

# Phylogenetics & Tree-Thinking

Simon J. Greenhill



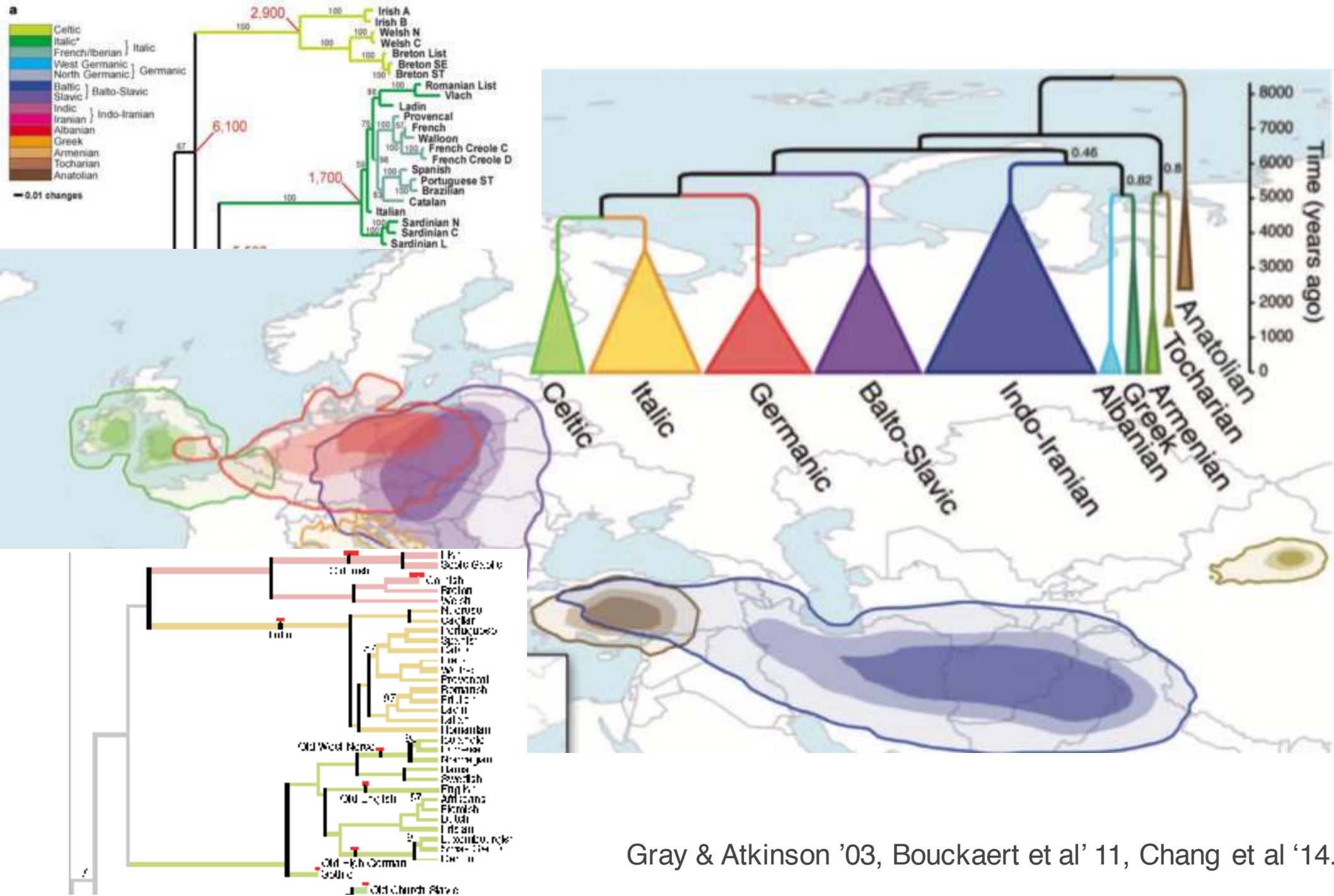
ARC CENTRE OF EXCELLENCE FOR  
**THE DYNAMICS OF LANGUAGE**

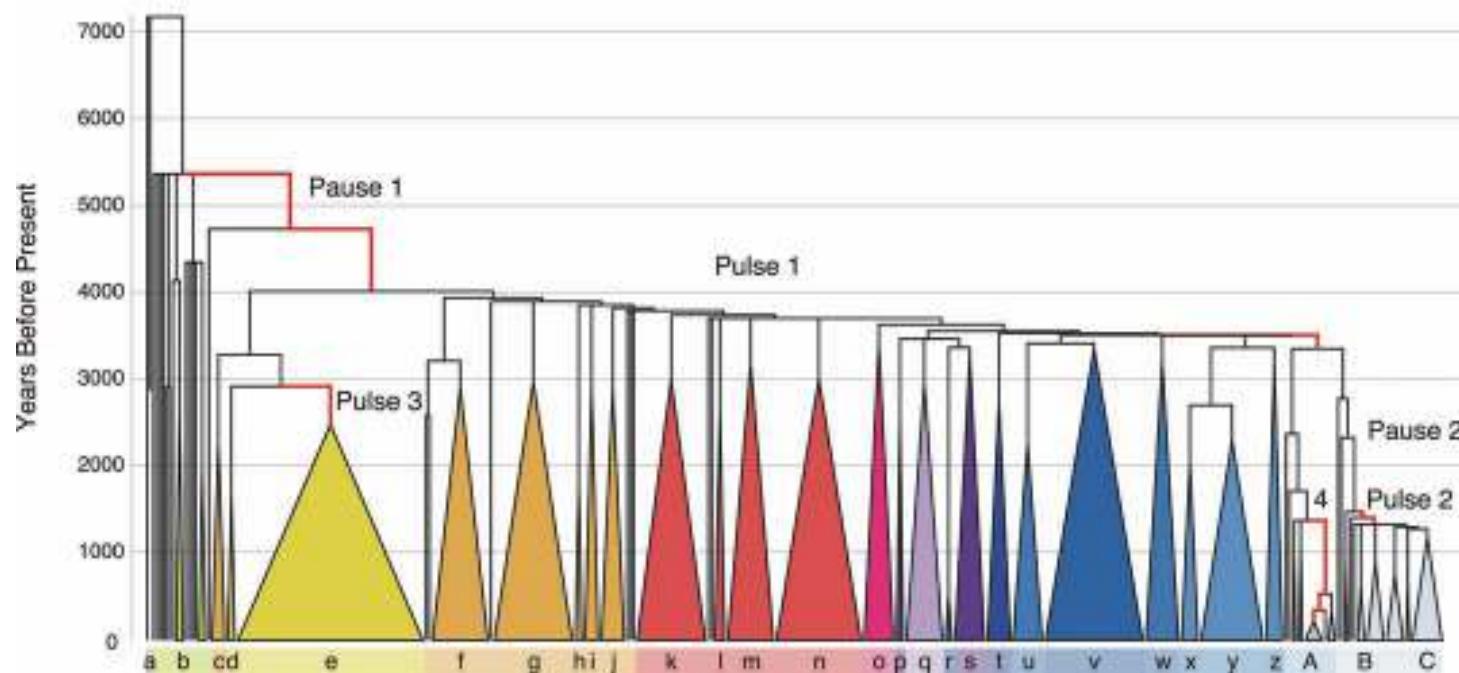
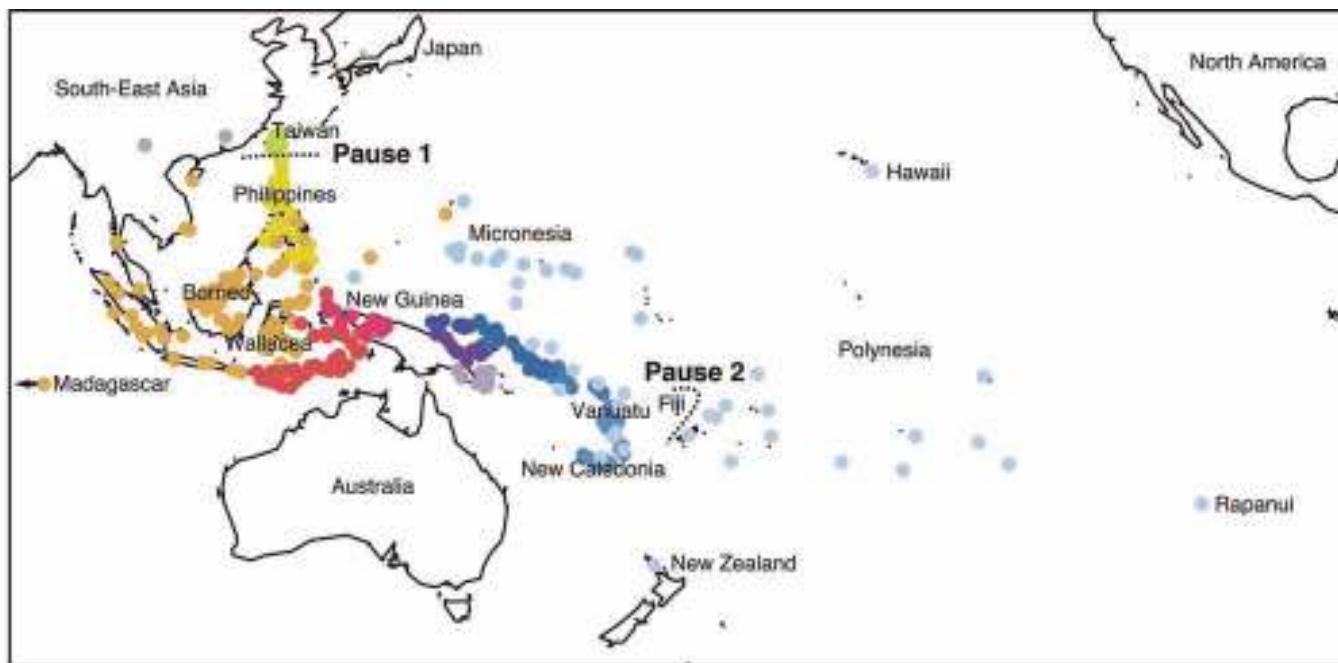


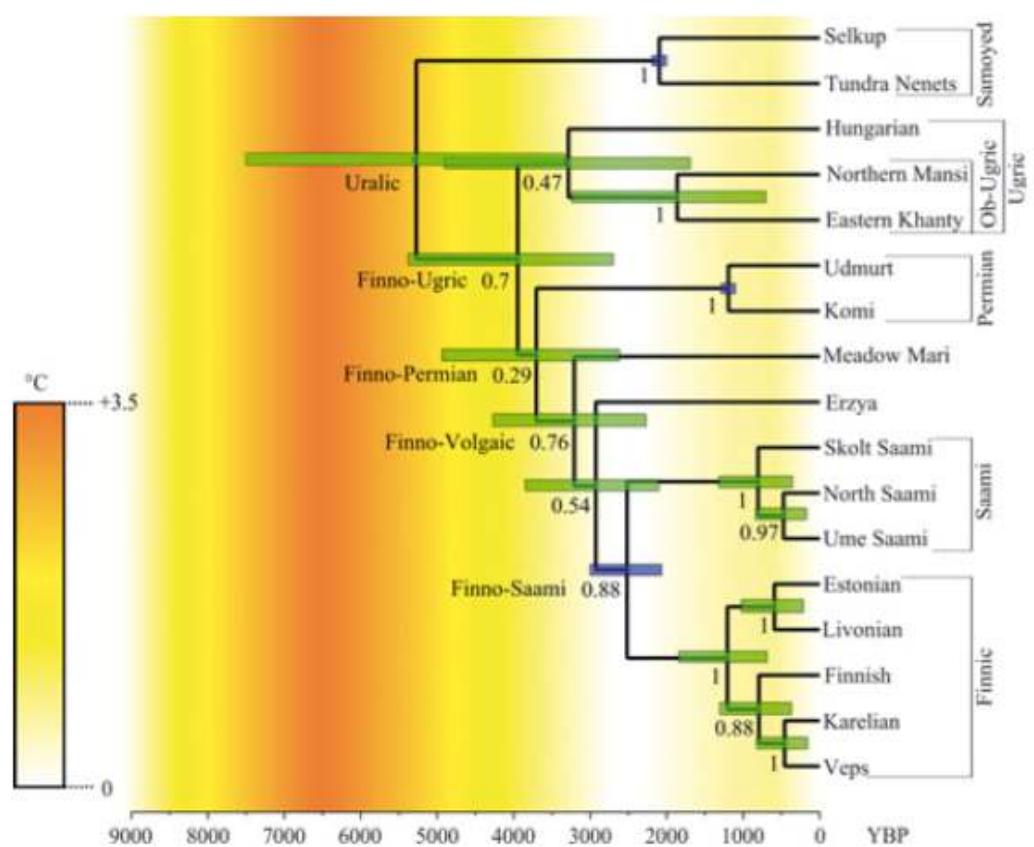
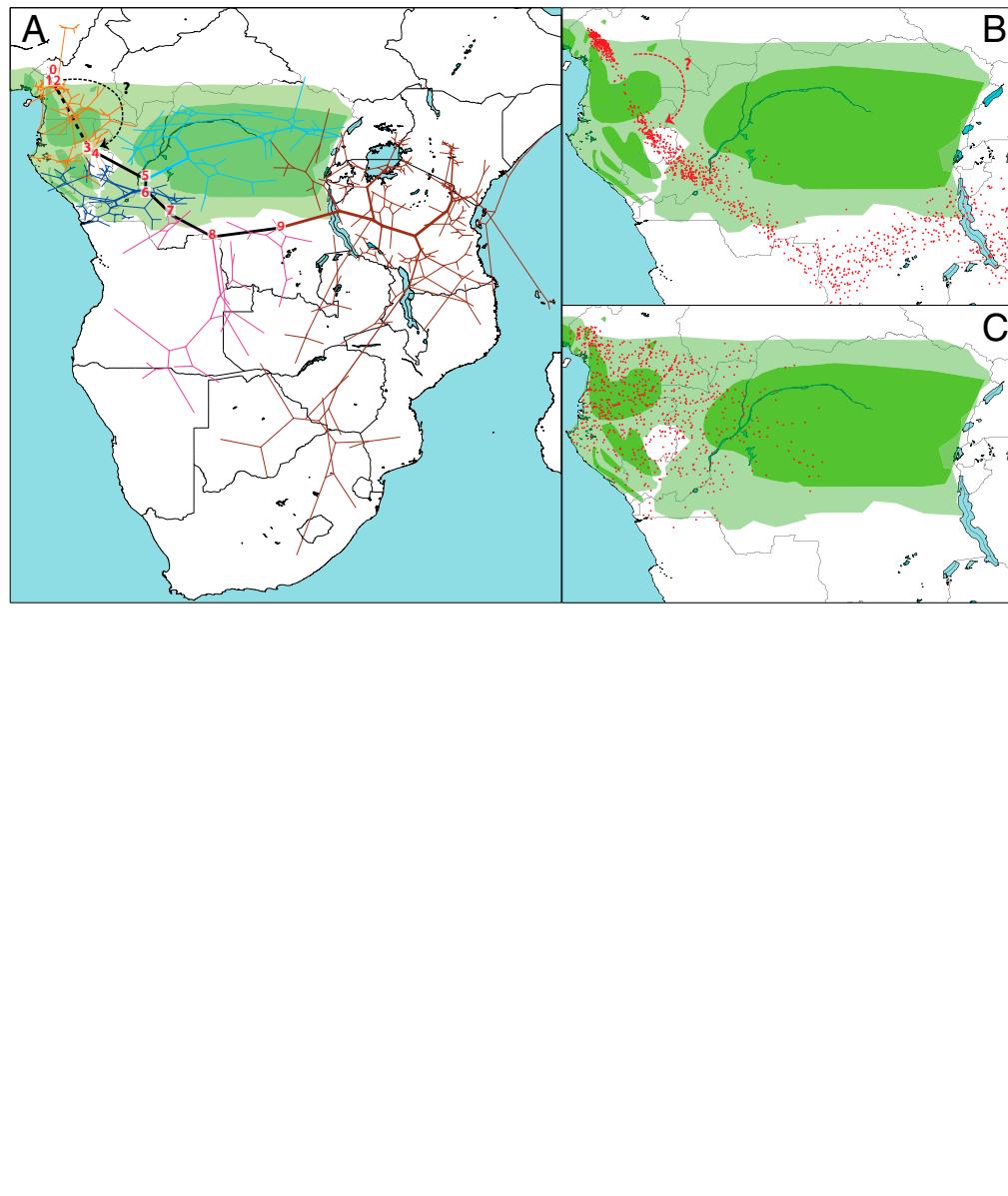
Max Planck Institute for the  
Science of Human History

# Phylogenetics & Tree-Thinking

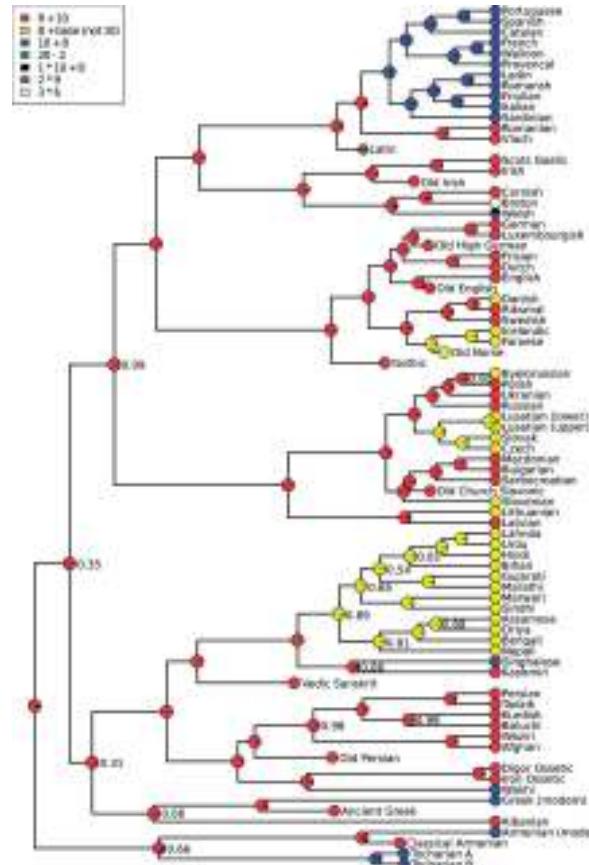
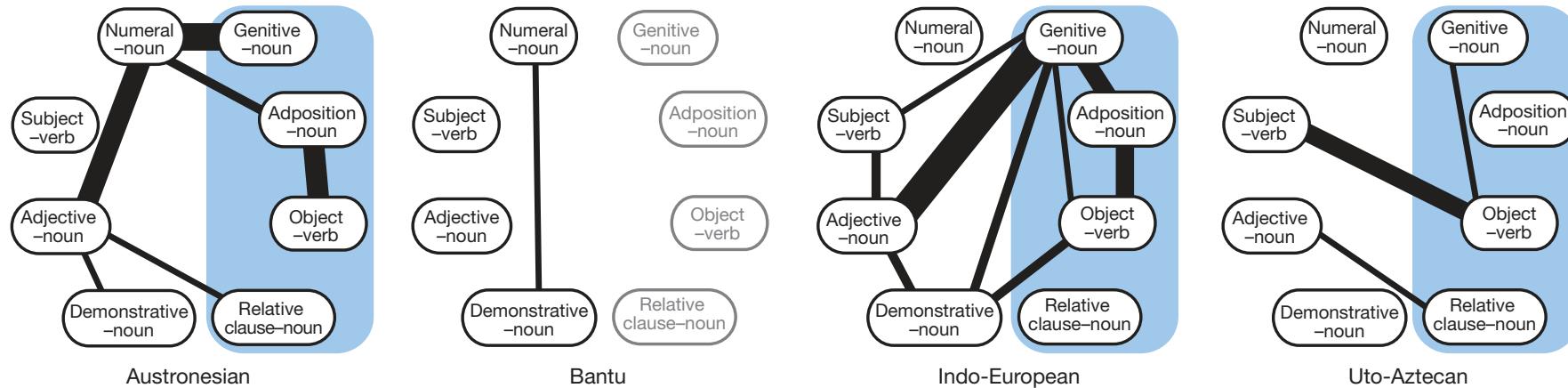
Range of methods in a robust, statistical/inferential framework for **testing** evolutionary hypotheses.



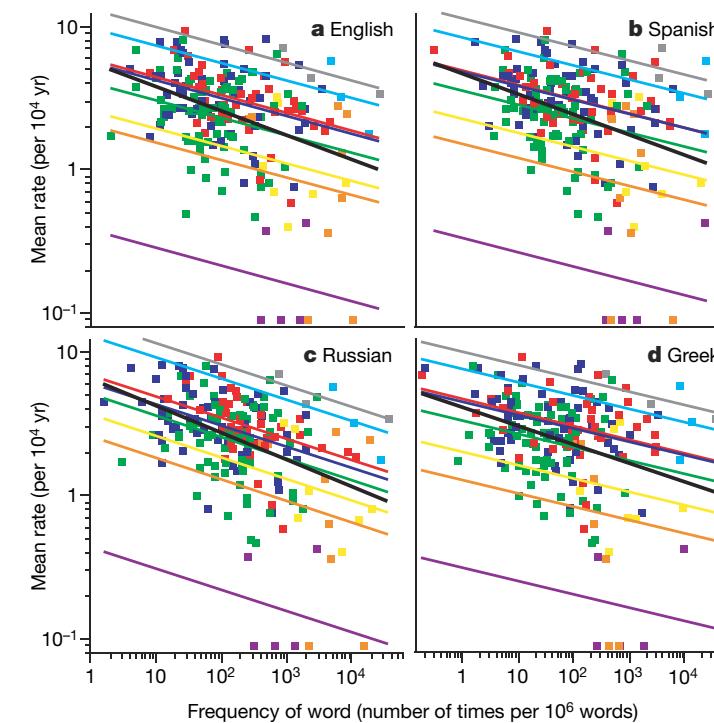




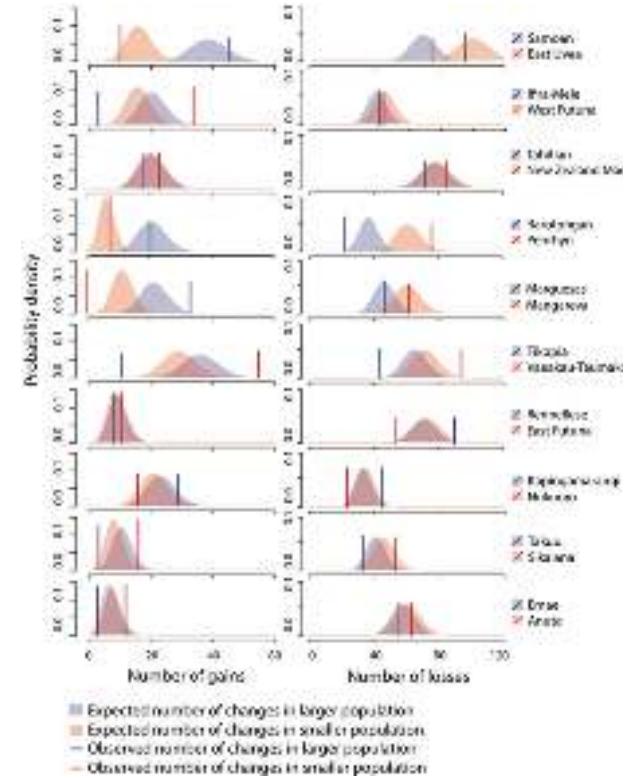
Dunn et al. '15



Calude and Verkerk '16

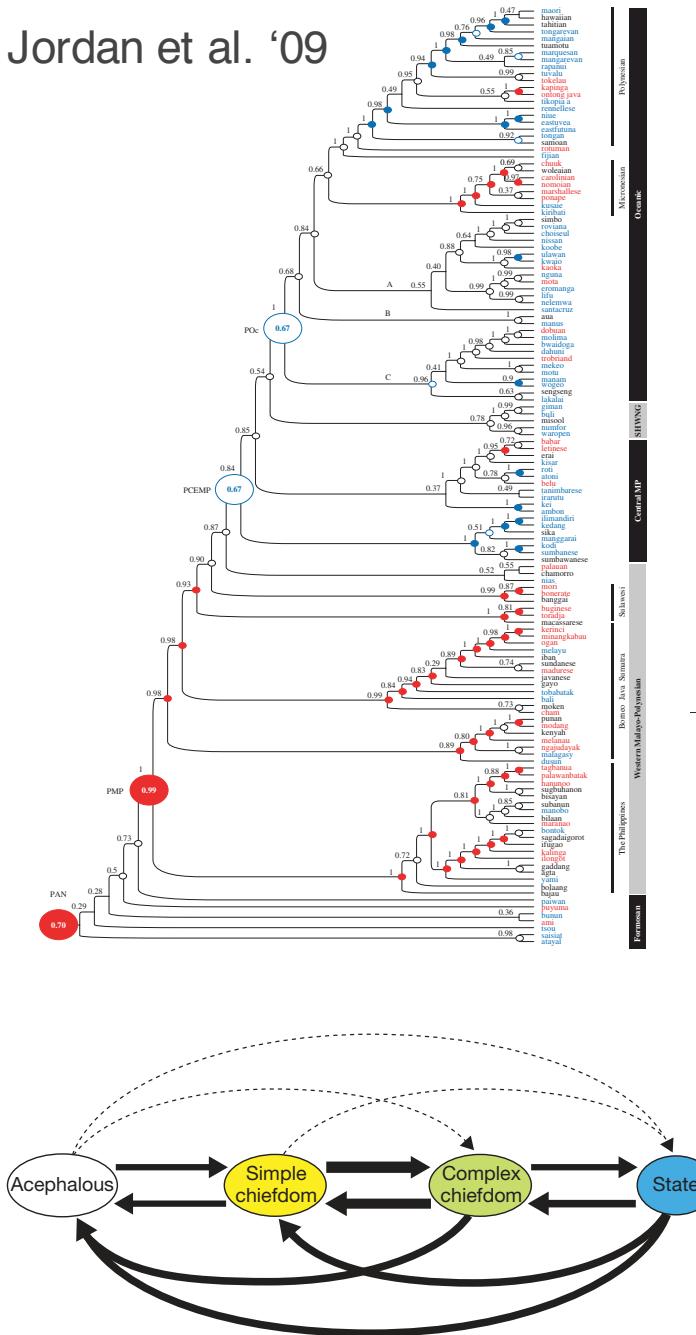


Pagel et al '07

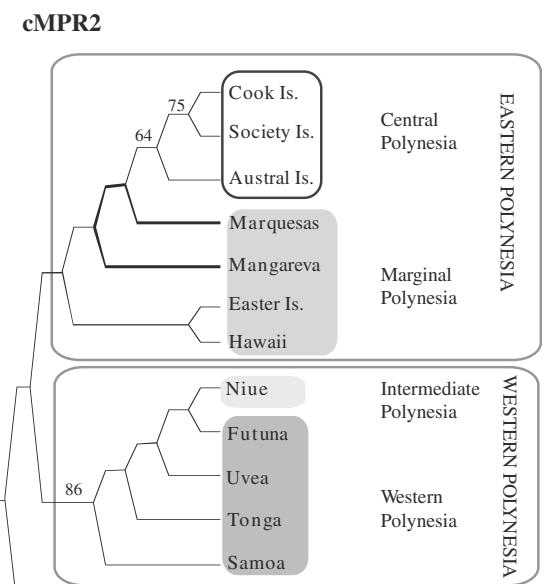


Bromham et al. '15

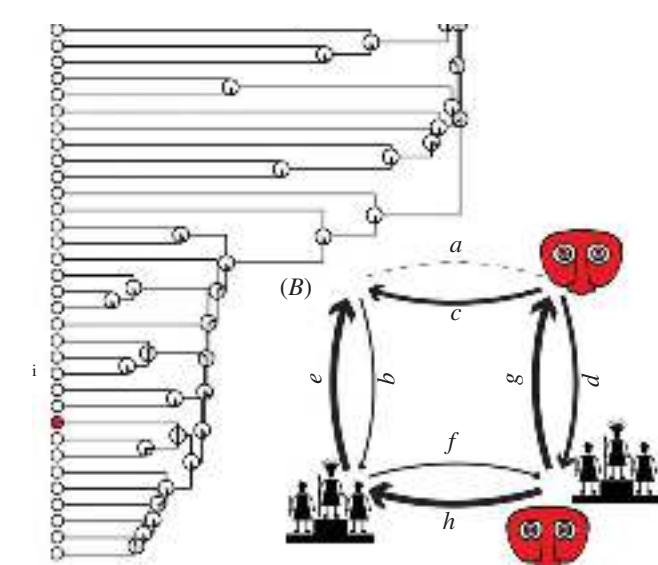
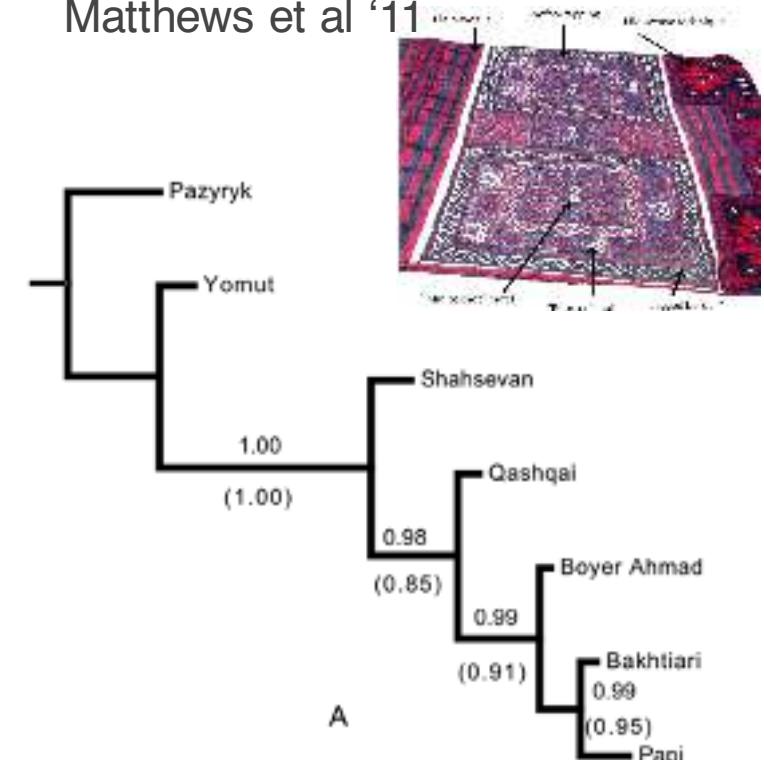
Jordan et al. '09



Larsen '11



Matthews et al '11



Da Silva & Tehrani '16

Currie et al. '10

# Controversial

“most vibrant stream  
of contemporary  
linguistics”

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

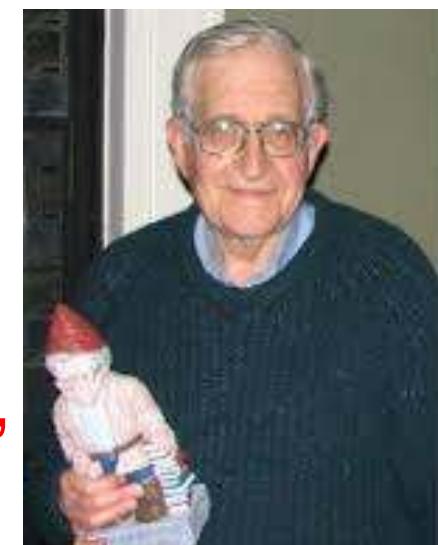
“languages and biomolecular  
sequences evolve in very  
different ways”

“more questions than  
answers”

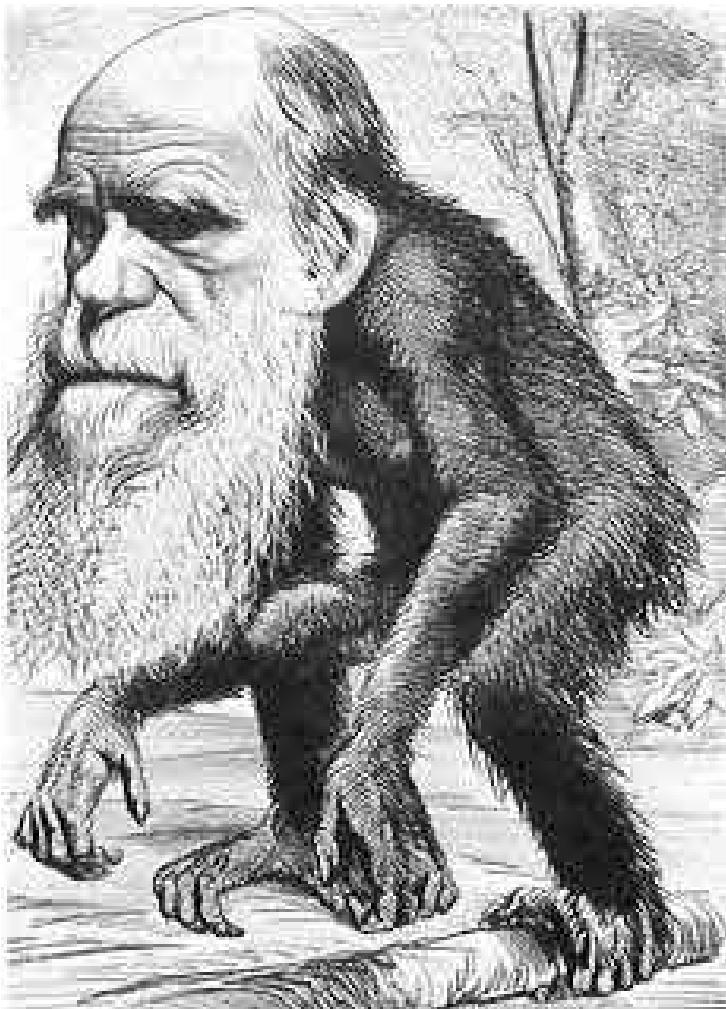
“utter bollocks”

“biggest intellectual fraud  
since Chomsky”

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



# What is evolution?



Variation

Heritability

Differential survival

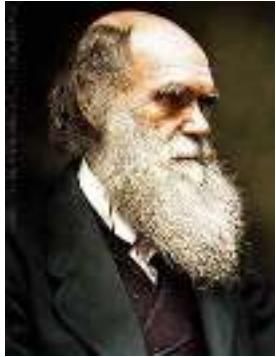
When and where did  originate?

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

How are  related to other 's?

What **processes** shaped 

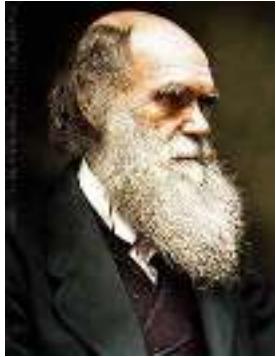
Can we infer what  were in the **past**?



## 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..."

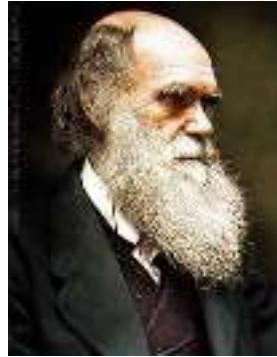
"The formation of different languages and of distinct species, and the proofs that both have been developed through a gradual process, are **curiously parallel**"



## 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)

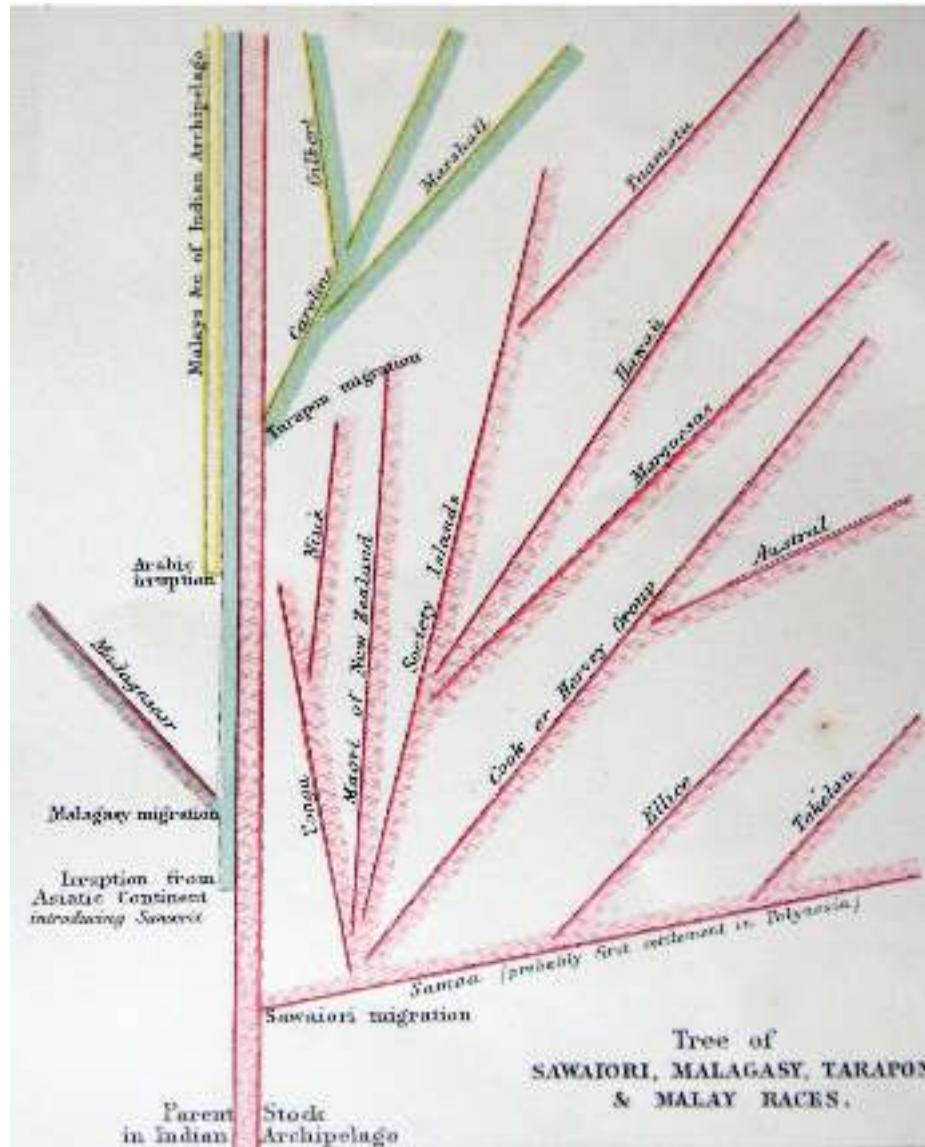
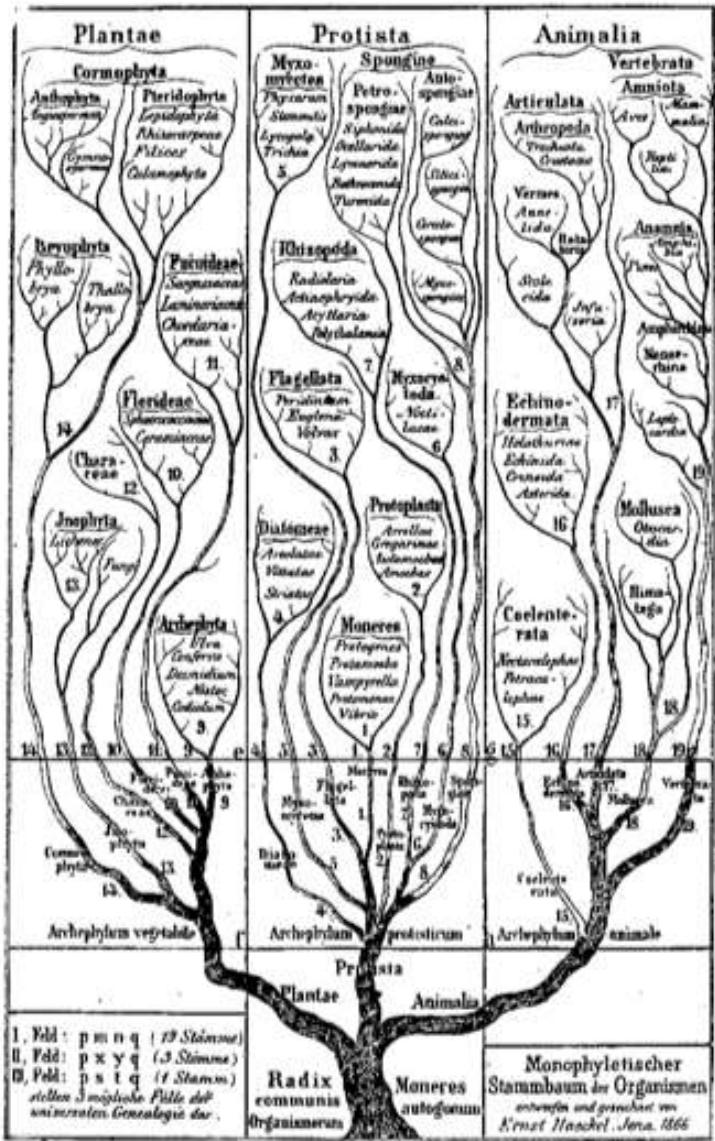
*Darwinism Tested by the Science of Language*

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

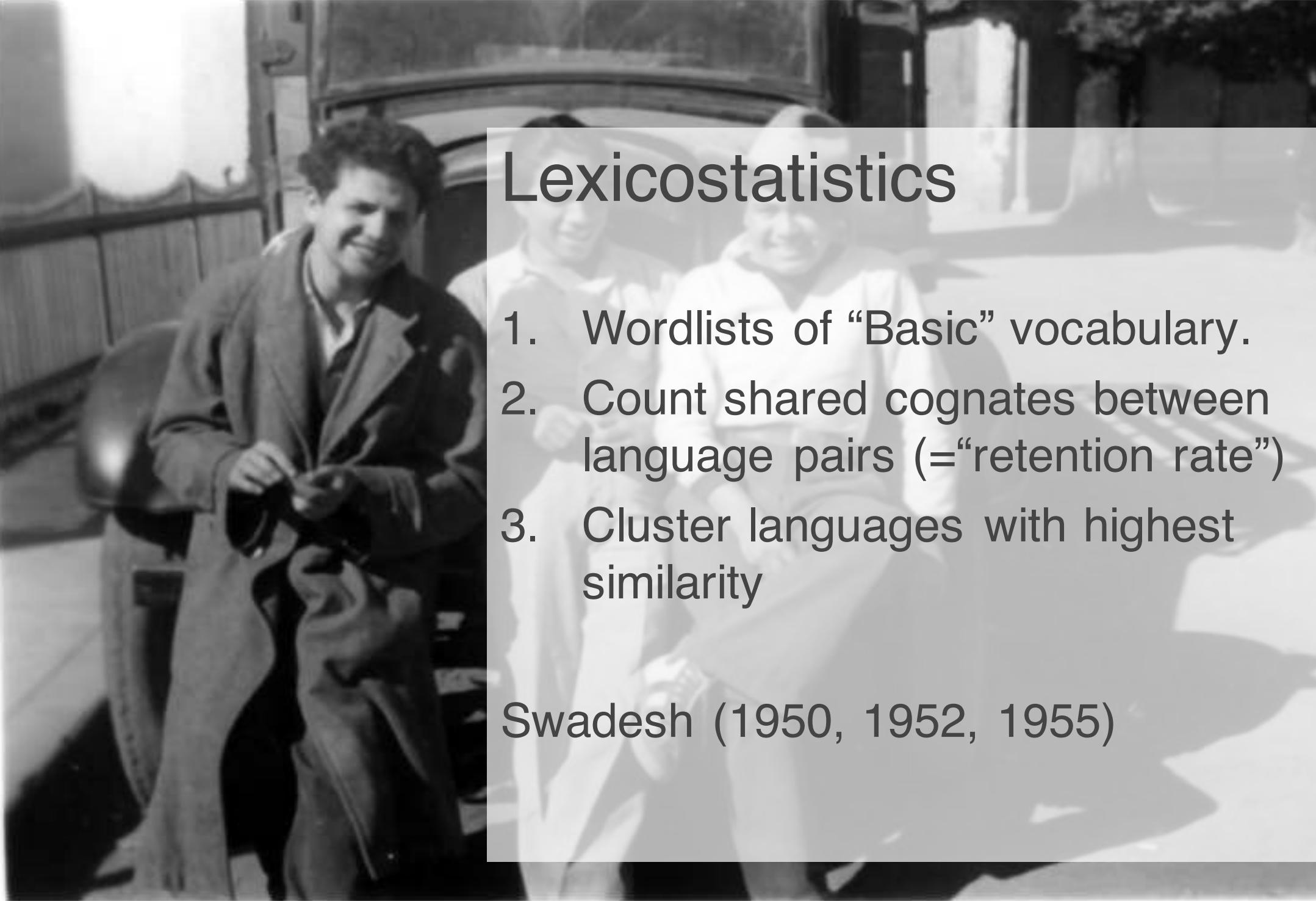


## Brugmann (1884)

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

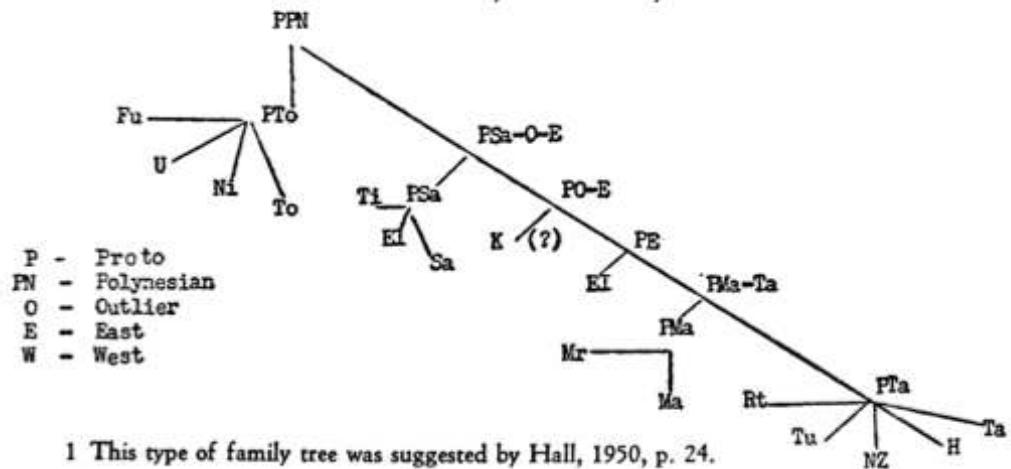
TABLE 2  
Polynesian cognate percentages

T	Hu	To	Ti	K	Ma	Pa <sup>1</sup>	Pu <sup>1</sup>	Pa <sup>2</sup>	Pu <sup>2</sup>	K	Ti	Hu	To	Ti	K	Ma	Ti	To
50	60	74	63	79	73	57	62	68	58	53	50	57	56	58	62	63	58	52
72	86	88	78	70	72	52	59	52	52	51	52	52	52	55	55	57	57	51
64	73	71	63	59	57	59	55	52	52	51	53	53	53	55	54	59	51	51
70	61	65	68	53	55	48	45	46	45	45	53	53	53	54	49	51	54	51
			61	76	66	66	63	59	59	62	67	62	71	68	71	67	68	61
			76	65	62	62	59	59	61	62	63	65	66	67	66	65	65	64
			64	63	62	55	53	53	53	55	52	57	62	57	59	55	55	54
				53	62	52	55	58	55	55	56	58	60	60	60	60	57	51
				57	52	52	54	55	55	51	55	51	56	55	55	51	51	51
				54	51	52	52	52	52	52	53	56	55	55	54	53	52	52
				53	59	52	52	52	52	52	53	59	52	52	52	53	52	52
				57	59	59	55	55	55	55	56	57	57	57	57	57	57	57
				53	59	59	55	55	55	55	56	57	57	57	57	57	57	57
				54	51	51	51	51	51	51	51	51	51	51	51	51	51	51
				56	61	51	54	52	52	52	53	53	53	53	53	53	53	51
				51	73	73	69	67	70	69	68	67	70	69	68	69	68	68
				63	61	73	73	69	67	70	69	68	67	70	69	68	68	68
				19	19	21	21	21	21	21	21	21	21	21	21	21	21	21
				71	71	72	72	72	72	72	72	72	72	72	72	72	72	72
				75	75	75	75	75	75	75	75	75	75	75	75	75	75	75

<sup>1</sup> Percentage based on incomplete data.

TABLE 4

A tentative family tree for Polynesia<sup>1</sup>



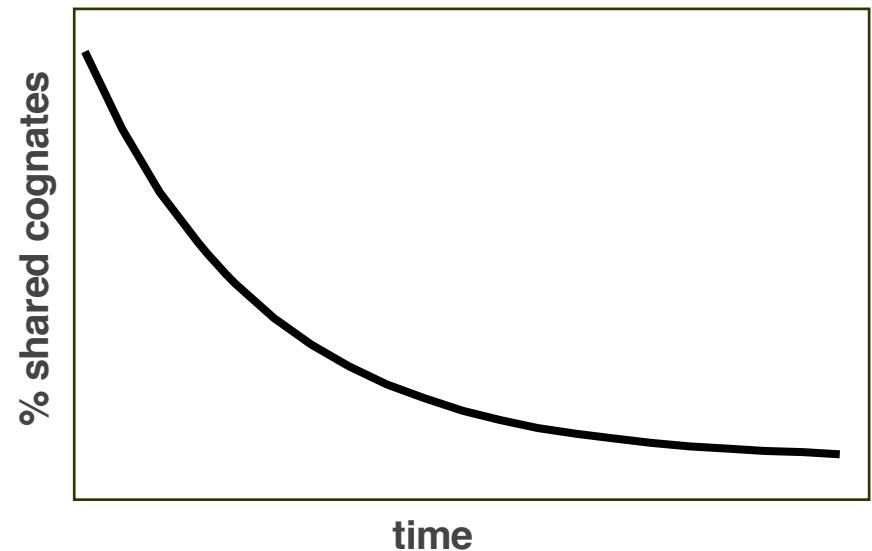


# Glottochronology

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

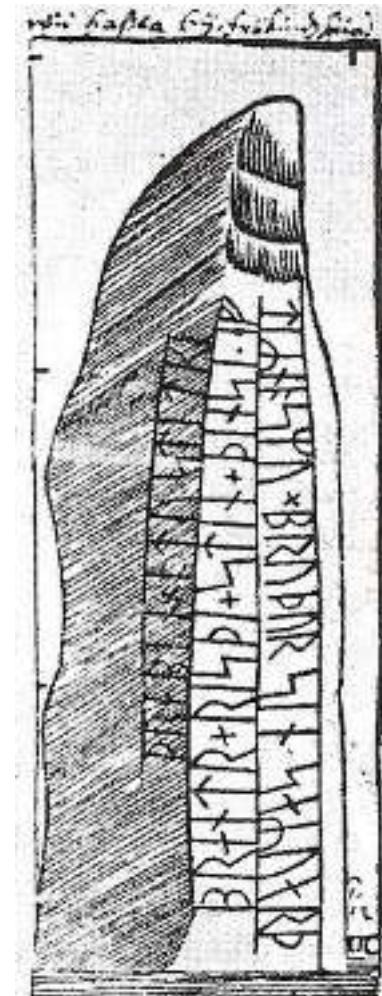
# ...and the fall of Lexicostatistics

## 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

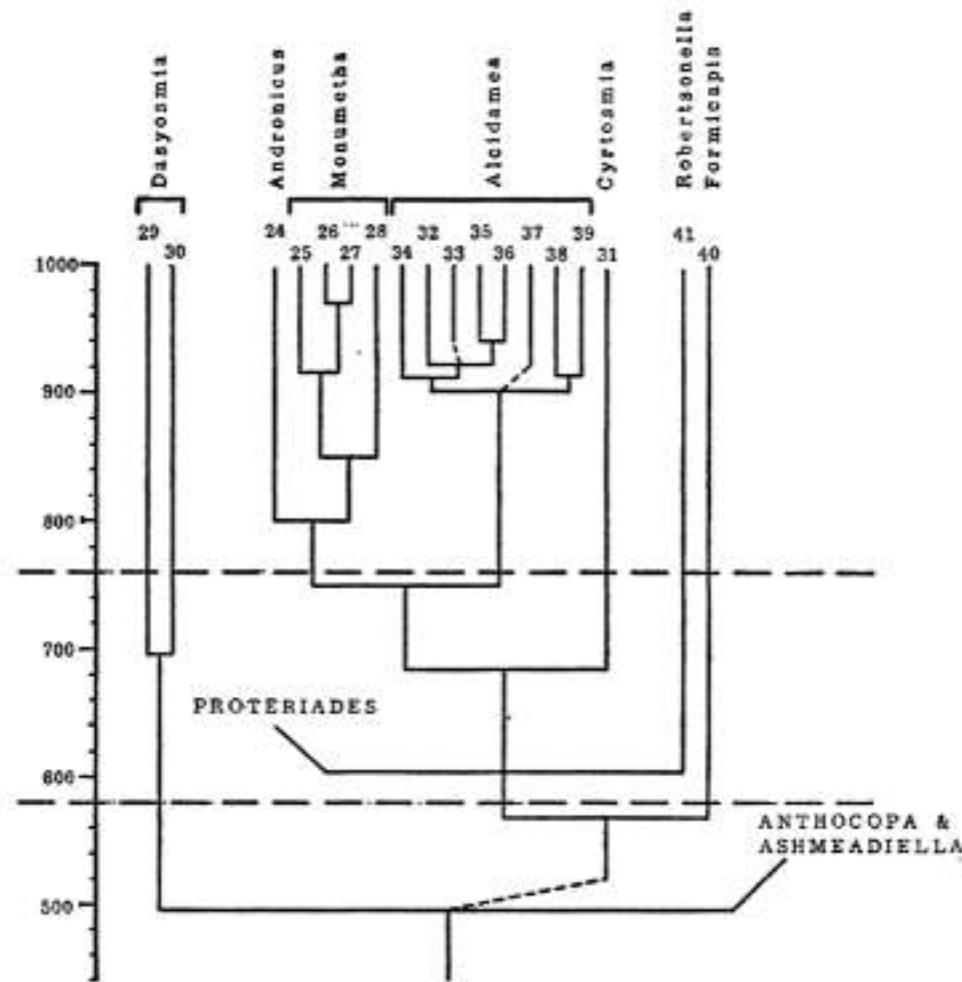


FIG. 6. Diagram of relationships for the genus *Hoplitis* obtained by the weighted variable group method.

A QUANTITATIVE APPROACH TO A PROBLEM  
IN CLASSIFICATION<sup>1</sup>

CHARLES D. MICHENER AND ROBERT R. SOKAL<sup>2</sup>

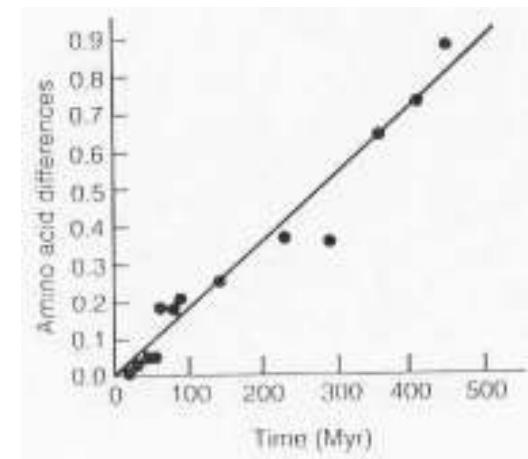
Department of Entomology, University of Kansas, Lawrence



# Molecular Clock

Zuckerkandl & Pauling 1962

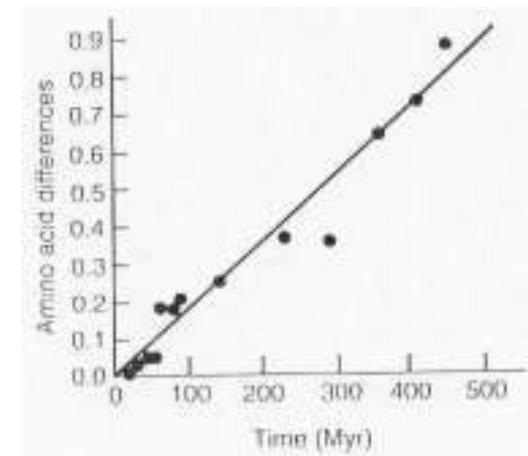
Number of AA differences were proportional to species divergence times.



# Molecular Clock

Zuckerkandl & Pauling 1962

Number of AA differences were proportional to species divergence times.



Kimura 1968

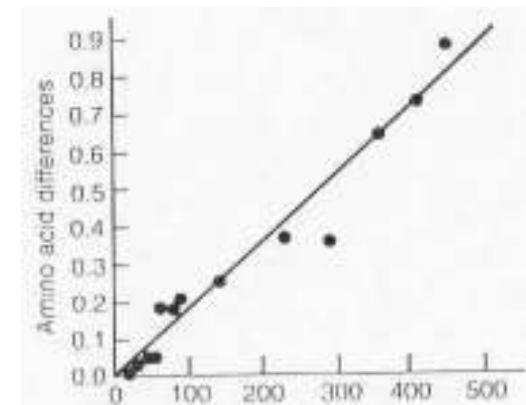
The average time taken for one base pair replacement within a genome is therefore

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

# 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

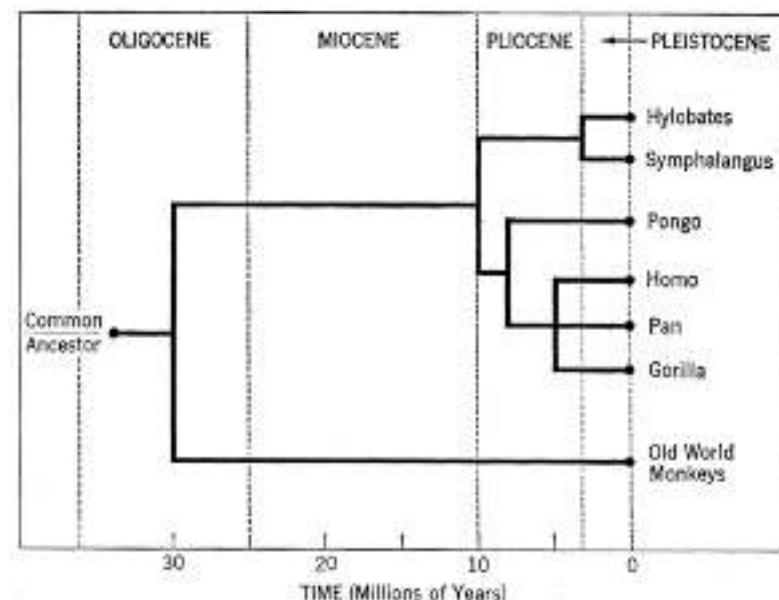


Fig. 1. Times of divergence between the various hominoids, as estimated from immunological data. The time of divergence of hominoids and Old World monkeys is assumed to be 30 million years.

# 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 MISLEADING<sup>1</sup>

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 Sites<sup>1</sup>**

*Ziheng Yang*

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

*Michael J. Sanderson*

## **Relaxed Phylogenetics and Dating with Confidence**

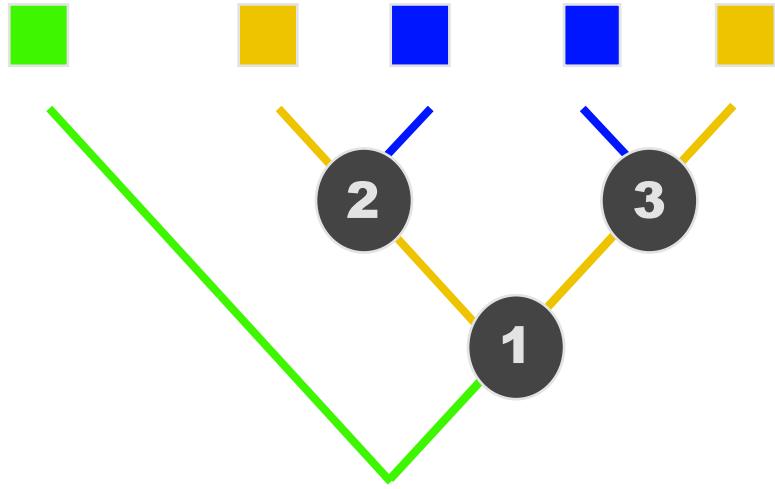
*Alexei J. Drummond<sup>✉</sup>, Simon Y. W. Ho, Matthew J. Phillips, Andrew Rambaut<sup>✉\*</sup>*

# How do we build trees?

1. Distance Methods (Chiara)
2. Maximum Parsimony.
3. Maximum Likelihood.
4. **Bayesian Phylogenetic Analyses.**



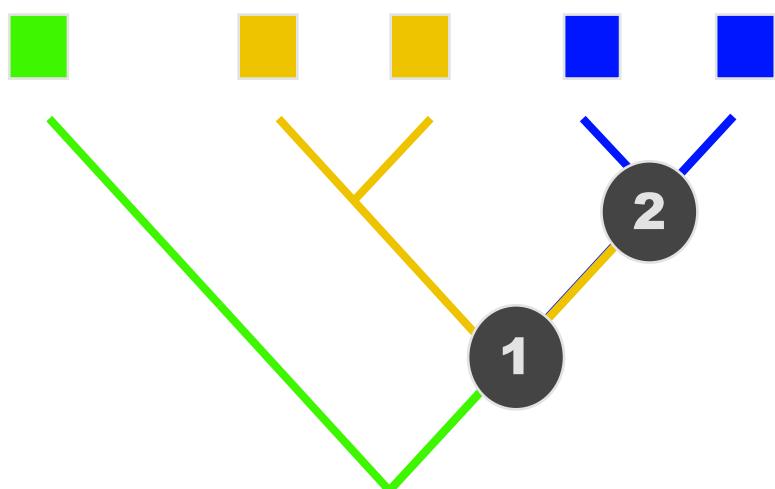
# Maximum Parsimony



Unlikely that complex traits  
should arise more than once

=> Best tree is maximally  
parsimonious

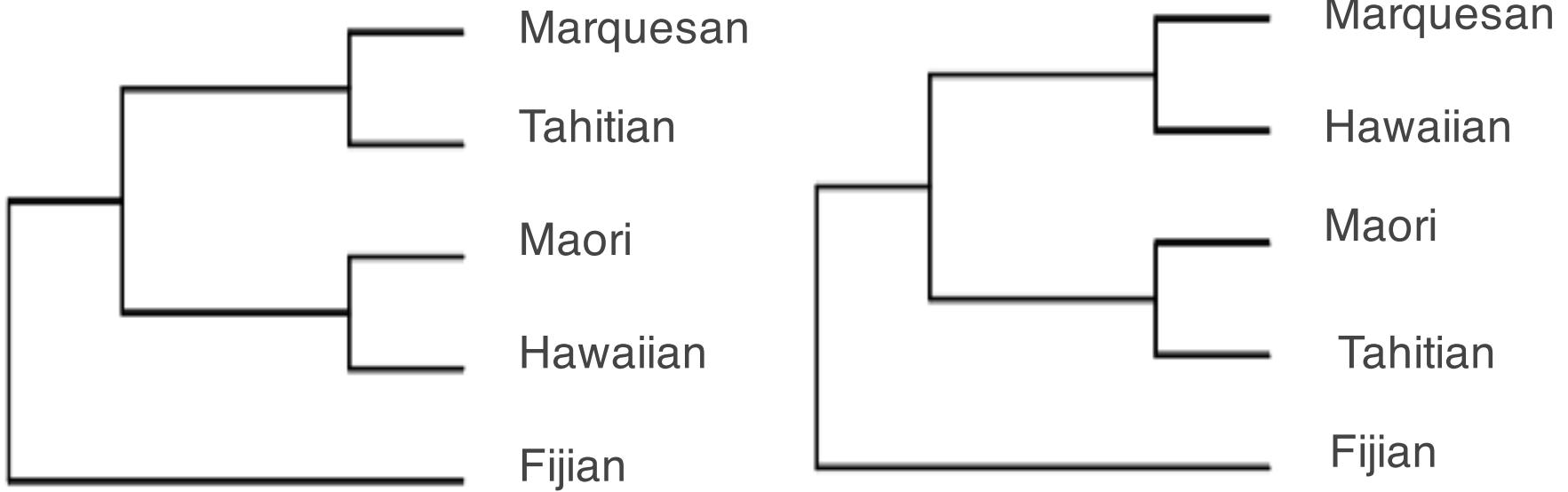
=> Smallest amount of  
evolution (Fewest number of  
character state changes)



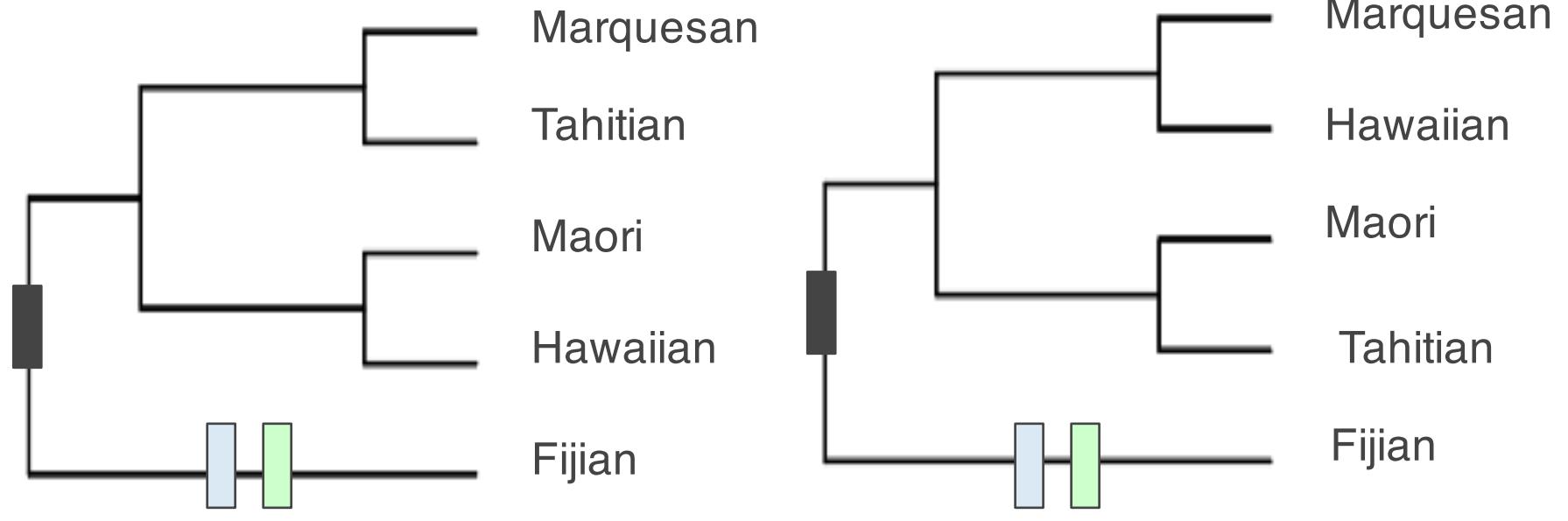
	Taboo	Blood	To Suck
Fijian	tabu	drā	sucu-ma
Tahitian	tapu	toto	ngote
Maori	tapu	toto	ngote
Hawaiian	kapu	koko	omo
Marquesan	tapu	toto	omo

	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



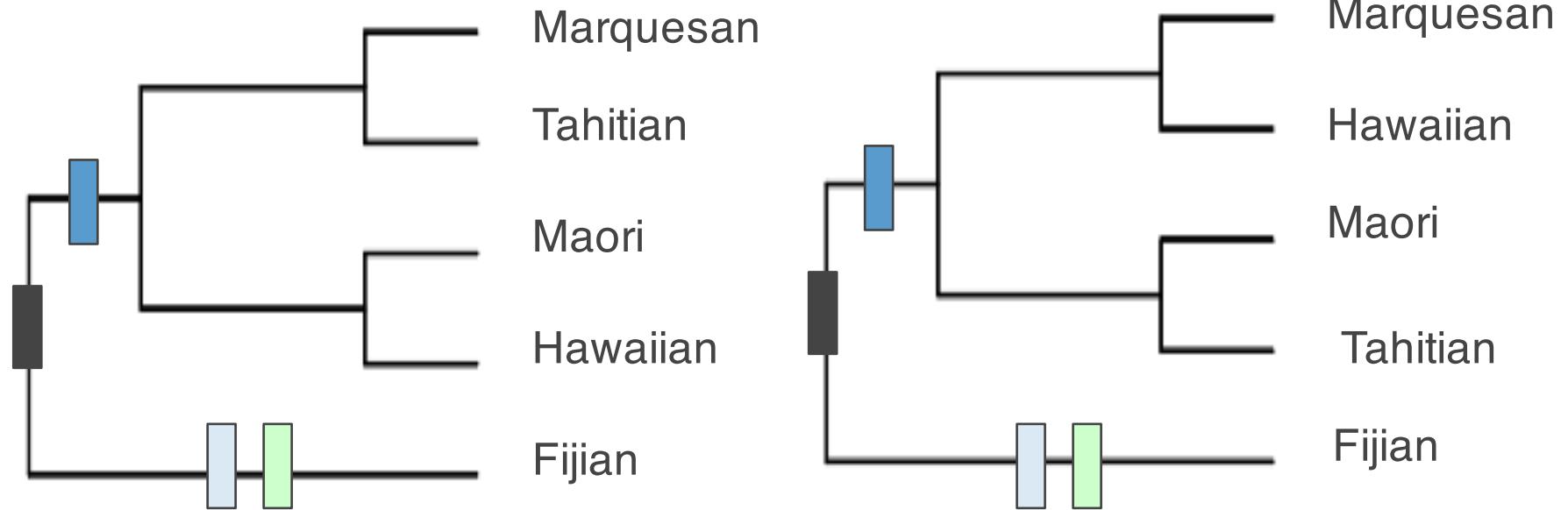
<b>Fijian</b>	1	1	0	1	0	0
<b>Tahitian</b>	1	0	1	0	1	0
<b>Maori</b>	1	0	1	0	1	0
<b>Hawaiian</b>	1	0	1	0	0	1
<b>Marquesan</b>	1	0	1	0	0	1



Length=3

Length=3

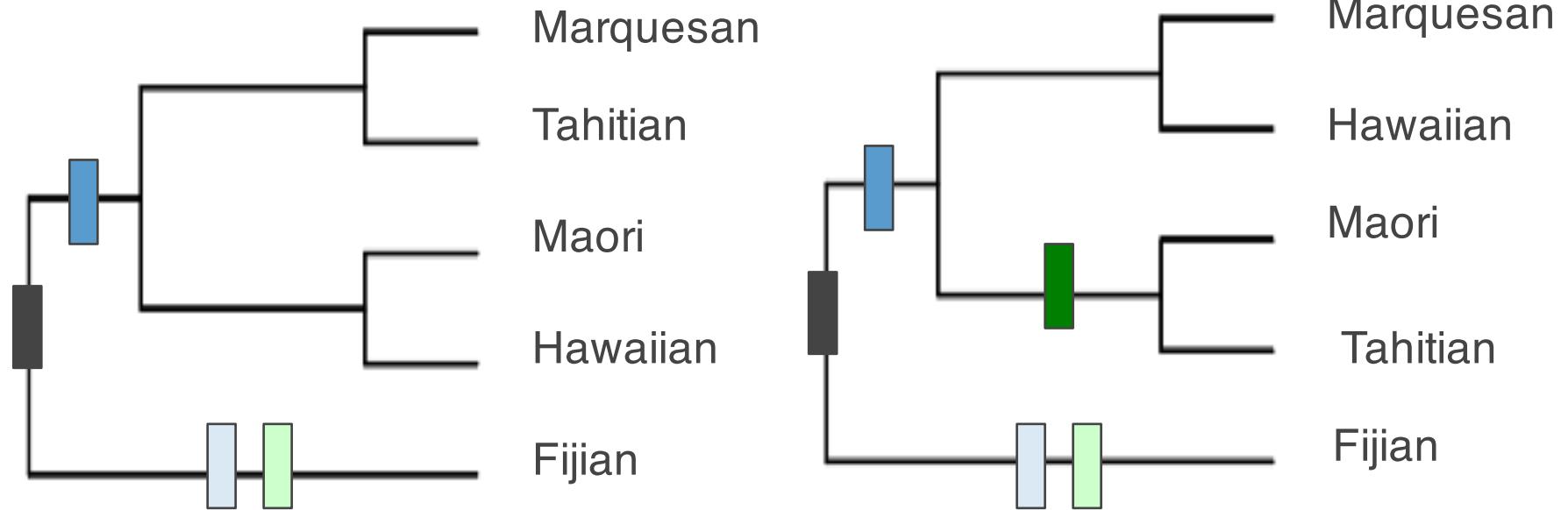
<b>Fijian</b>	1	1	0	1	0	0
<b>Tahitian</b>	1	0	1	0	1	0
<b>Maori</b>	1	0	1	0	1	0
<b>Hawaiian</b>	1	0	1	0	0	1
<b>Marquesan</b>	1	0	1	0	0	1



Length=4

Length=4

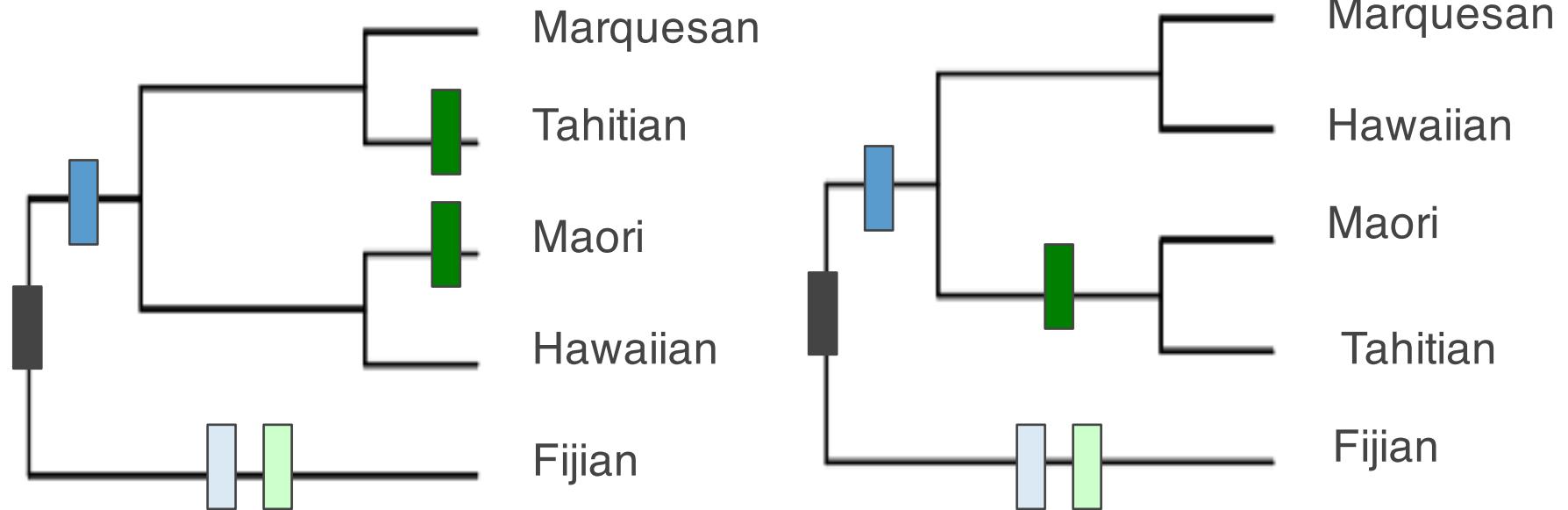
<b>Fijian</b>	1	1	0	1	0	0
<b>Tahitian</b>	1	0	1	0	1	0
<b>Maori</b>	1	0	1	0	1	0
<b>Hawaiian</b>	1	0	1	0	0	1
<b>Marquesan</b>	1	0	1	0	0	1



Length=4

Length=5

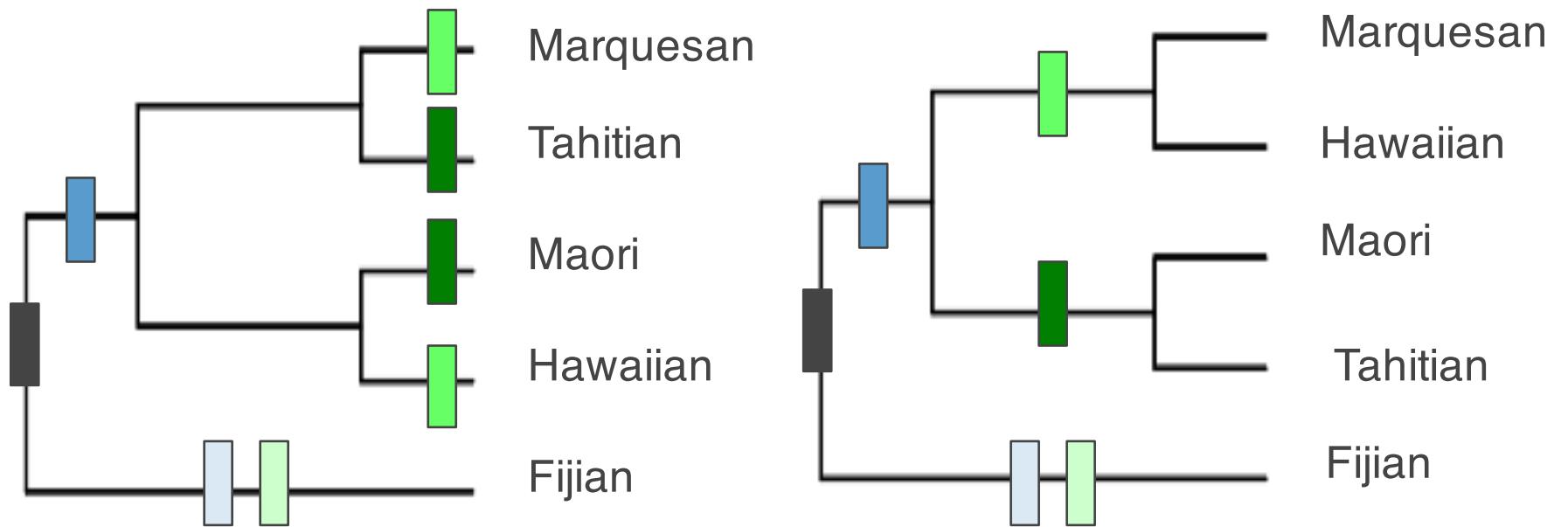
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



Length=6

Length=5

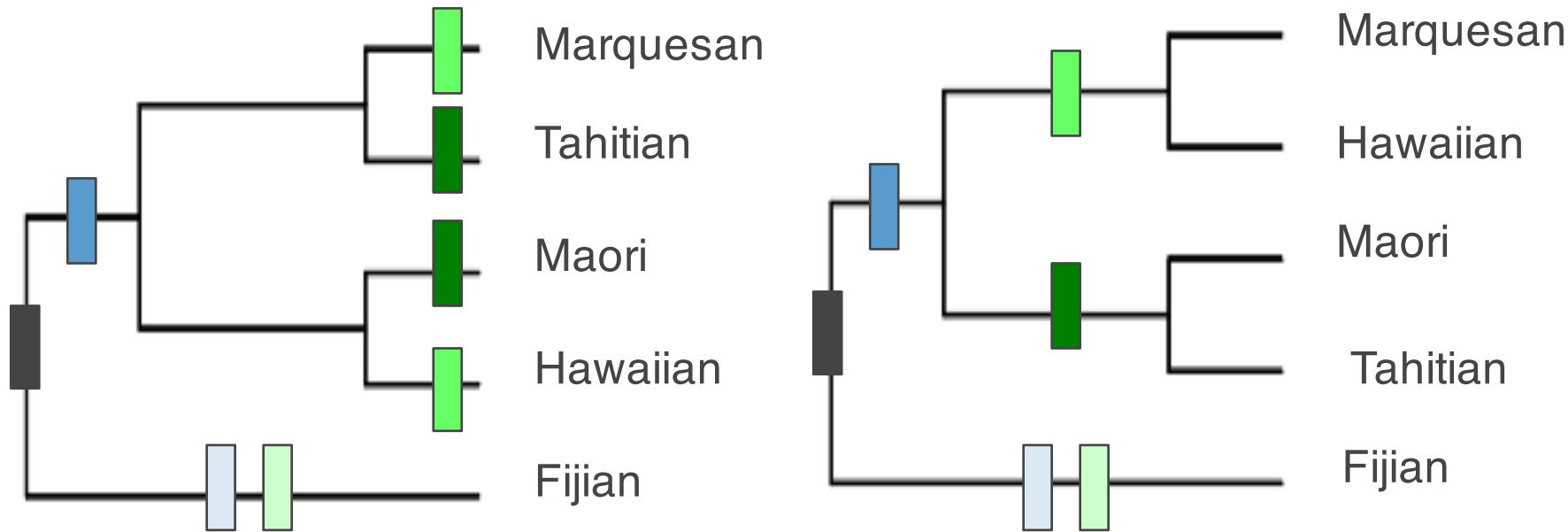
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



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

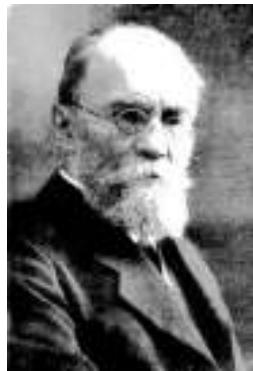


Length=8

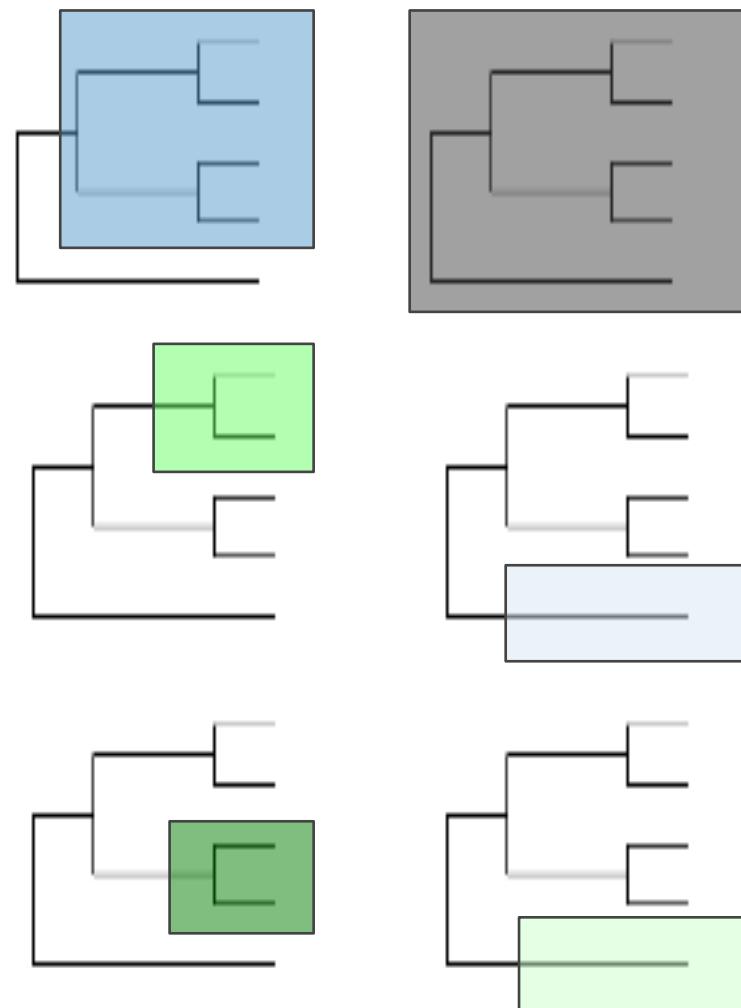
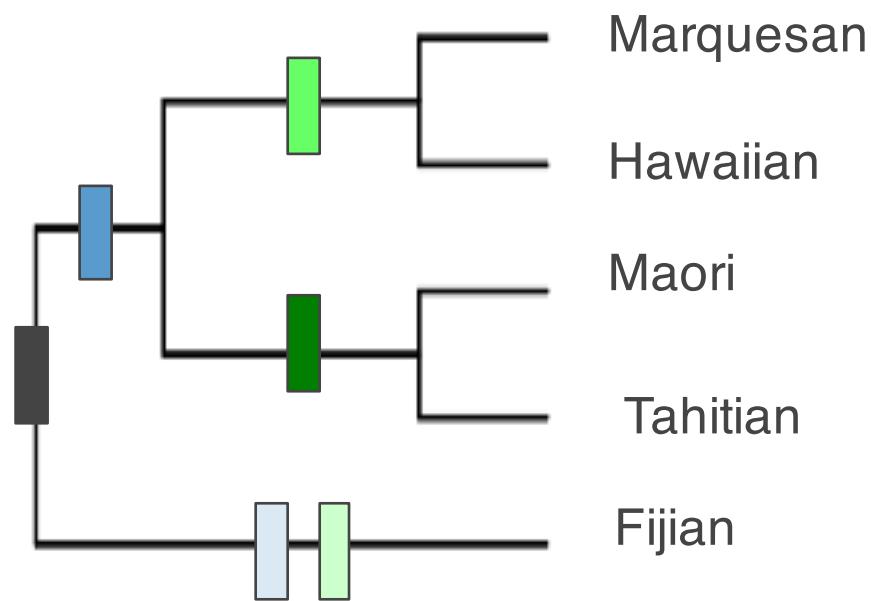
Length=6



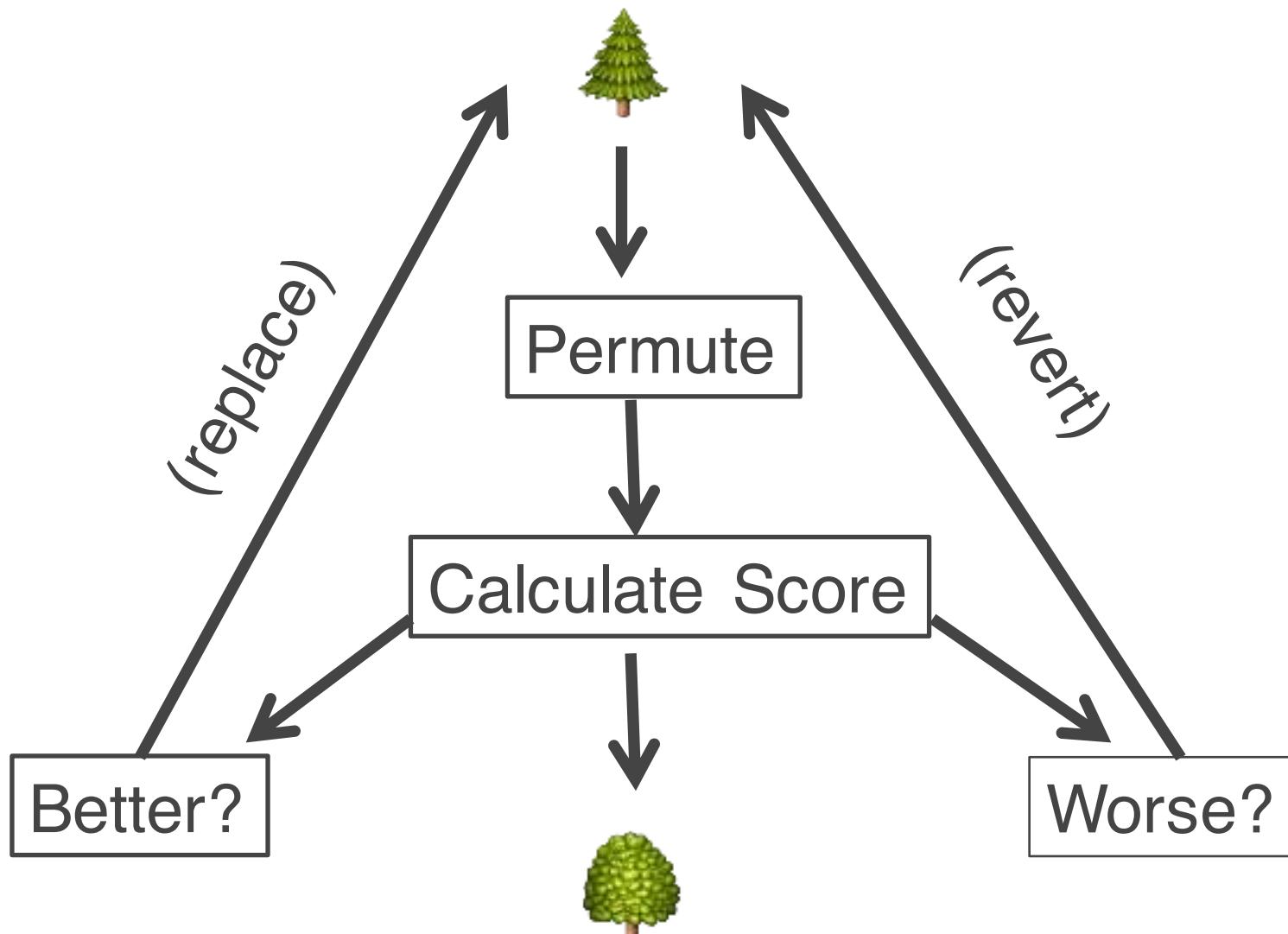
Fijian	1	1	0	1	0	0
Tahitian	1	0	1	0	1	0
Maori	1	0	1	0	1	0
Hawaiian	1	0	1	0	0	1
Marquesan	1	0	1	0	0	1



# Innovations vs. Retentions



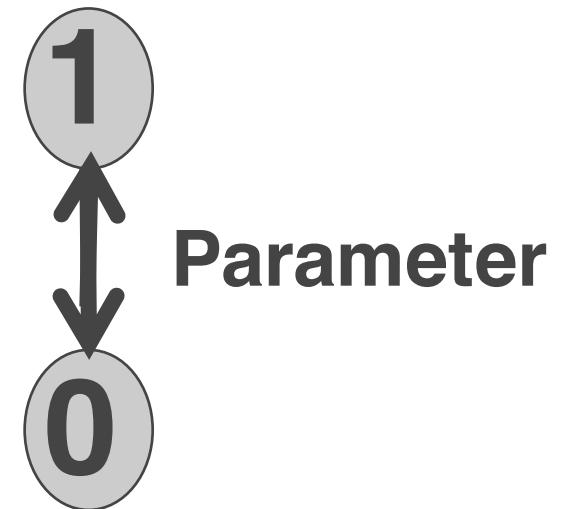
# Algorithm

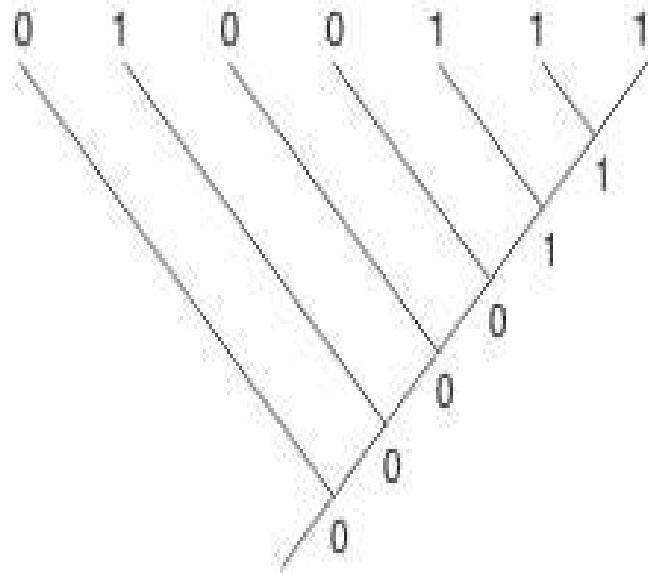


“best” tree

# Maximum Likelihood

- Builds on Max. Parsimony
- Stochastic **model** of change  
(=assumptions about how evolution works)
- => Find tree that maximises the likelihood
- **Likelihood** = fit of data to tree under a model.
  - Very small number =  $\log(L_h)$
  - Closer to zero = better fit.





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

$L_h =$

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.

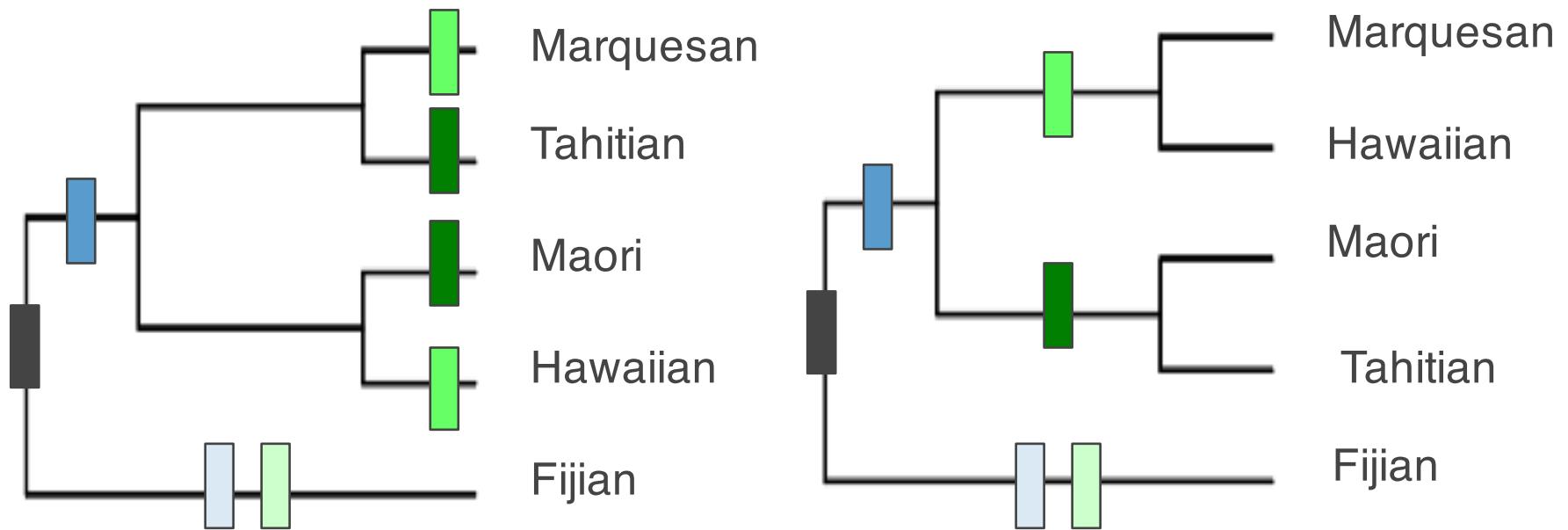
x

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

.... etc

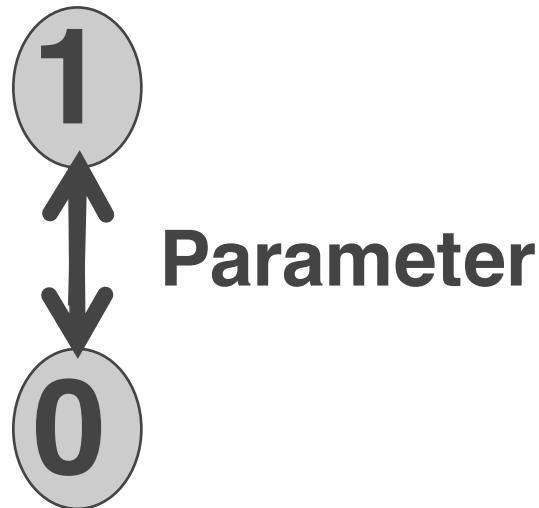
$$\text{Site Unlinked}(s) = \left( \begin{array}{c} \diagup \\ \diagdown \end{array} \right) \dots \left( \begin{array}{c} \diagup \\ \diagdown \end{array} \right) \dots \left( \begin{array}{c} \diagup \\ \diagdown \end{array} \right)$$

Site Unlinked( $s'$ ) = Preconstruction 1 ... &lt;> ... Preconstruction 6 &lt;> ... Preconstruction n

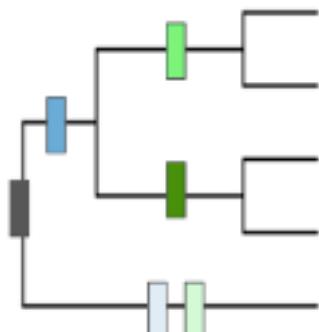

 $\text{Ln}(L) = -14.804$ 
 $\text{Ln}(L) = -12.007$  ←

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

# Model Comparison with Lh

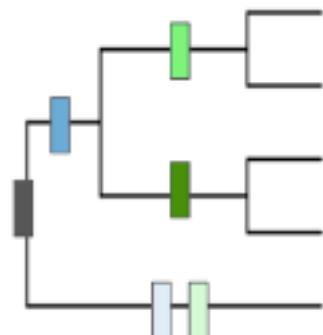
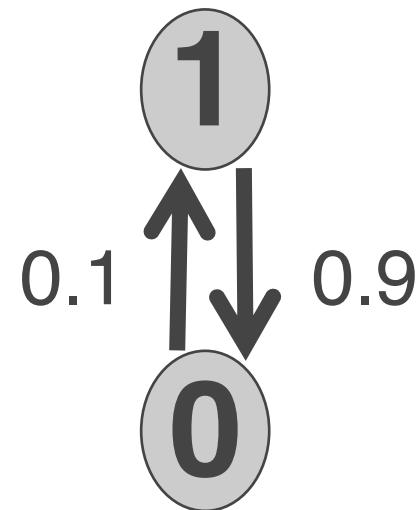
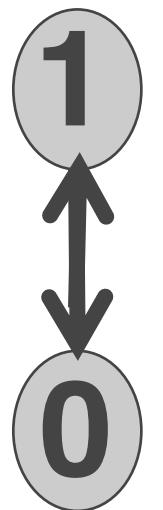


Can we modify our  
**assumptions** and get a **better**  
explanation of the data?



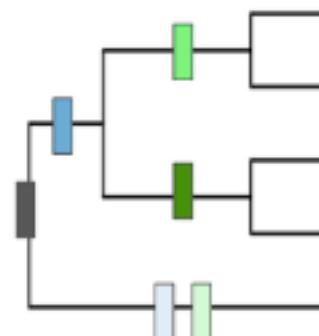
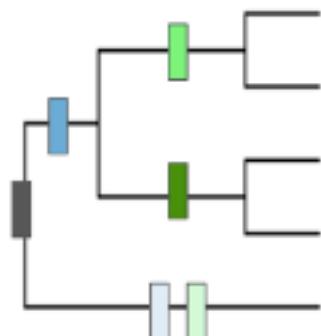
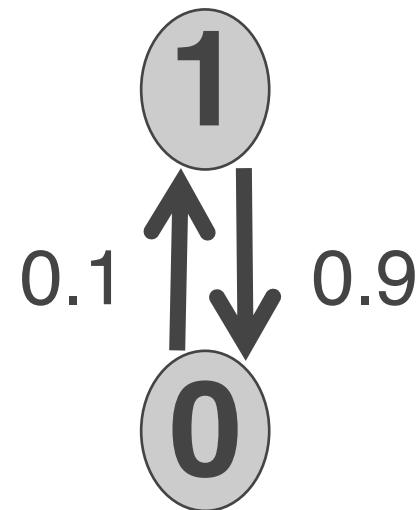
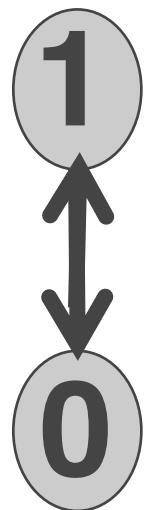
$$\ln(L) = -12.007$$

# Models



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

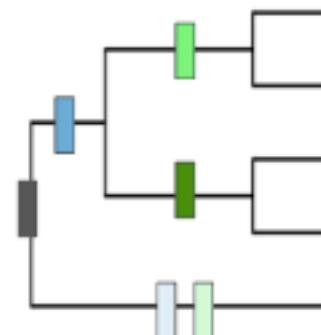
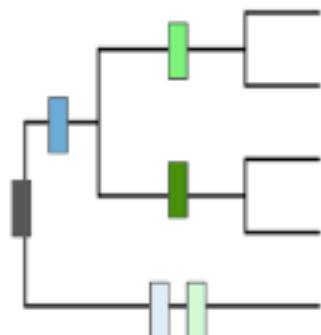
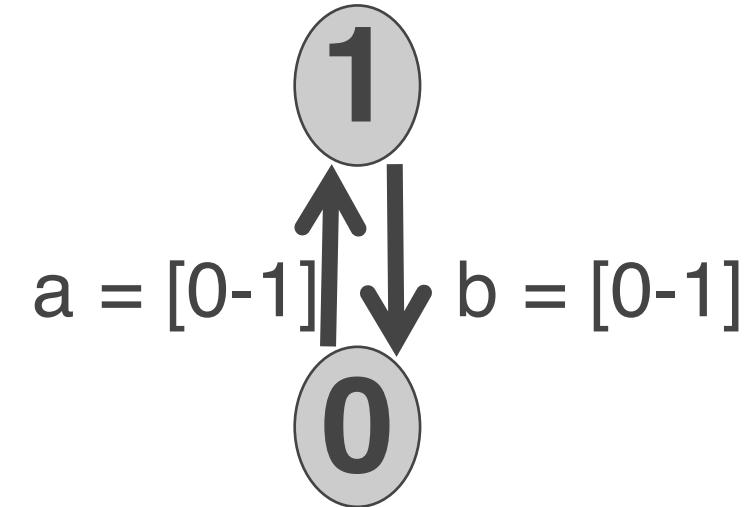
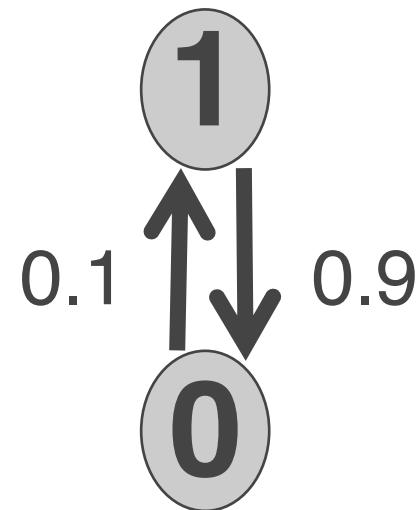
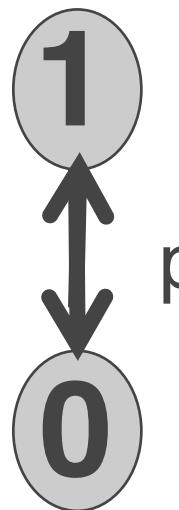
# Models



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

$$\ln(L) = -11.310$$

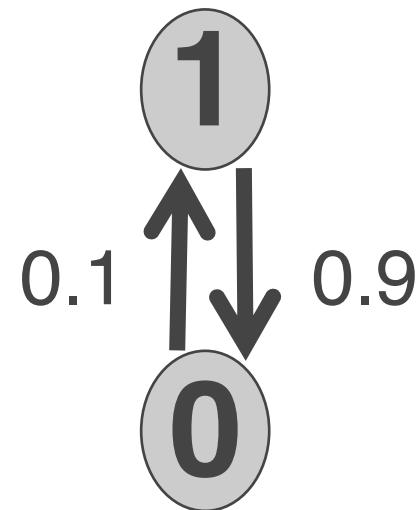
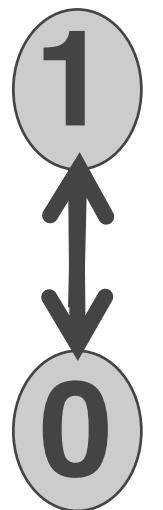
# Models



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

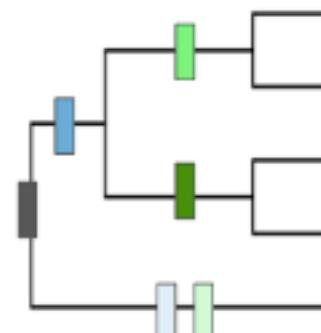
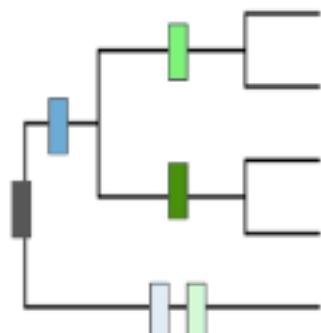
$$\ln(L) = -11.310$$

# Models



$a = [0-1]$     $b = [0-1]$

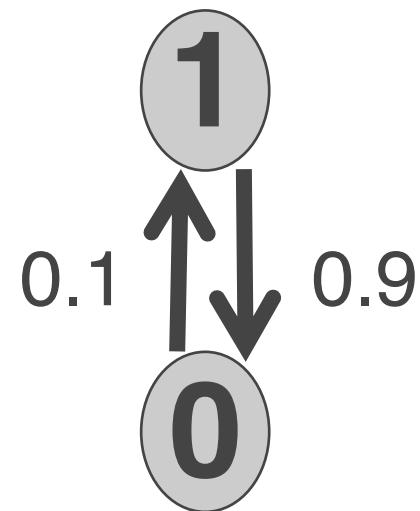
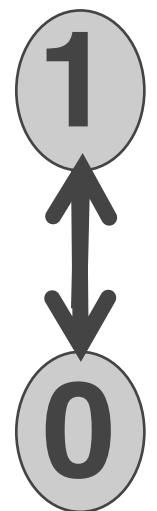
$a = 0.92, b = 0.08$



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

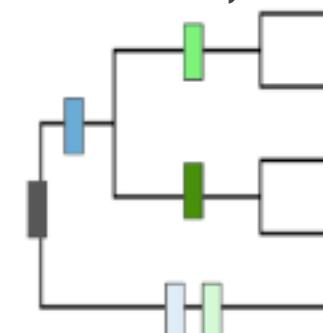
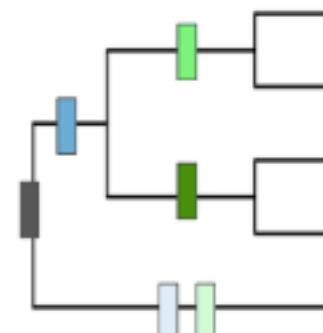
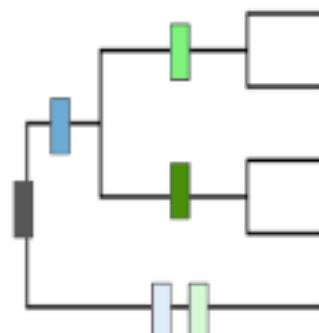
$$\ln(L) = -11.310$$

# Models



$a = [0-1]$   $b = [0-1]$

$a = 0.92, b = 0.08$

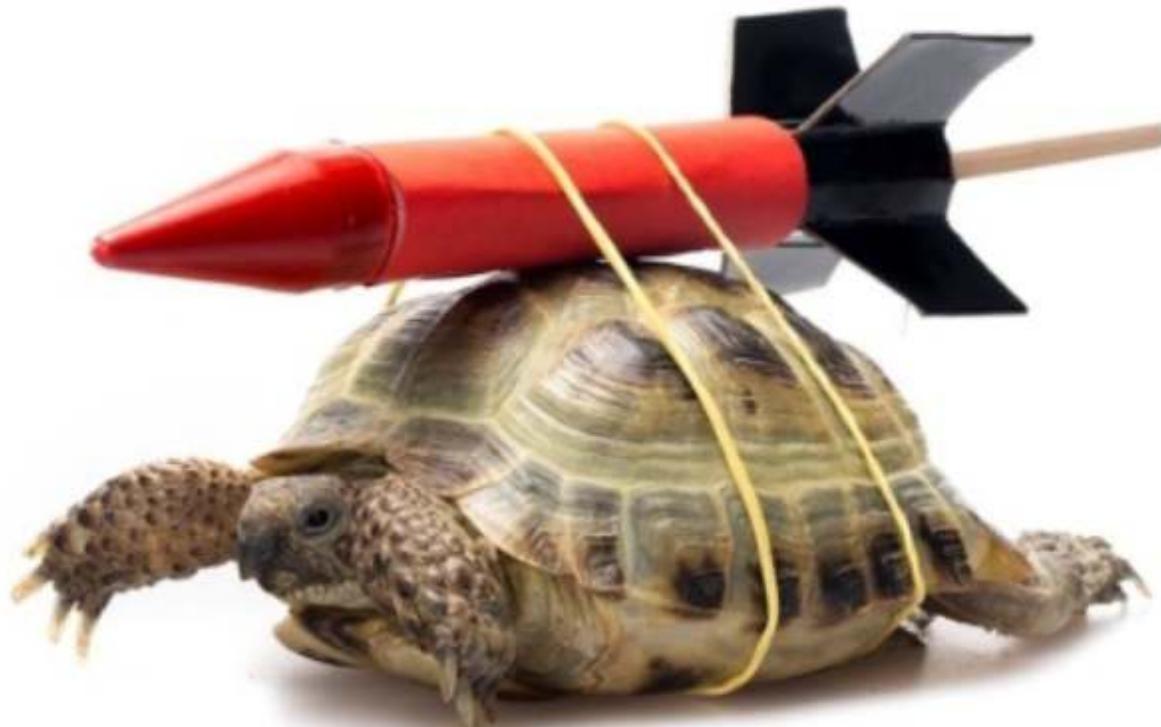


$$\ln(L) = -12.007$$

$$\ln(L) = -11.310$$

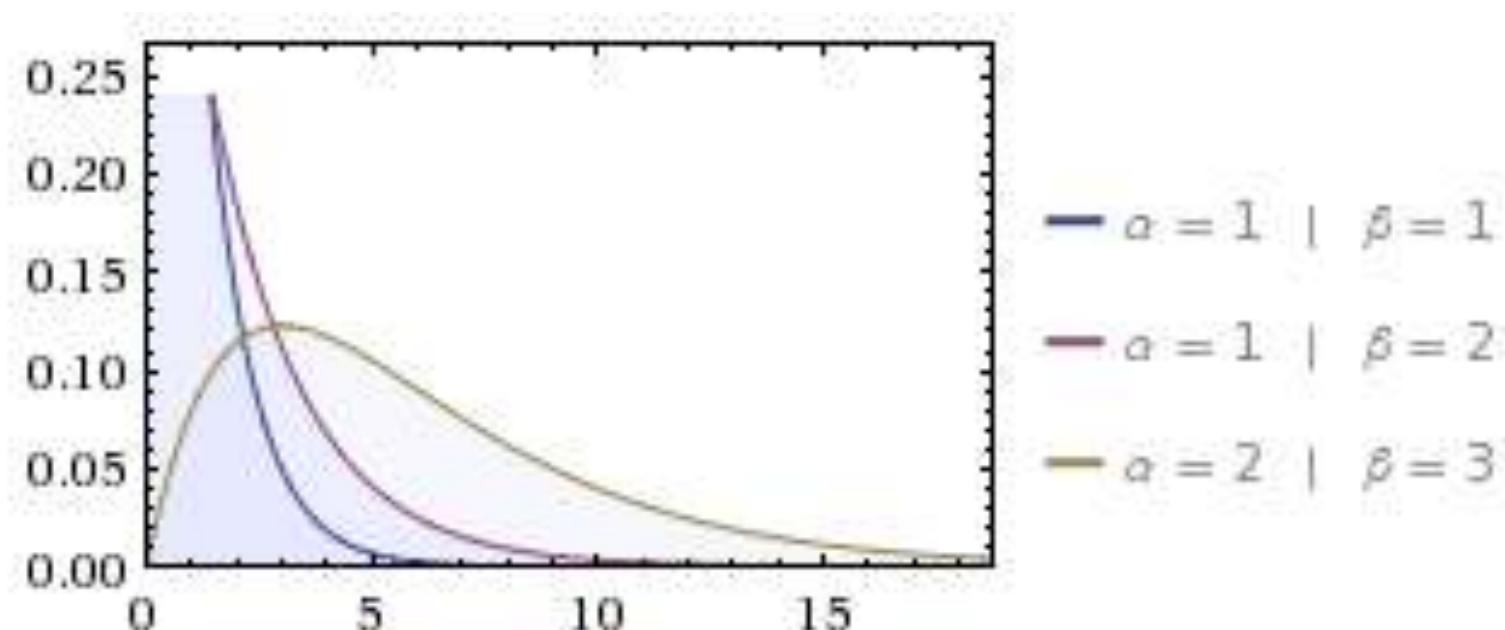
$$\ln(L) = -9.072$$

# But Rates Vary!



# Rate Variation - Gamma

- “Site-Specific Rates”: Rates vary across **sites**
- Gamma Distribution: One parameter,  $\alpha$ , controls the shape
- Discretized into  $n$  categories and Characters assigned a category
- Estimate the best value for  $\alpha$  using Lh.



# Rate Variation - Covarion

- Rates change across **phylogeny**
- Each site is either “on” or “off”
  - In “on” state it evolves according to one set of rates
  - In “off” state it evolves according to another (often 0)
- Analysis estimate when sites are “on” and “off” and switching rates.
- = Each site can evolve at different rates at different times on the tree.

Among-Site Rate Variation



Site-Specific Rate Variation



Gautier 2001

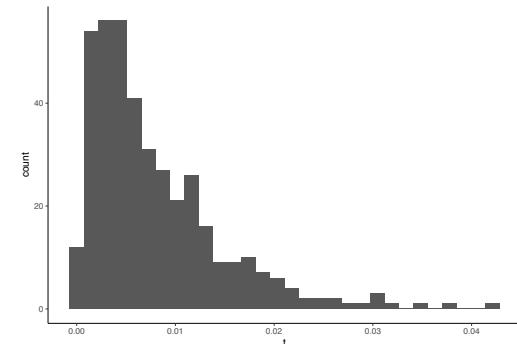
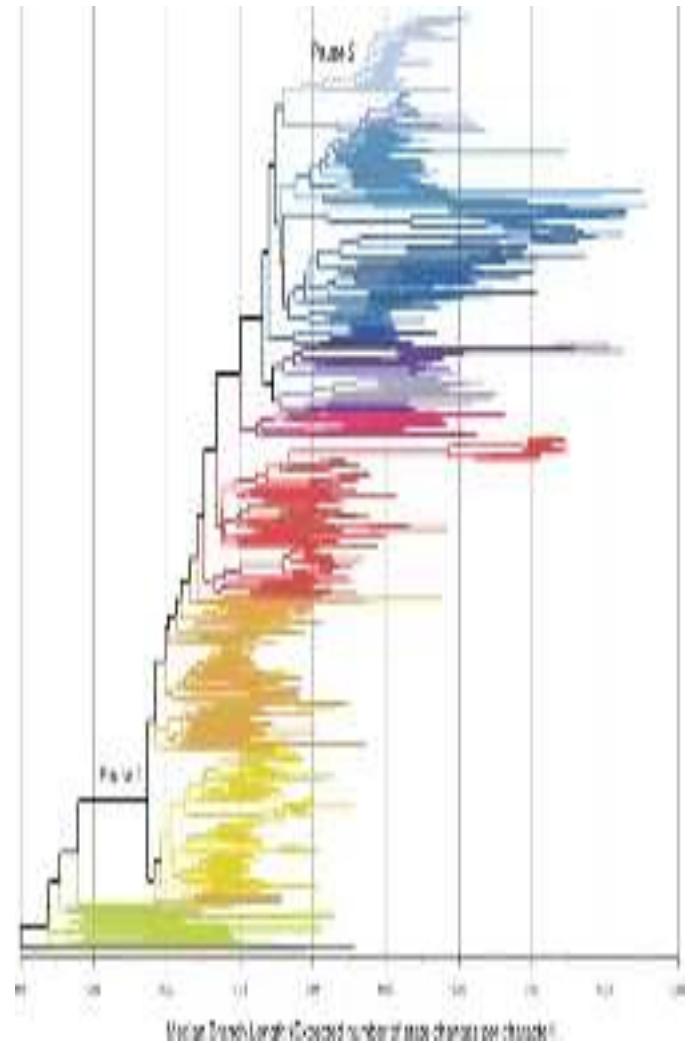


“Linguists don’t do dates.”

# Phylogenetic Dating

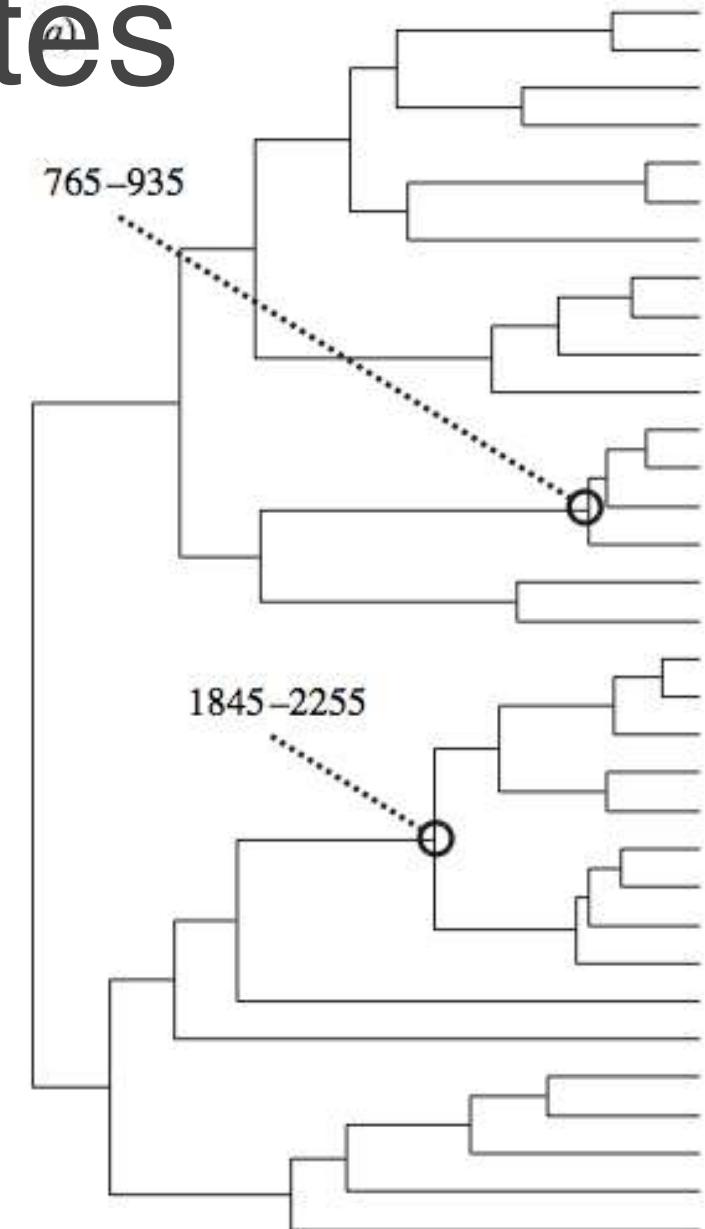
ML estimates amount of change along lineage  
= (number of changes per character)

Not a global retention rate but a **per-lineage** estimate of the amount of change.



# Convert Rates to Dates

- (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
- $\sim$ = glottochronology

# Strict Clock

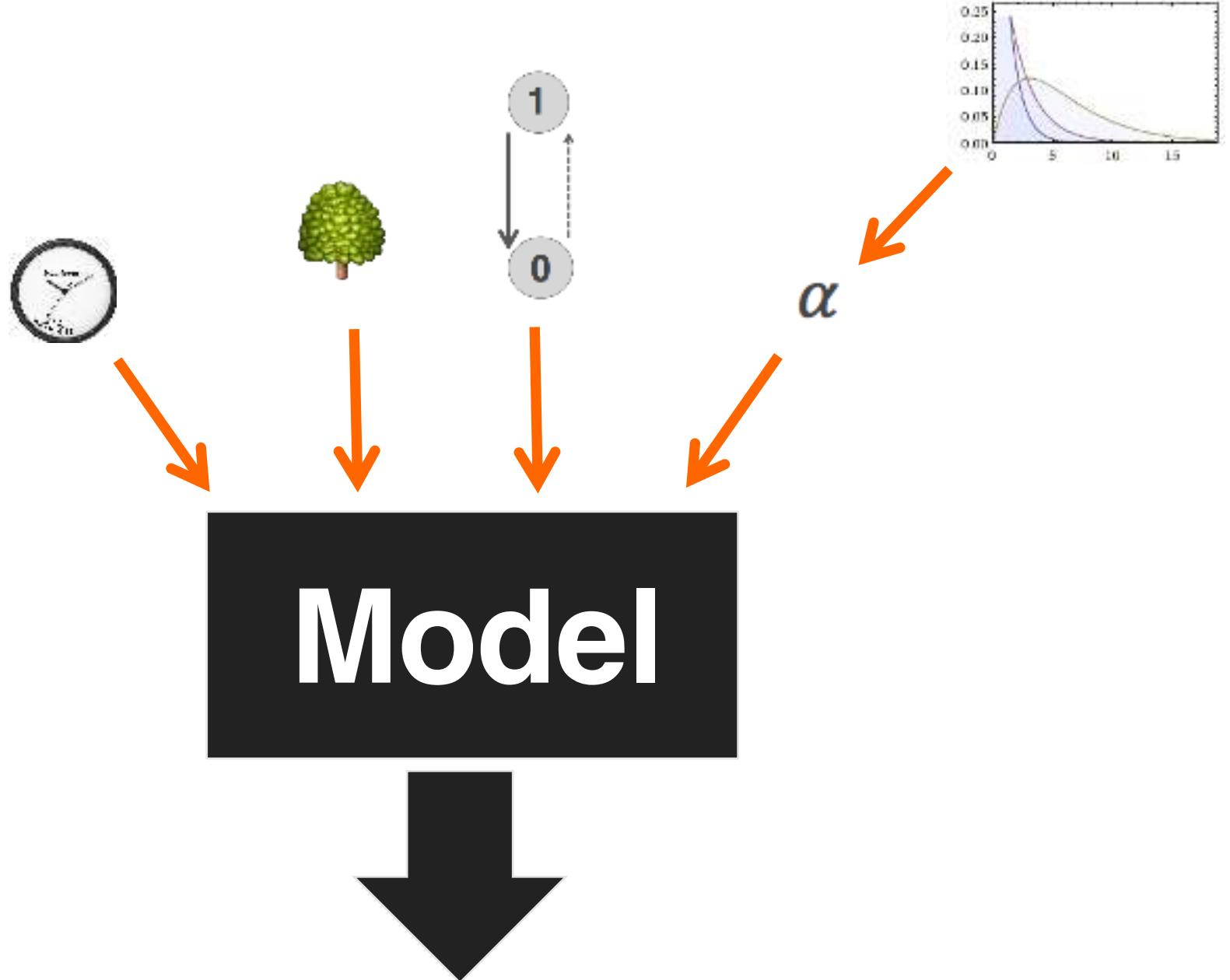


- One rate for all languages.
- No variation
- $\sim$ = glottochronology

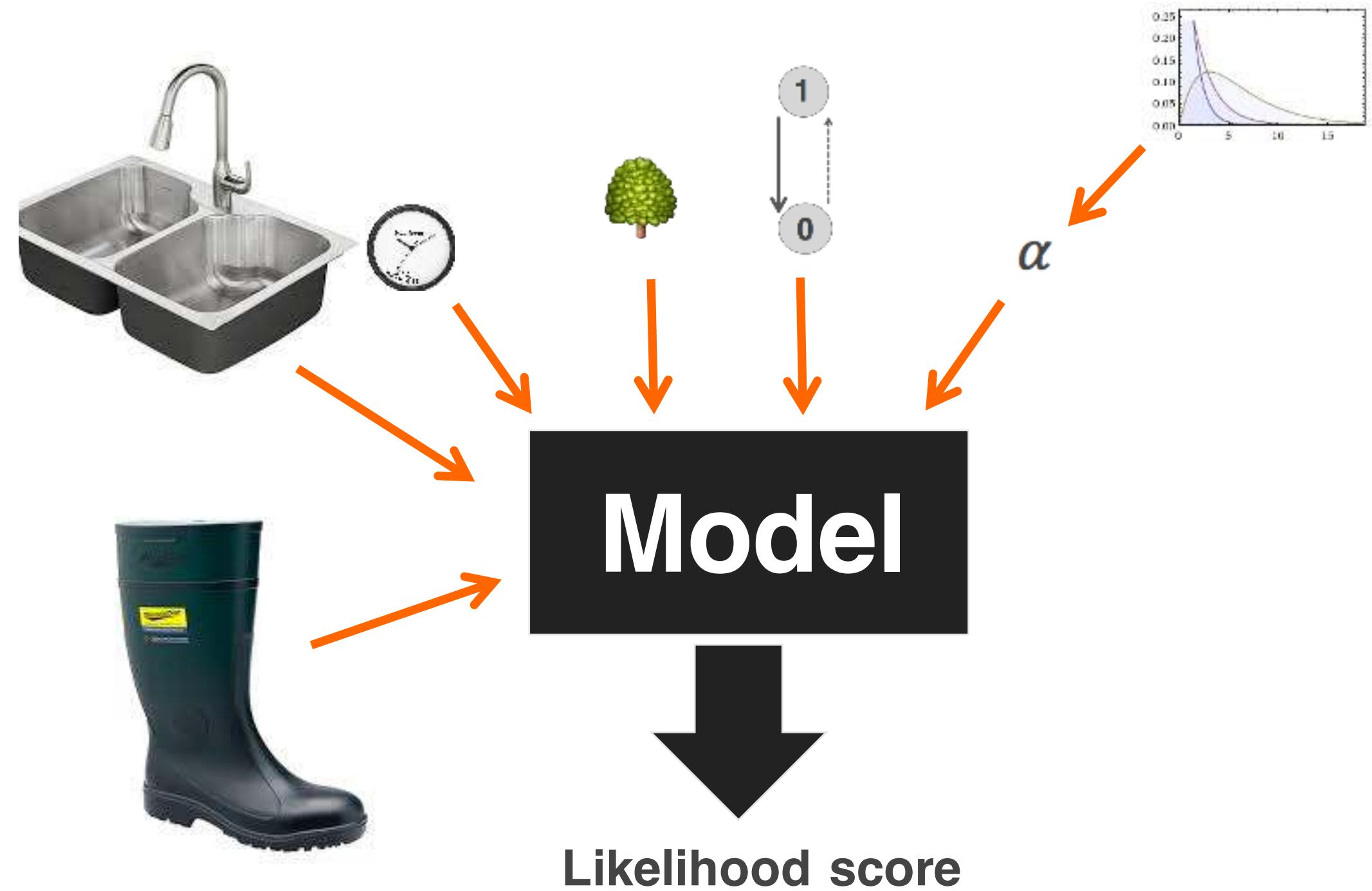
# Relaxed Clock

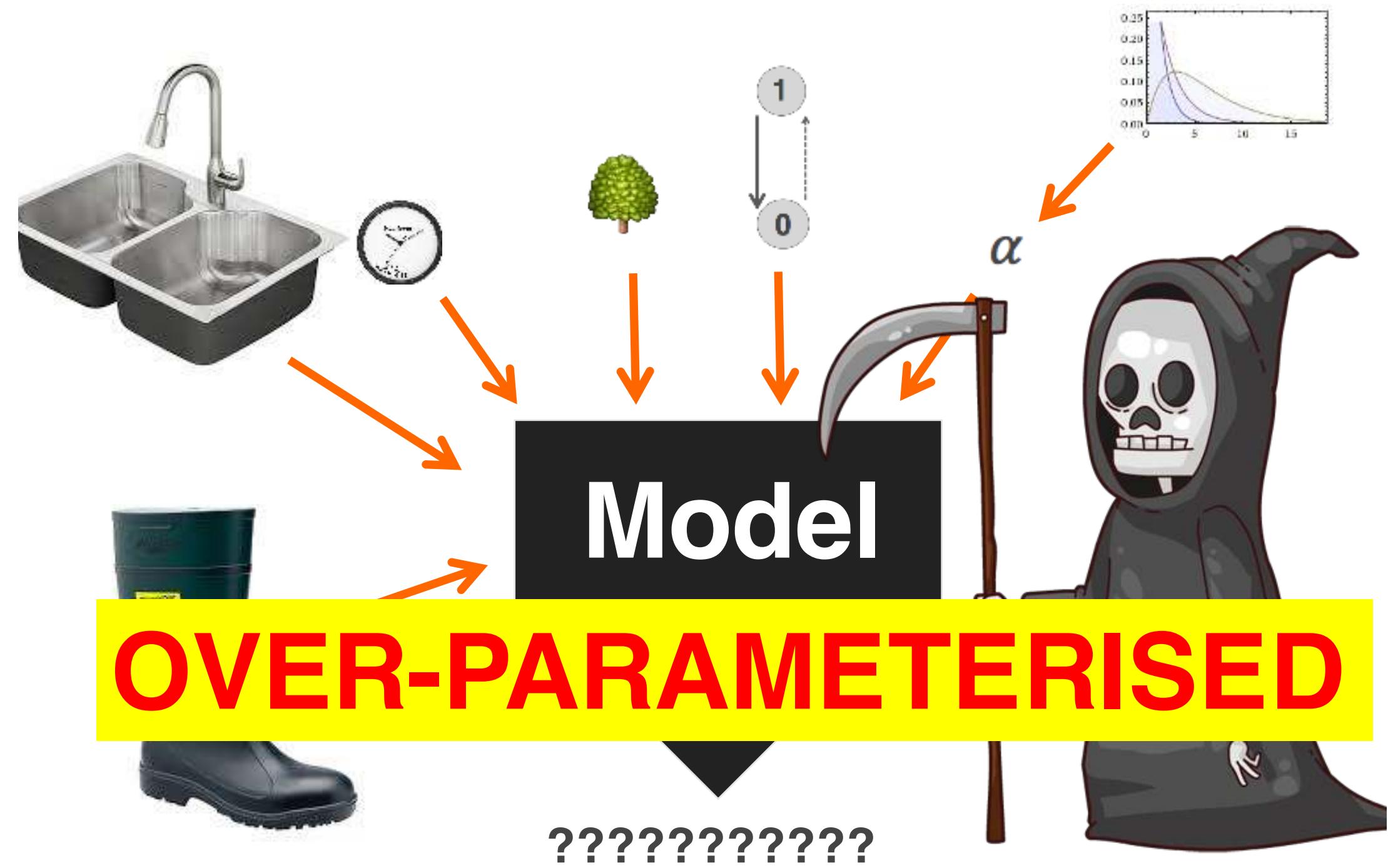


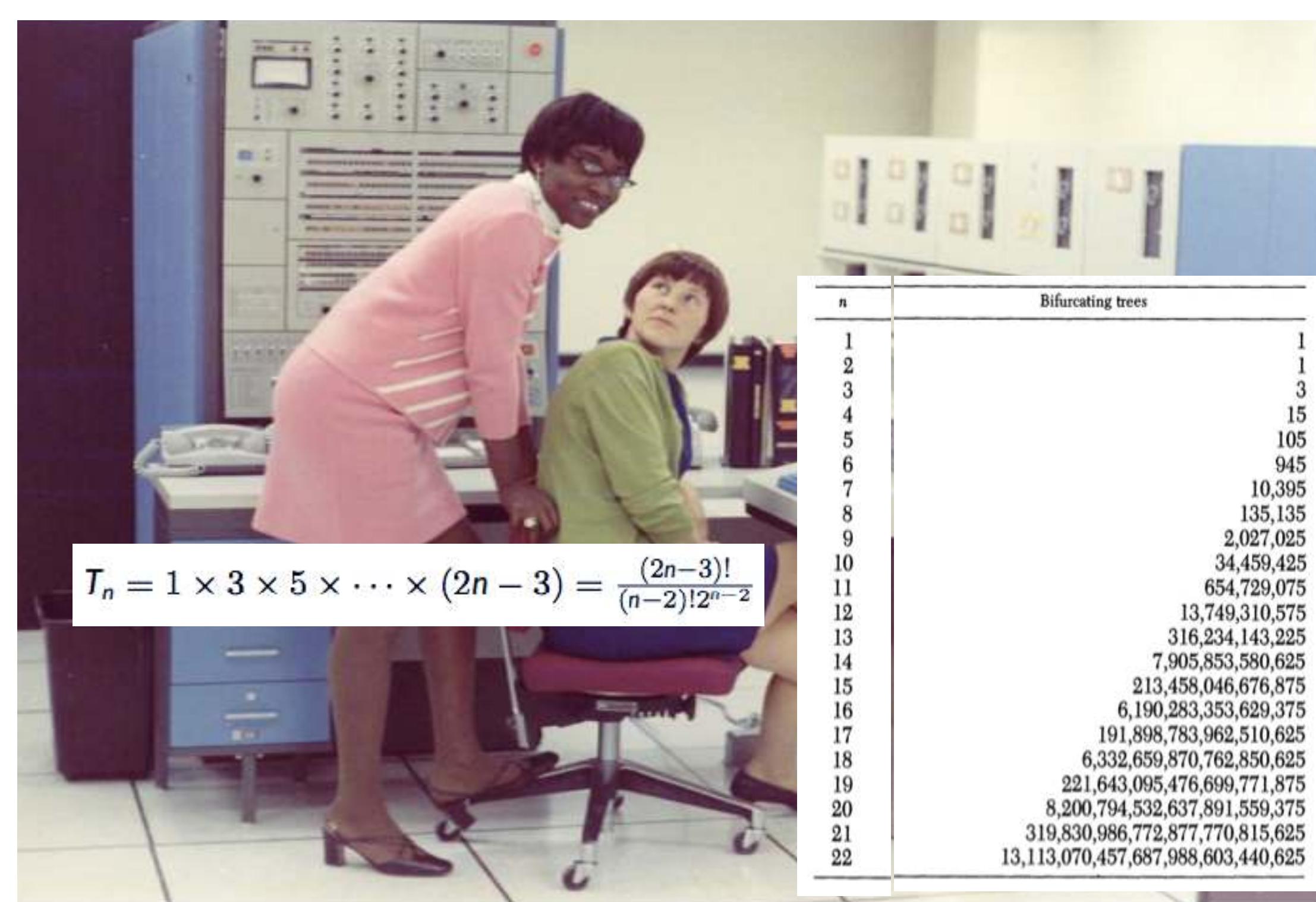
- Estimate distribution of *branch* rates from data
- Rates drawn from a parametric distribution estimated from the data
  - ( $\sim$ =gamma)
  - **LogNormal**, Exponential
- Allows rate to vary across branches (each branch could have own rate)



Likelihood score







$$T_n = 1 \times 3 \times 5 \times \cdots \times (2n - 3) = \frac{(2n-3)!}{(n-2)!2^{n-2}}$$

$n$	Bifurcating trees
1	1
2	1
3	3
4	15
5	105
6	945
7	10,395
8	135,135
9	2,027,025
10	34,459,425
11	654,729,075
12	13,749,310,575
13	316,234,143,225
14	7,905,853,580,625
15	213,458,046,676,875
16	6,190,283,353,629,375
17	191,898,783,962,510,625
18	6,332,659,870,762,850,625
19	221,643,095,476,699,771,875
20	8,200,794,532,637,891,559,375
21	319,830,986,772,877,770,815,625
22	13,113,070,457,687,988,603,440,625

# Need to handle complexity

- 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.





# Bayesian methods

- Extend ML methods to explicitly account for uncertainty

$$\Pr(H|D) = \frac{\Pr(H) \Pr(D|H)}{\Pr(D)}$$

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

Phylogenetics: Brown “The State of Bayesian Phylogenetics: Bayes for the Uninitiated” (<https://tinyurl.com/ly2cvfd>)



# Bayesian methods

- Extend ML methods to explicitly account for uncertainty

$$\Pr(H|D) = \frac{\text{prior} \quad \Pr(H) \Pr(D|H)}{\text{Lh} \quad \Pr(D)}$$

(Pr hypothesis given data)

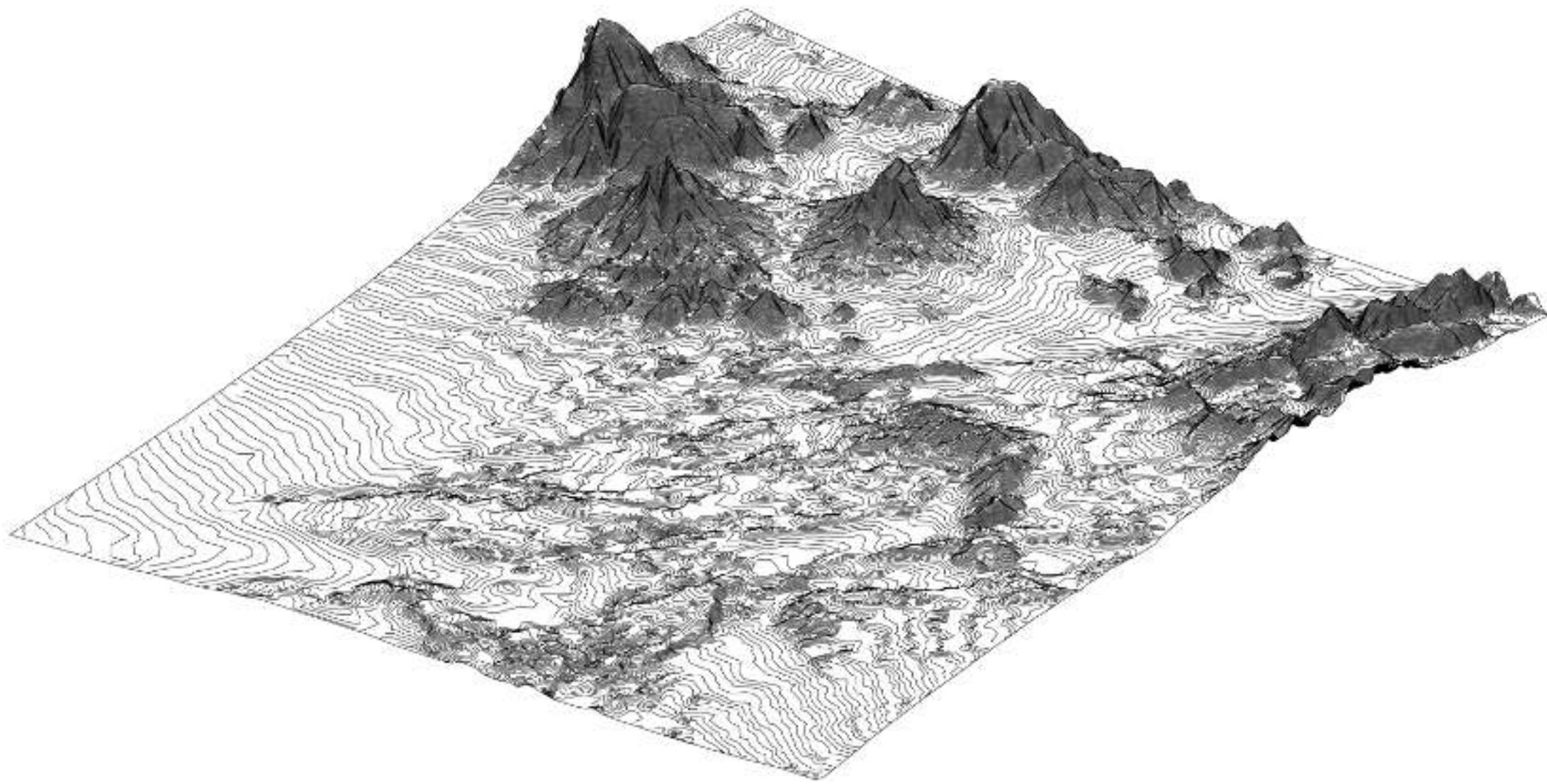
**Prior Pr that data is true**

The diagram illustrates the Bayes' theorem formula. A purple bracket on the left groups the term  $\Pr(H|D)$  and the text "(Pr hypothesis given data)". Above the fraction, the word "prior" is written in orange, with a yellow bracket above it grouping  $\Pr(H)$  and  $\Pr(D|H)$ . To the right of the fraction, the word "Lh" is written in blue, with a blue bracket above it grouping  $\Pr(D)$ . A green bracket at the bottom groups the entire denominator  $\Pr(D)$ .

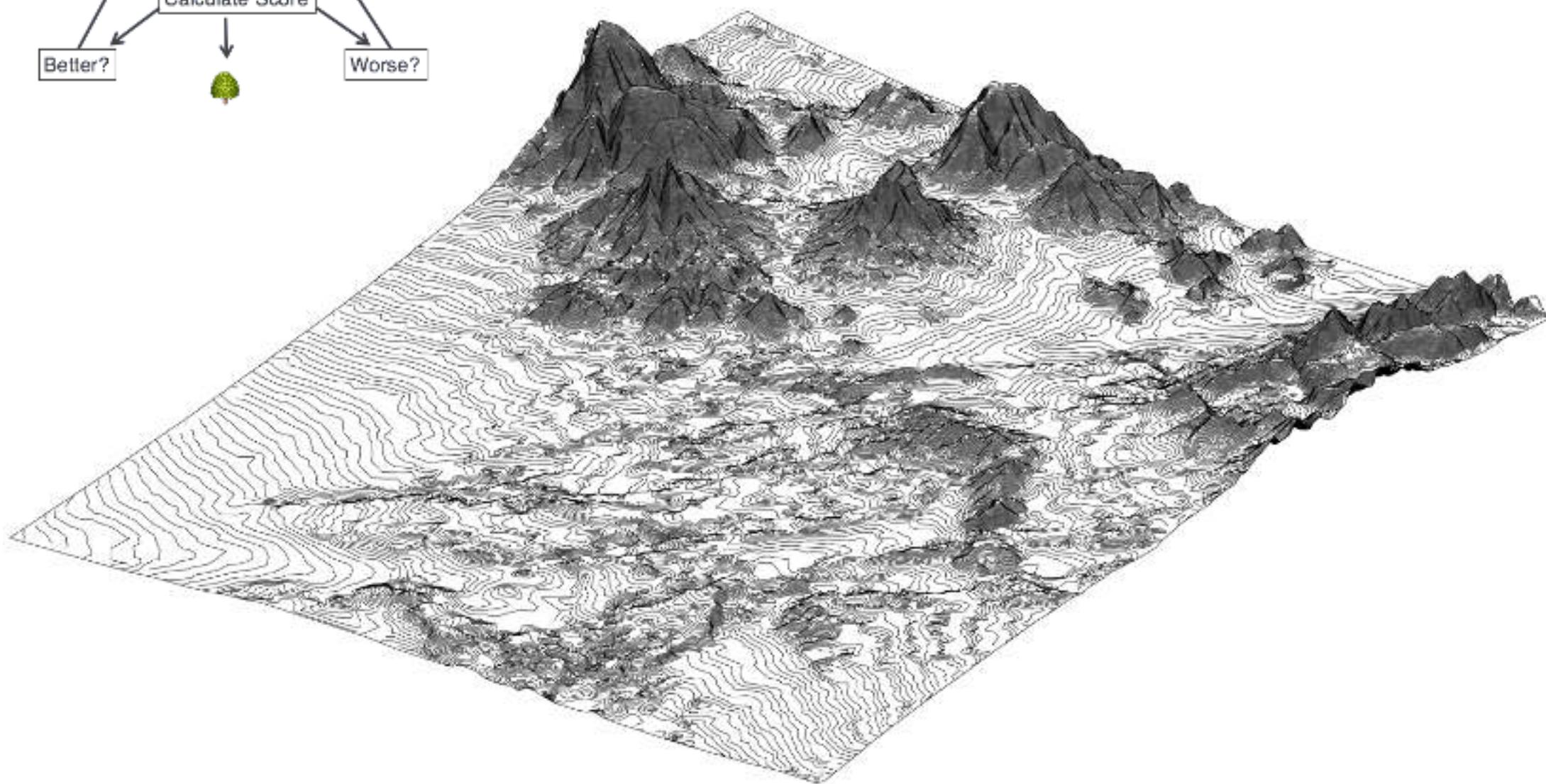
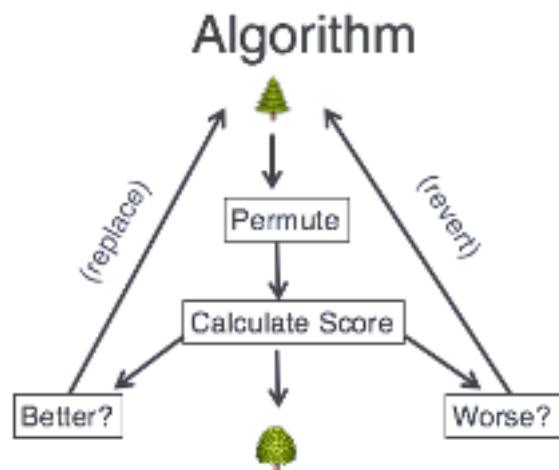
General intro: McElreath “Statistical Rethinking” (<http://xcelab.net/rm/statistical-rethinking/>)

Phylogenetics: Brown “The State of Bayesian Phylogenetics: Bayes for the Uninitiated” (<https://tinyurl.com/ly2cvfd>)

# Treespace

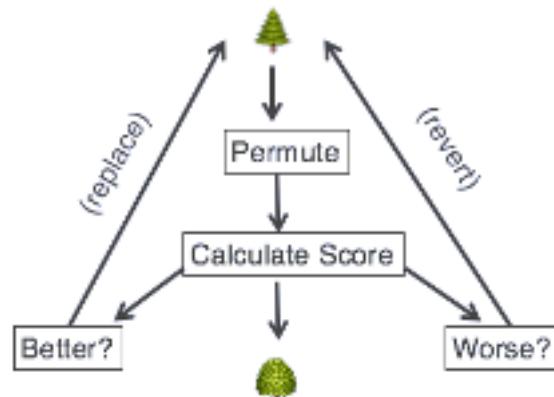


# MCMC

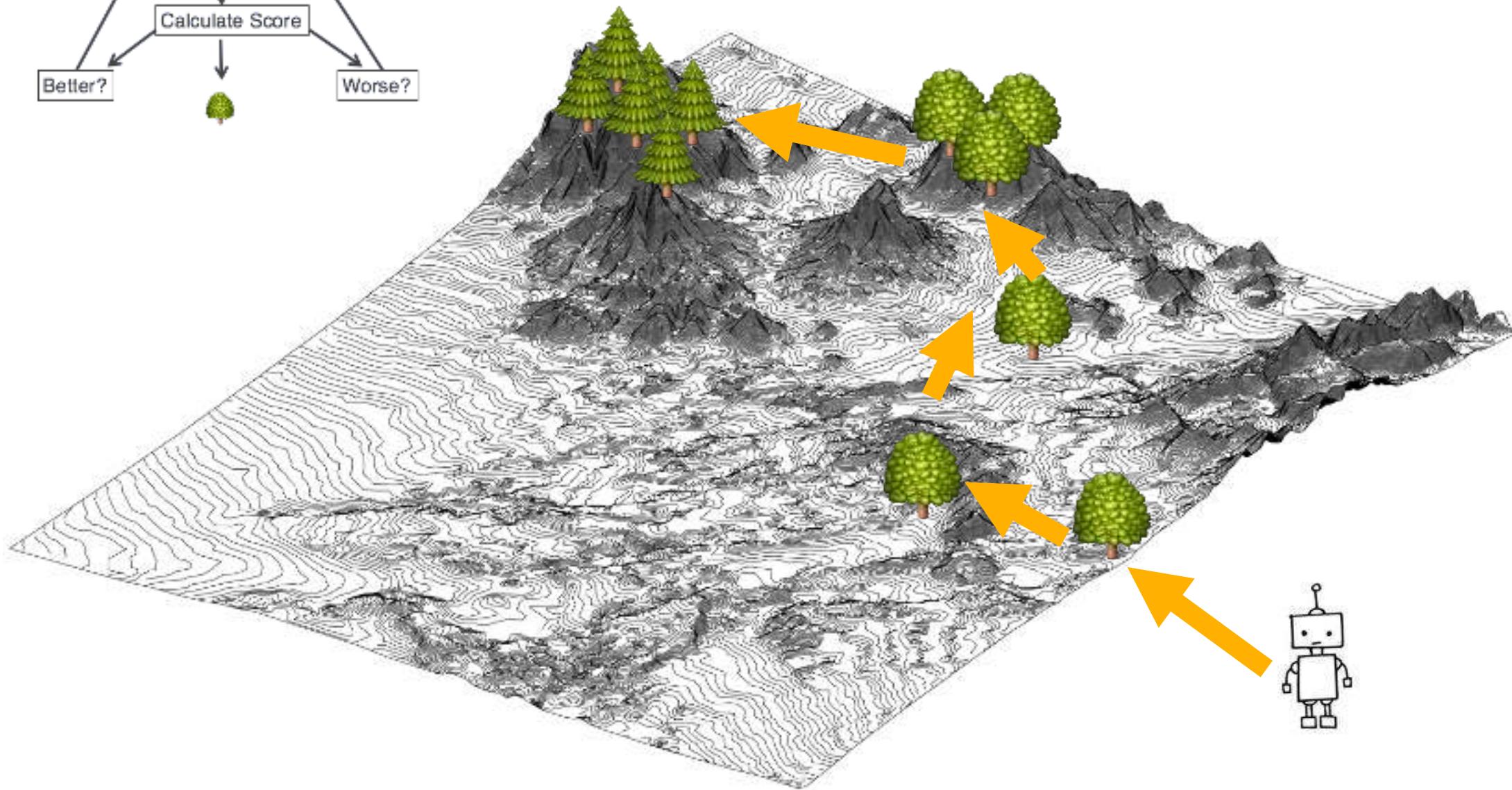
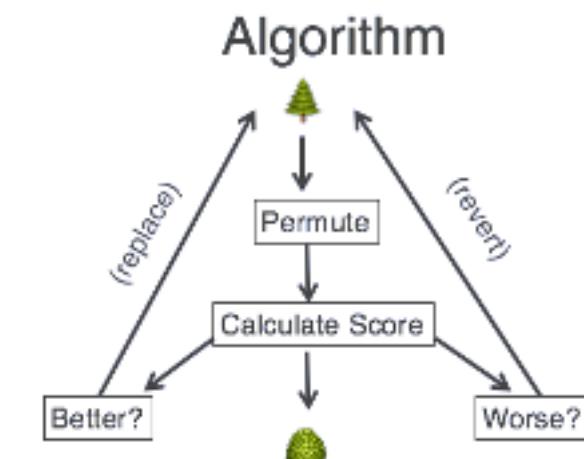


# MCMC

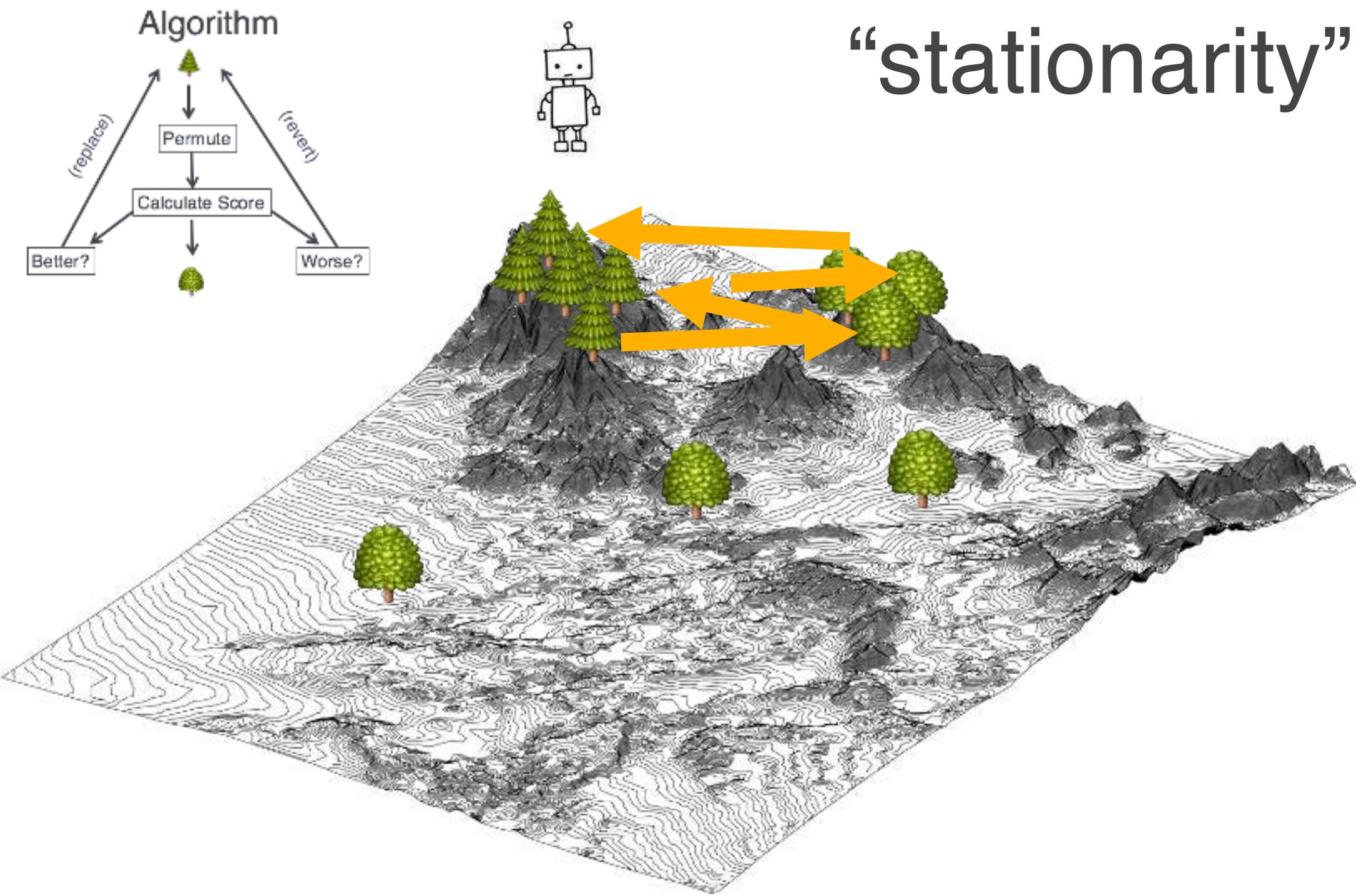
Algorithm



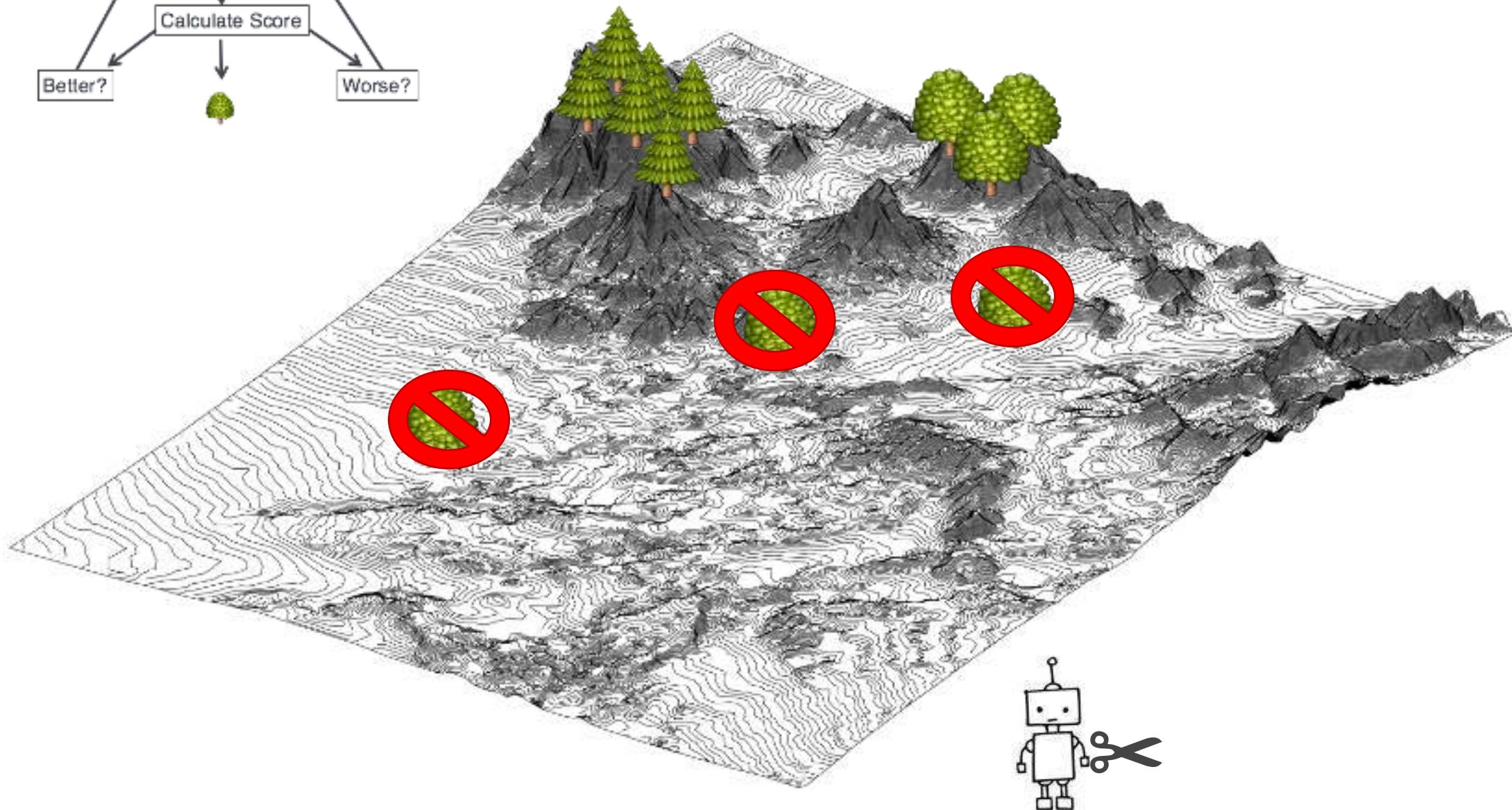
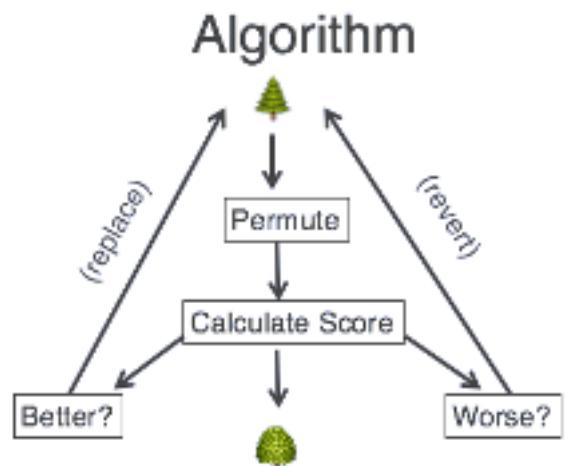
# MCMC



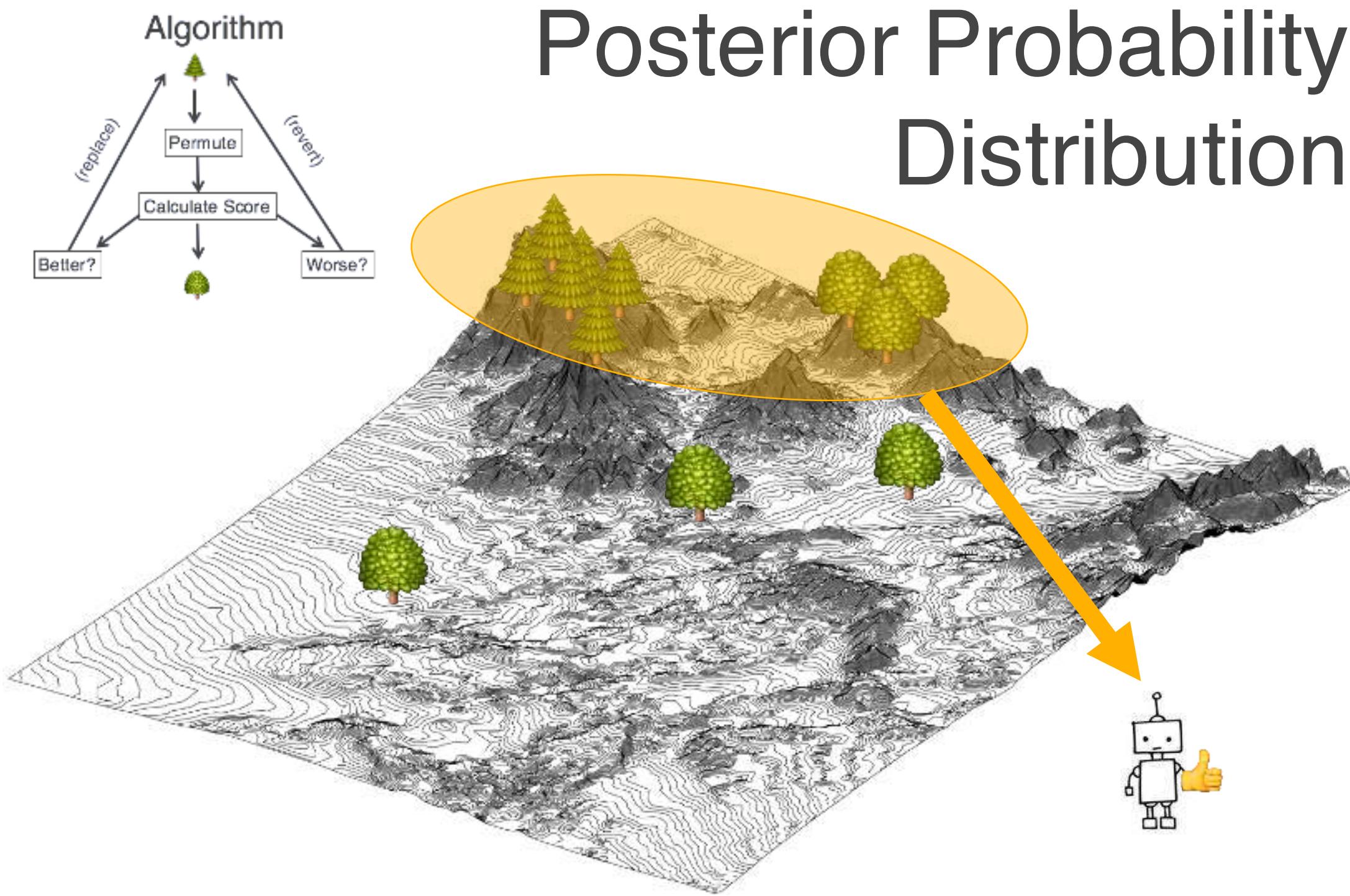
# “stationarity”



# “burn-in”



# Posterior Probability Distribution



# Posterior Probability Distribution: Uto-Aztecán Languages

Covarion & Relaxed Clock + etc  
Sampled 10,000 trees.



# Posterior Probability Distribution: Uto-Aztec Languages

Covarion & Relaxed Clock + etc

Sampled 10,000 trees.

Draw each one (Densitree, Bouckaert '10)

Some well supported regions

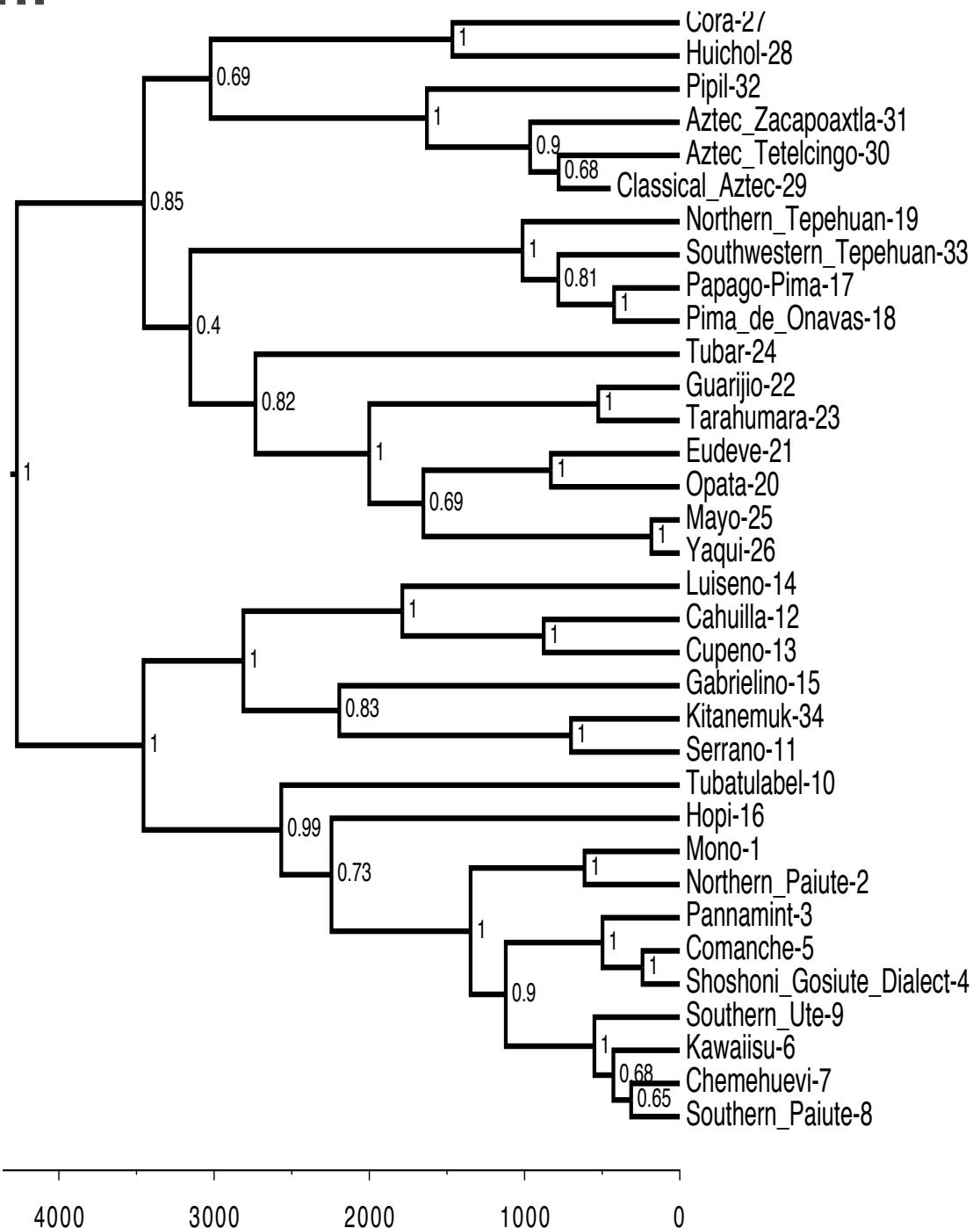
Some less...

Some conflict



# Posterior Probability Distribution: Uto-Aztec Languages

Maximum Clade Credibility Tree  
Reduce posterior to single summary tree  
Values are “Posterior probabilities”  
- 1.00 = present in ALL trees in P.  
- 0.50 = only half of the trees have this.



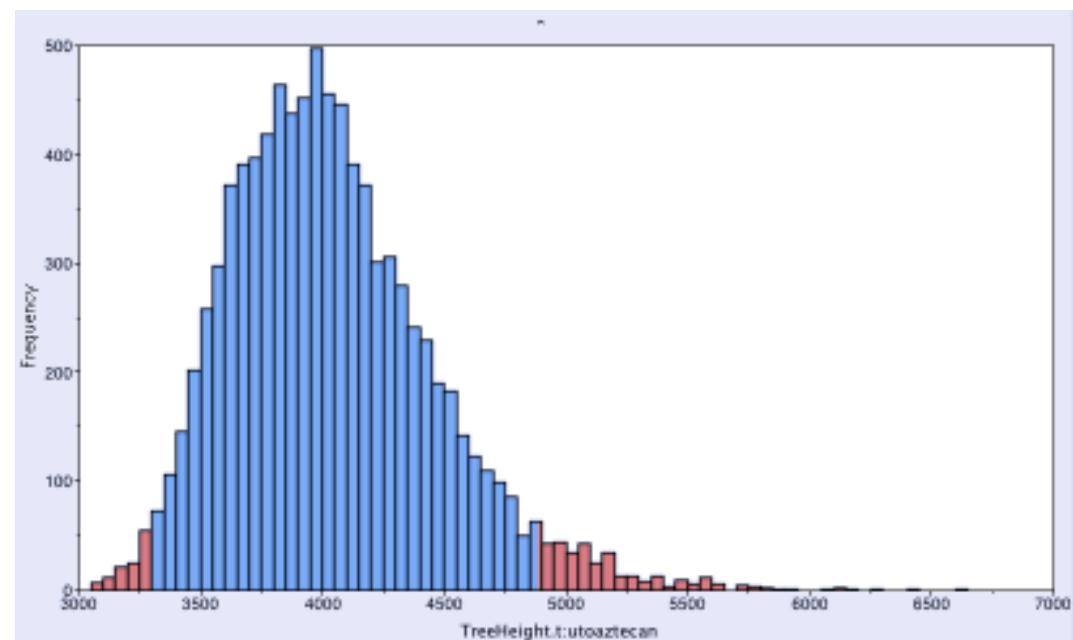
# Posterior Probability Distribution: Uto-Aztec Languages

...and any statistic we calculate  
has a PP distribution.

- mean
- 95% highest posterior density interval  
(95% probability that  $\mu$  falls in range)



Mean = 4334 years



95% HPD: 3302 – 4884 years

# Whirlwind tour

- Disciplinary parallels & similarity of questions
- Lexicostatistics and Glottochronology
- Parallel evolution of methods in Biology
- Maximum Parsimony
- Maximum Likelihood
- Bayesian Phylogenetics
- Theory
- Practical: My BEAST tutorial tomorrow.
  - <https://taming-the-beast.github.io/>

