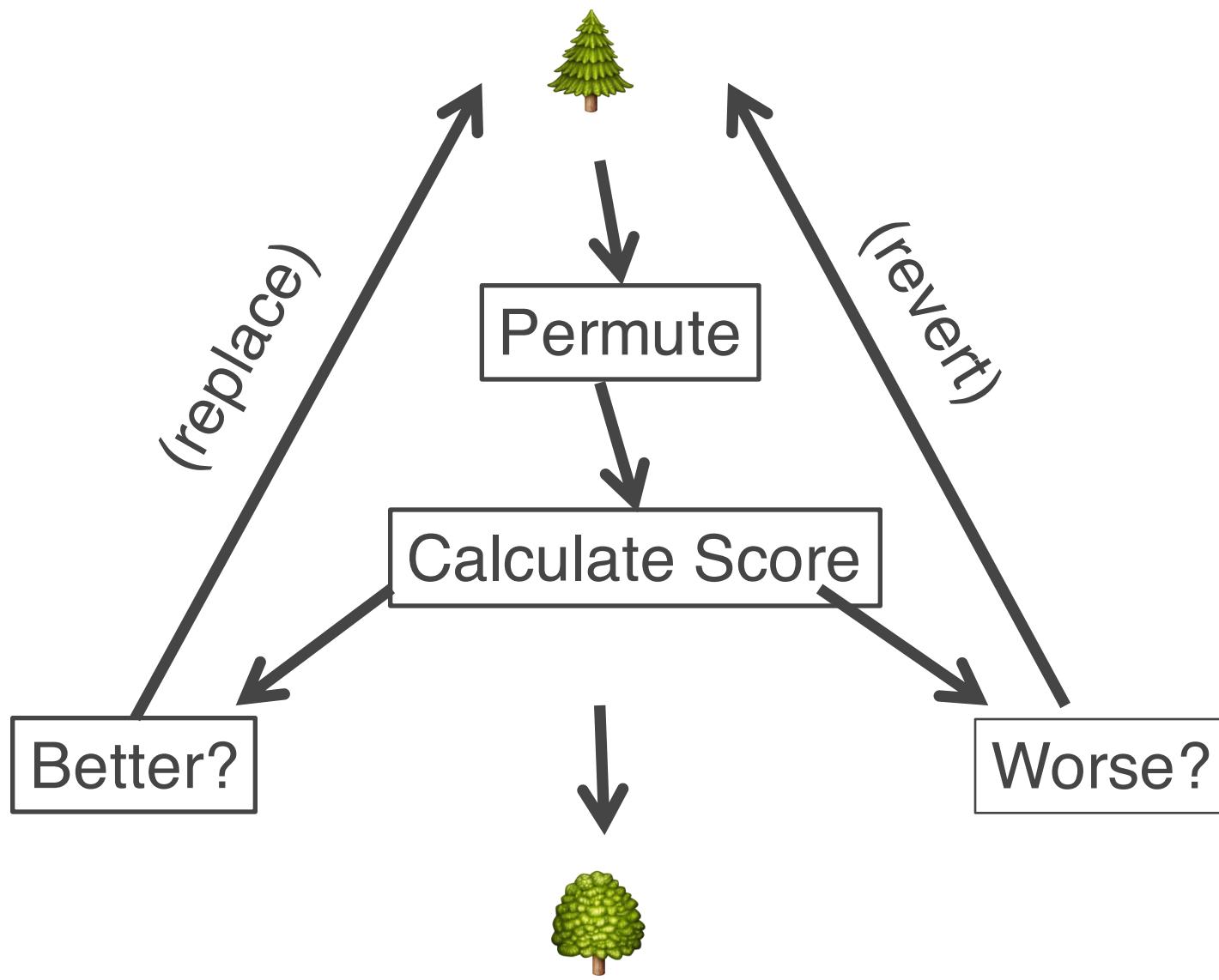
Length=8

Length=6



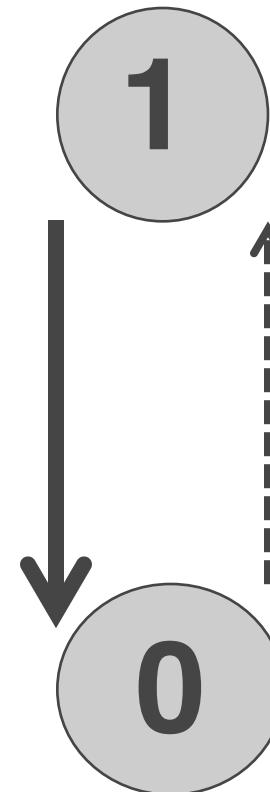
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1

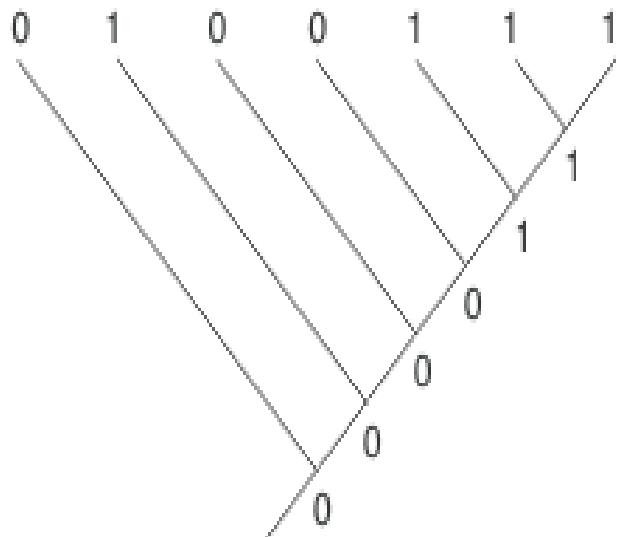
Algorithm



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.





$$L(a) = P(0 \rightarrow 0|b_1) \times P(0 \rightarrow 0|b_2) \times P(1 \rightarrow 1|b_3) \times P(1 \rightarrow 0|b_4) \times \\ P(0 \rightarrow 0|b_5) \times P(0 \rightarrow 0|b_6) \times P(0 \rightarrow 0|b_7) \times P(0 \rightarrow 1|b_8) \times P(1 \rightarrow 1|b_9) \\ \times P(1 \rightarrow 1|b_{10}) \times P(1 \rightarrow 1|b_{11}) \times P(1 \rightarrow 1|b_{12})$$

Site Likelihood(a) = $\left(\begin{array}{c} \diagdown \\ \diagup \\ \diagdown \\ \diagup \\ 0 \\ \diagdown \\ 0 \end{array} \right) \dots \times \dots \left(\begin{array}{c} \diagdown \\ \diagup \\ \diagdown \\ \diagup \\ 0 \\ \diagdown \\ 1 \end{array} \right) \dots \times \dots \left(\begin{array}{c} \diagdown \\ \diagup \\ \diagdown \\ \diagup \\ 1 \\ \diagdown \\ 1 \end{array} \right)$

Site Likelihood(a) = P(reconstruction 1) ... x ... P(reconstruction 5) ... x ... P(reconstruction n)

$L_h =$

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

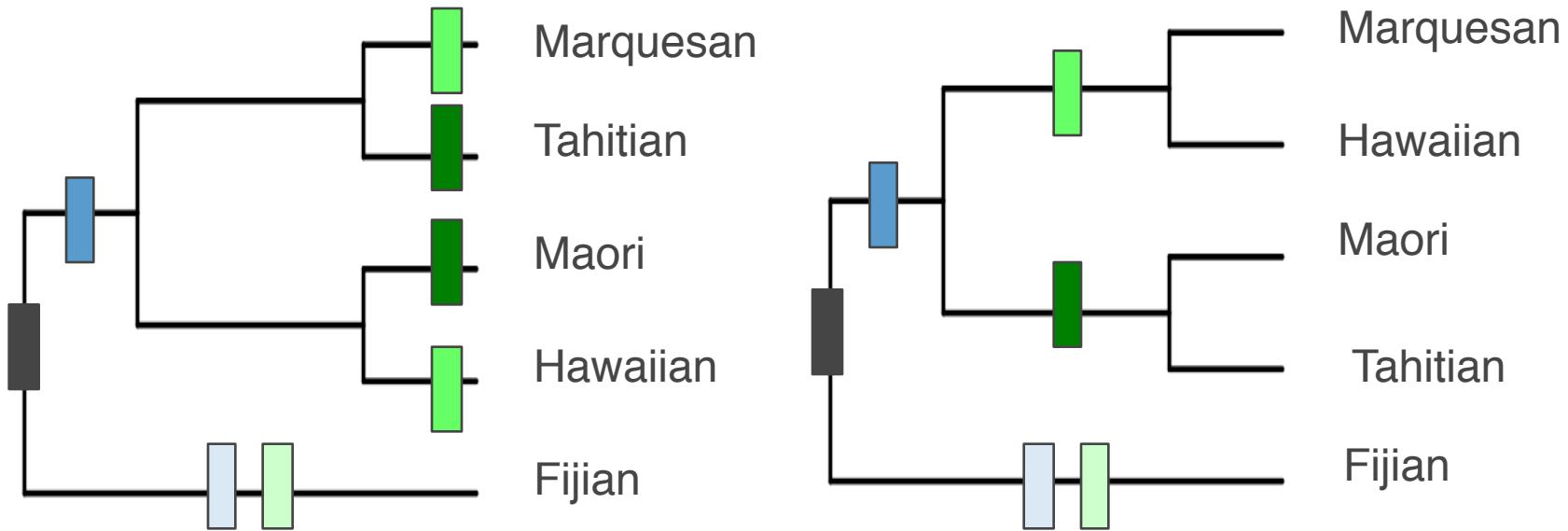
x

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

x

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

.... etc



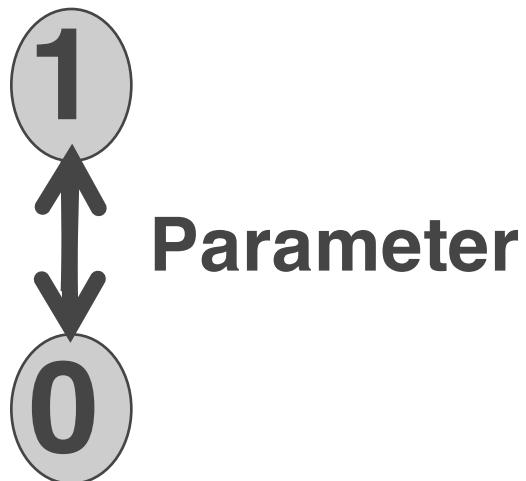
$$\ln(L) = -14.804$$

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

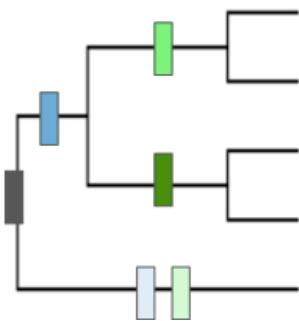
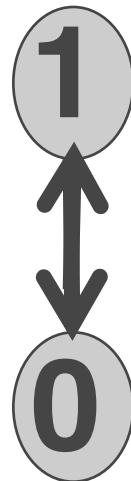


Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1

Models

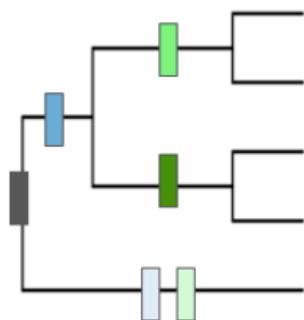
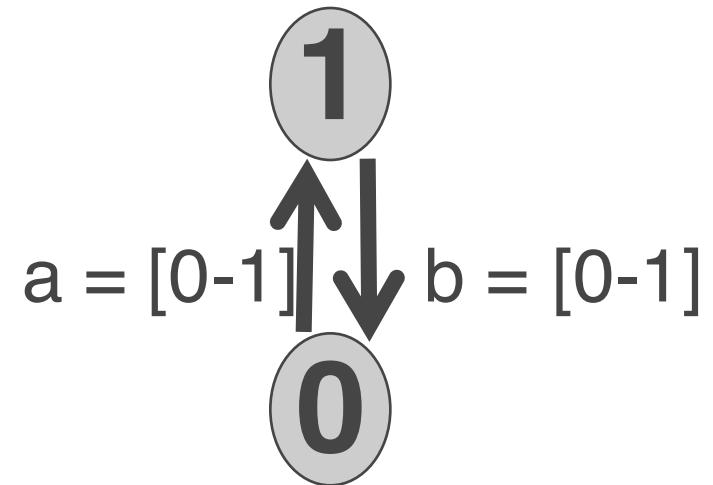
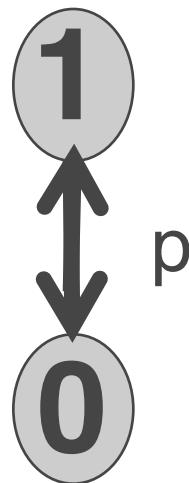


Models



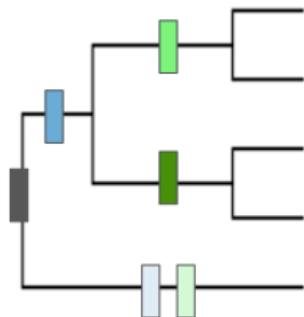
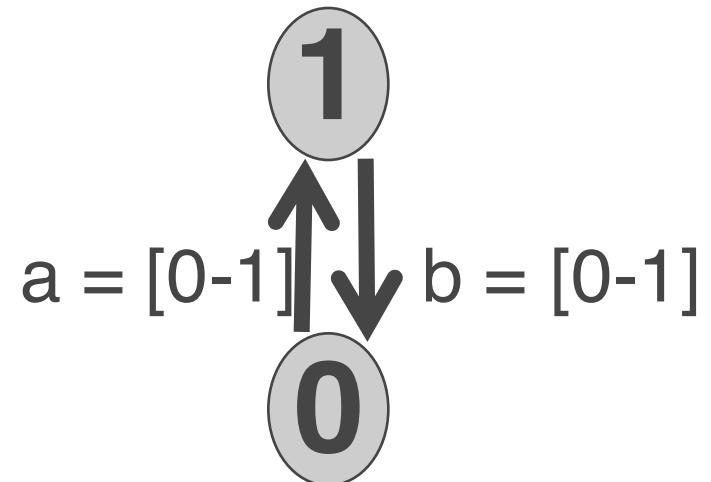
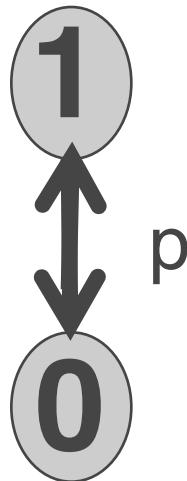
$$\ln(L) = -12.007$$

Models

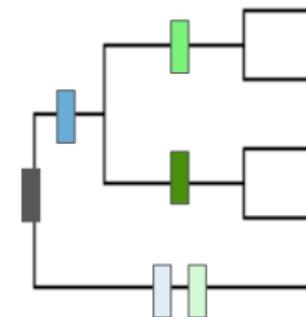


$$\ln(L) = -12.007$$

Models



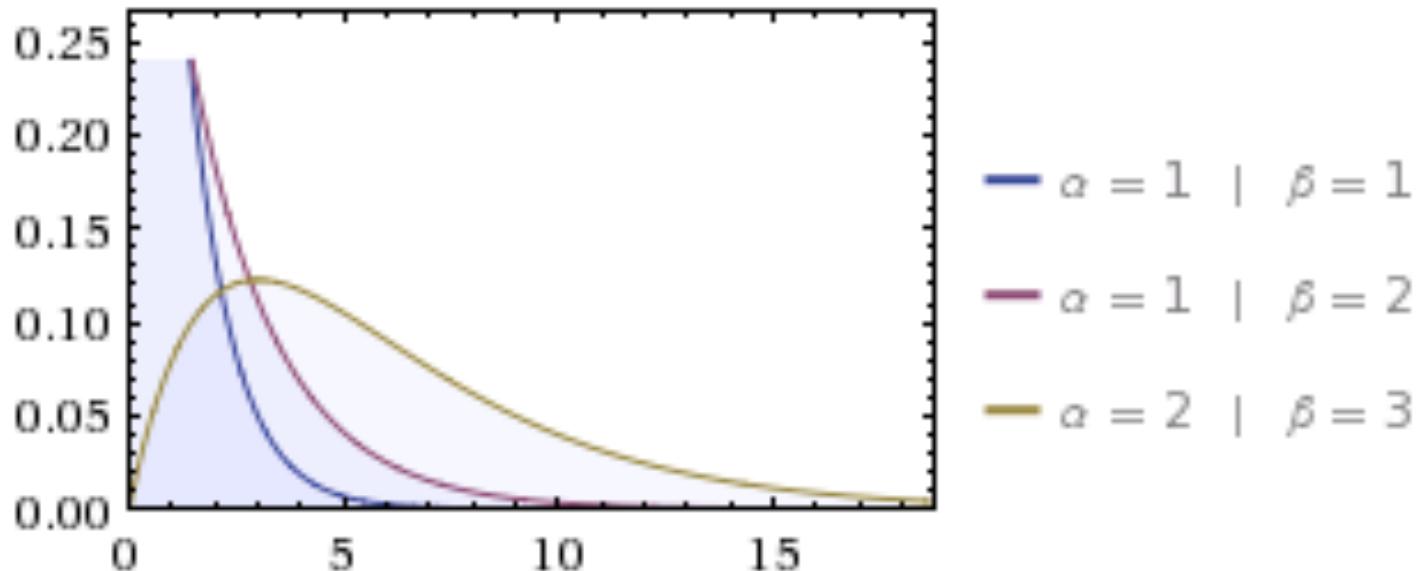
$$\ln(L) = -12.007$$

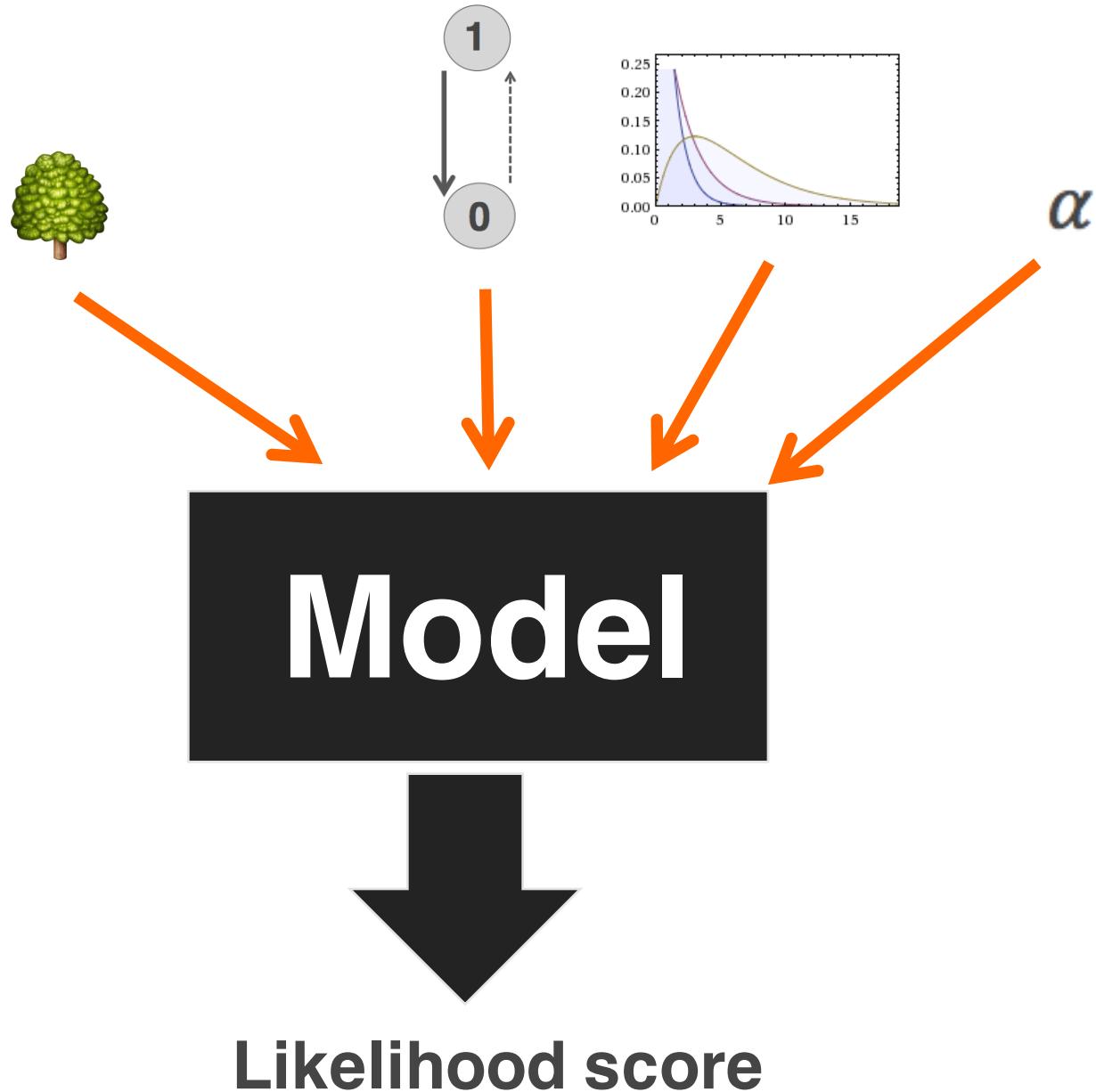


$$\ln(L) = -9.072$$

Rate Variation - Gamma

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







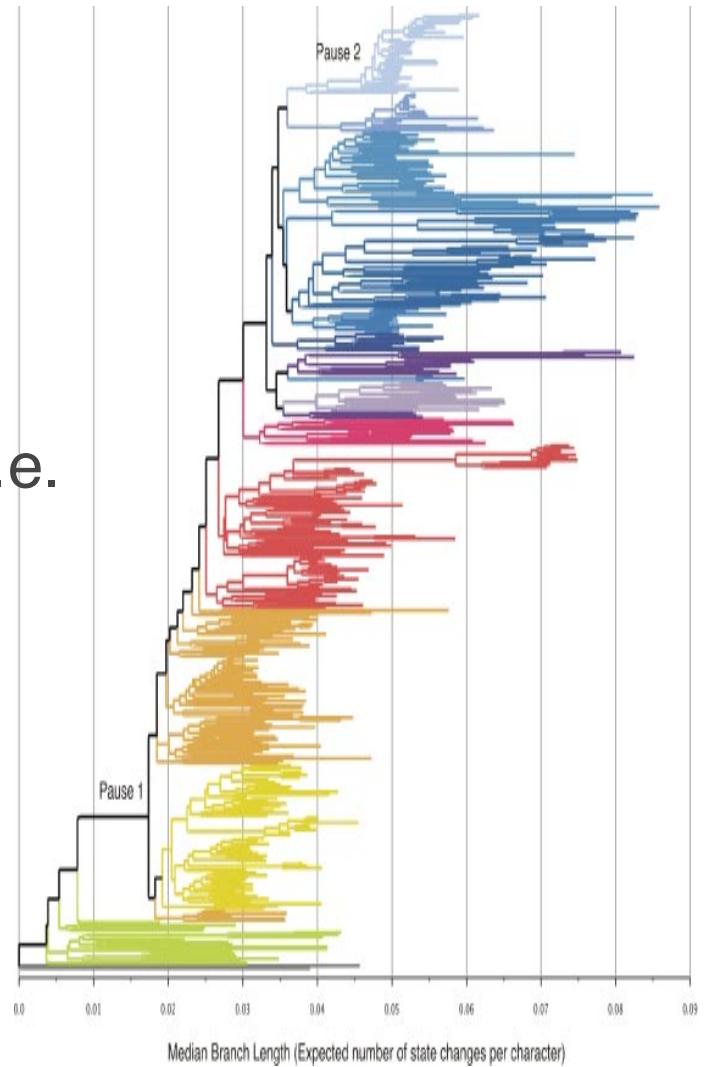
Linguists don't do dates.

Phylogenetic Dating

ML: 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 at time Z
- Smooth rates over these calibrations

Strict Clock

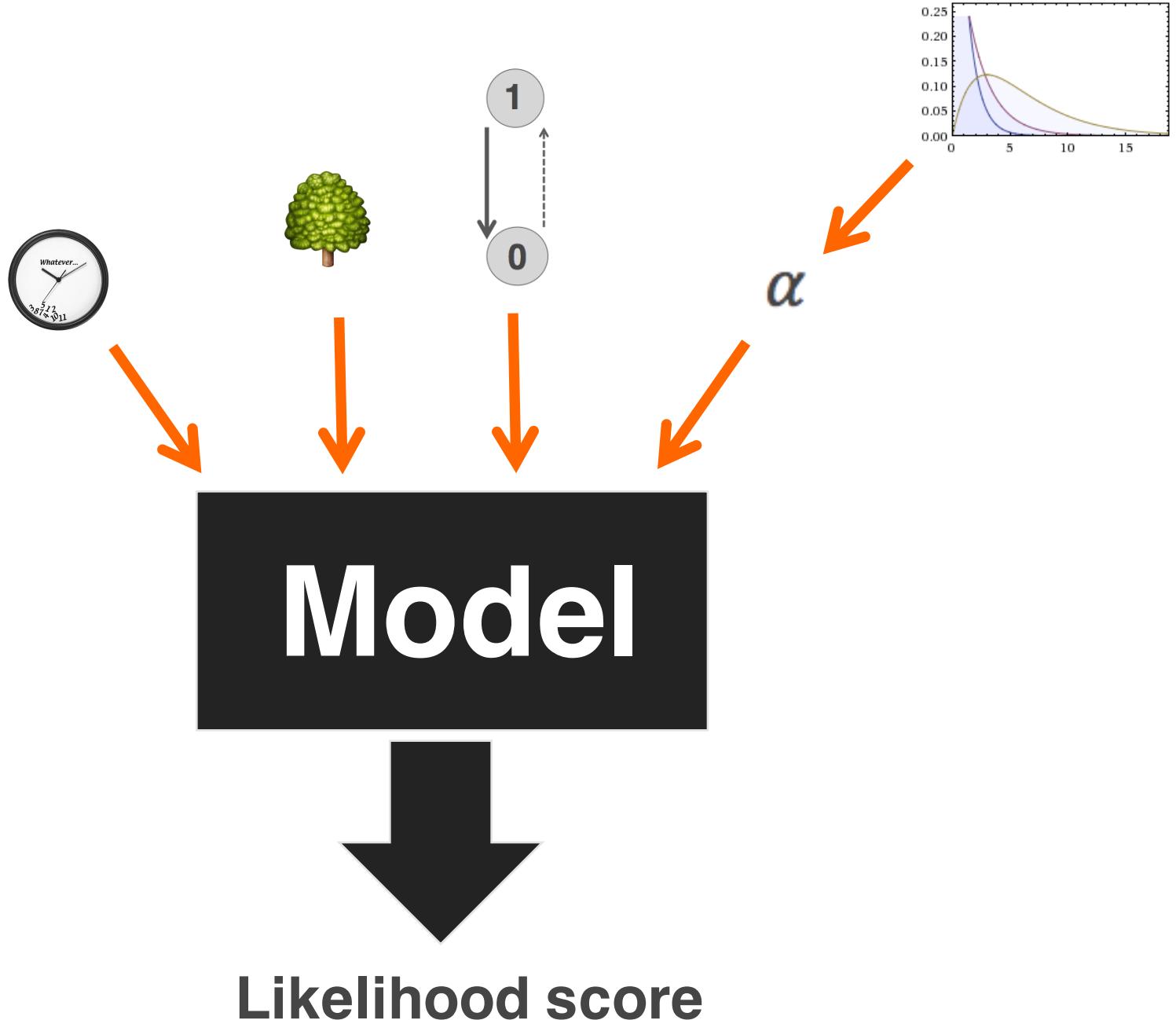


- One rate for all languages.
- No variation
- = glottochronology

Relaxed Clock



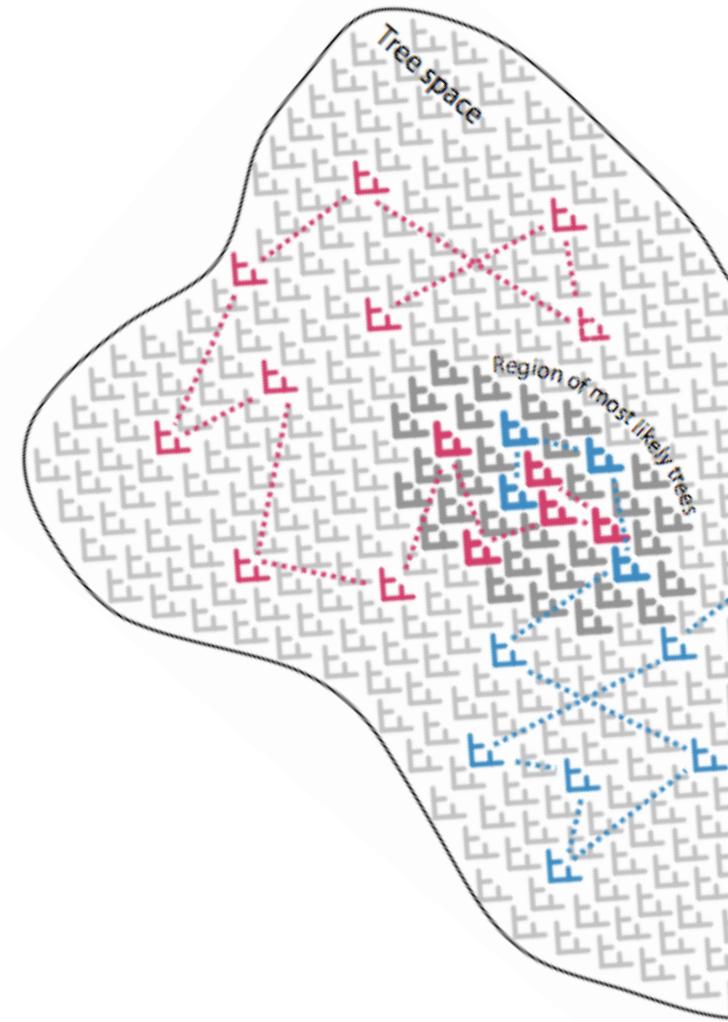
- Allows rate to vary across branches.
- Rates are drawn from a parametric distribution with parameters estimated from the data

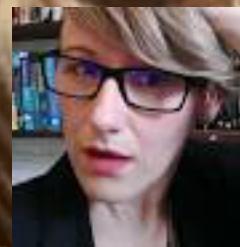


Bayesian Phylogenetics

1. Data & Model & Tree
2. Calculate the **Likelihood** of that tree
3. Modify The Tree or a Model Parameter
4. Repeat (MCMC “Walk” through treespace)

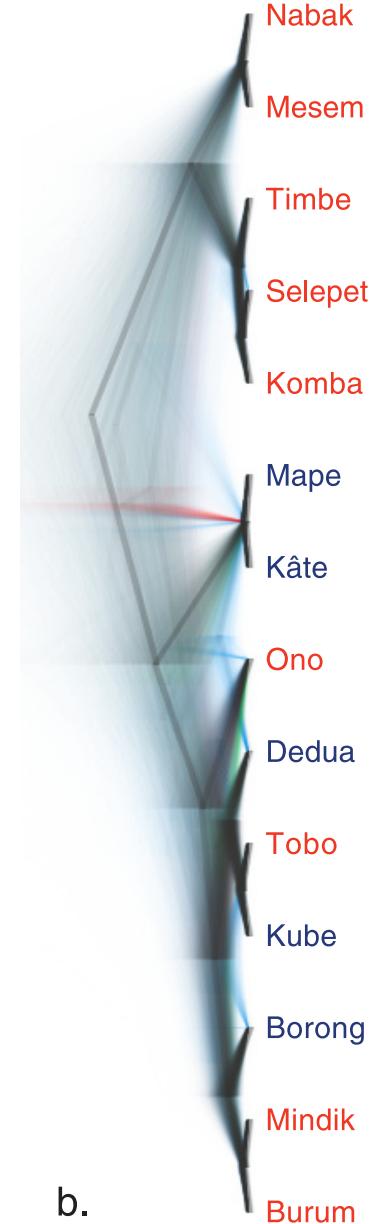
-> samples best fitting trees from all possible trees.
-> "Posterior Probability Distribution"





Posterior Probability Distribution

- Not just 1 tree but sample
- Uncertainty, conflicting signal come out.
- Simulations: robust to reasonable levels of borrowing (15% / 1000 y)



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).
- Currently are incredibly powerful tools for answering certain types of questions (not all!).
- Range of methods available. Bayesian P.M. 
- Fun.

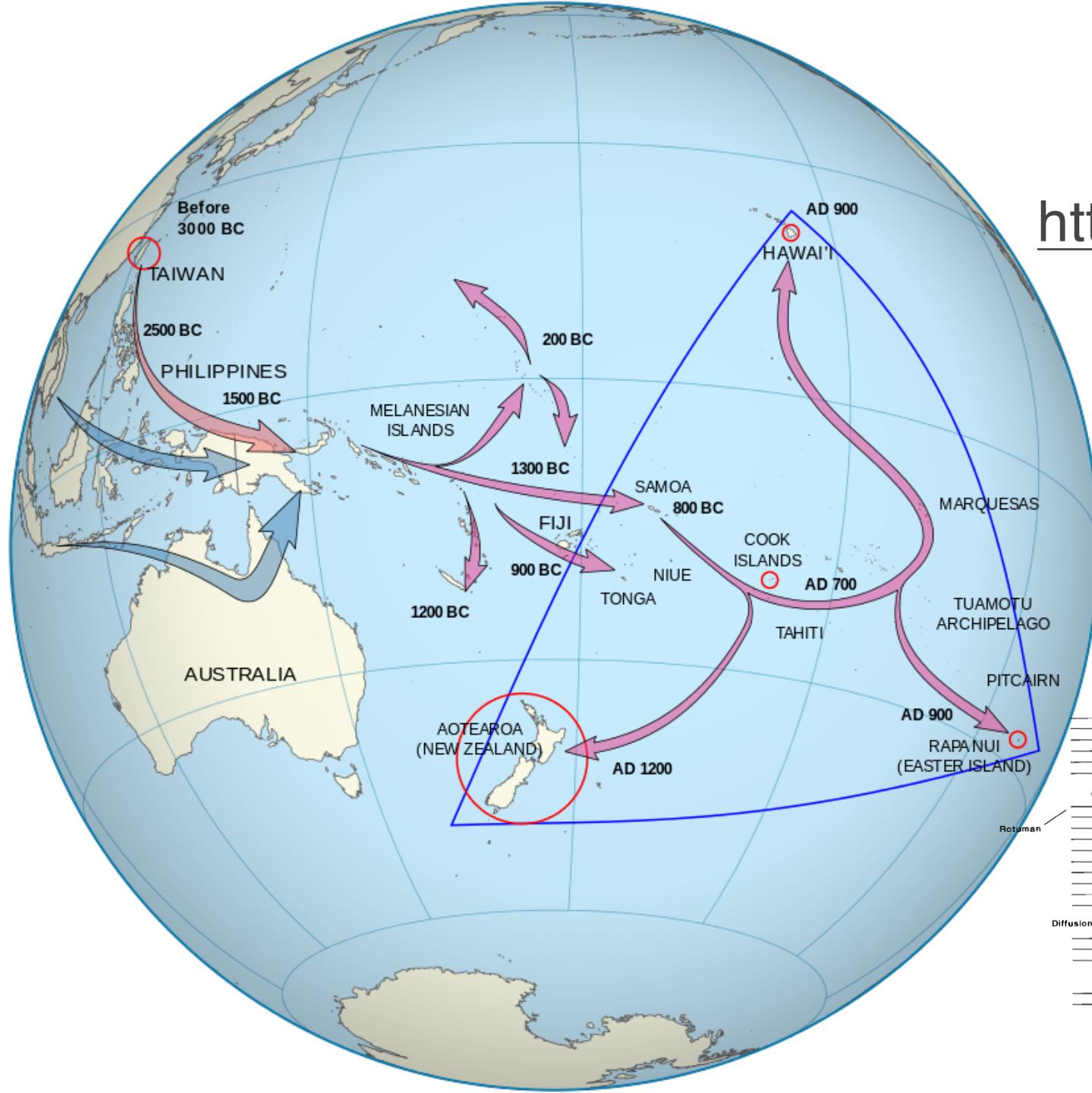
Note: H.O.P by Fiona on Bayesian methods

Questions?



DATA

<http://bit.ly/252YNMj>



Proto Central Pacific dialect chain		Tokelau Fijian–Polynesian dialect chain	
CP dialects of Western & Central Fiji	Tokelau Fijian	S	N
Proto-Polynesian		PTO	PNP
Diffusion of innovations across Fijian dialect chain			
Western Fijian dialect chain	Eastern Fijian dialect chain	EFU,EUV	PEC
			TOK.TUV
		TON NIU	SAM
			PEP
		WFU,WUV,EFU, TIK,REN,EUV	
			KAP.SIK,TOK, NUK,ONG, TUV, PCE