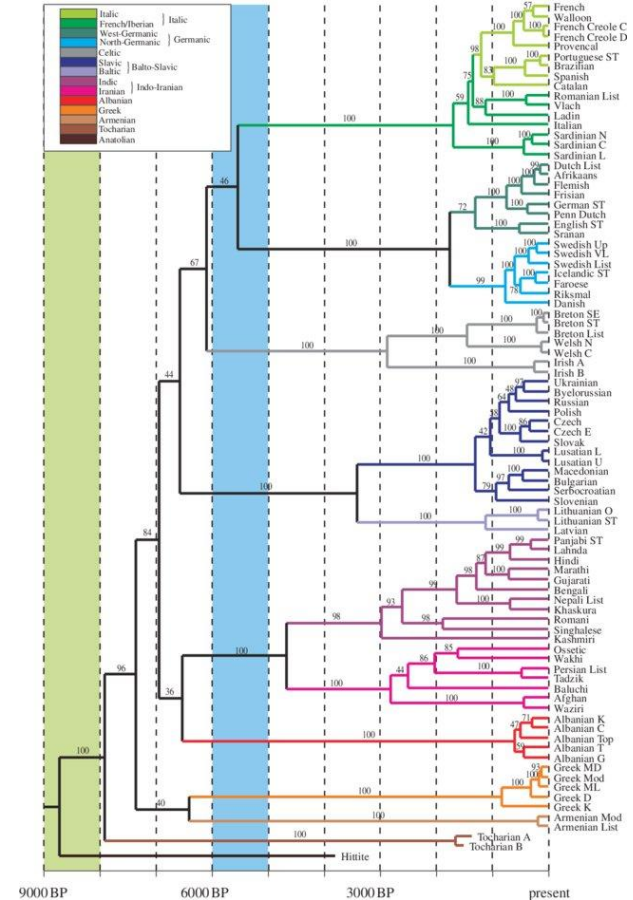

Illustrating Bayesian inference for stemmatology

— Tiago Tresoldi —
CEoT - Uppsala universitet

Introduction

- Advances in computational historical linguistics with evolutionary methods from biology
- Very different from NJ and UPGMA, a bit more similar to Maximum Likelihood
- A certain skepticism in stemmatology, which is right: the *tools* are not ready



What it does

- Given the *data* and *evolutionary model*, the tools collect thousands of trees, constantly looking for the best one
- The “goodness” of a tree is given by the probability of obtaining the *data* we observed when assuming an *underlying model*
- Everything is in terms of probability: not how much to “trust” a single result considered the best, but how *likely* that result is
- The collection of trees is reduced to a single tree with different strategies

Bayesian inference - I

Priors: original probability for the model parameters

All parameters have priors!

$$P(\text{model} \mid \text{data}) = \frac{P(\text{data} \mid \text{model}) P(\text{model})}{P(\text{data})}$$

Bayesian inference - II

Likelihood: probability of data given the parameters

$$P(\text{model} \mid \text{data}) = \frac{P(\text{data} \mid \text{model}) P(\text{model})}{P(\text{data})}$$

Bayesian inference - III

Evidence: probability in any combination of parameters

$$P(\text{model} \mid \text{data}) = \frac{P(\text{data} \mid \text{model}) P(\text{model})}{P(\text{data})}$$

Bayesian inference - IV

Posterior: updated probability

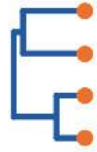
$$P(\text{model} \mid \text{data}) = \frac{P(\text{data} \mid \text{model}) P(\text{model})}{P(\text{data})}$$

Data

$$P(\text{model} \mid \begin{matrix} \text{ACAC} \dots \\ \text{TCAC} \dots \\ \text{ACAG} \dots \end{matrix}) = \frac{P(\begin{matrix} \text{ACAC} \dots \\ \text{TCAC} \dots \\ \text{ACAG} \dots \end{matrix} \mid \text{model}) P(\text{model})}{P(\begin{matrix} \text{ACAC} \dots \\ \text{TCAC} \dots \\ \text{ACAG} \dots \end{matrix})}$$

Model

MODEL



Genealogy



Population Model



Site Model



Molecular Clock Model

Model - II

Tree prior: How likely is a tree, given a demographic model

TREE

Realisation of a
stochastic process

$$P(\text{Tree} \mid \text{Population Model})$$

POPULATION MODEL

Describes the population
dynamics (growth of the tree)

Model - III

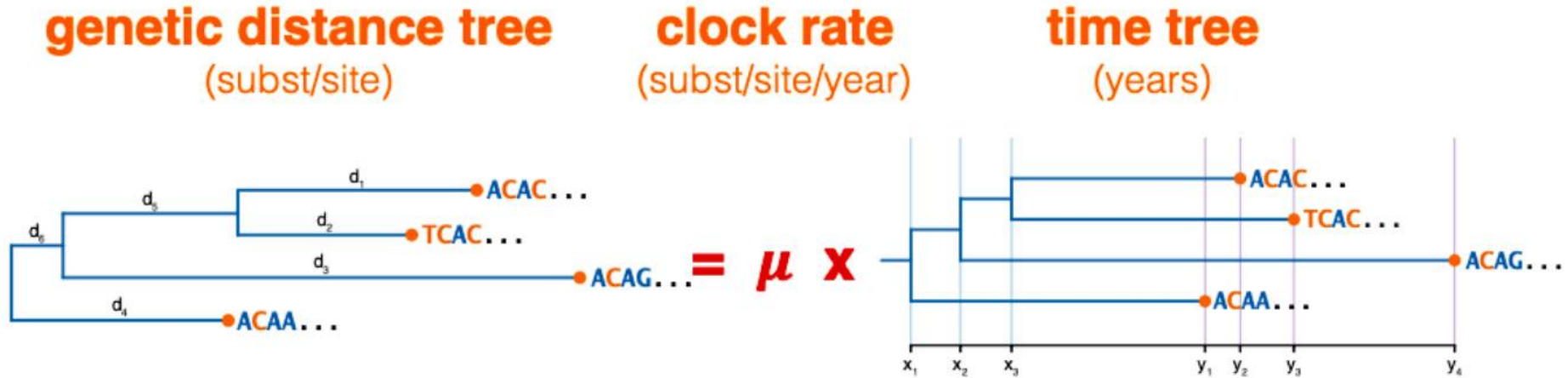
Substitution model: tree likelihood sums all substitution histories leading to the observed sequences



$$P\left(\begin{array}{l} ACAC\dots \\ TCAC\dots \\ ACAG\dots \end{array} \mid \begin{array}{c} \text{Tree} \end{array} \begin{array}{c} \text{Matrix} \end{array} \right)$$

Model - IV

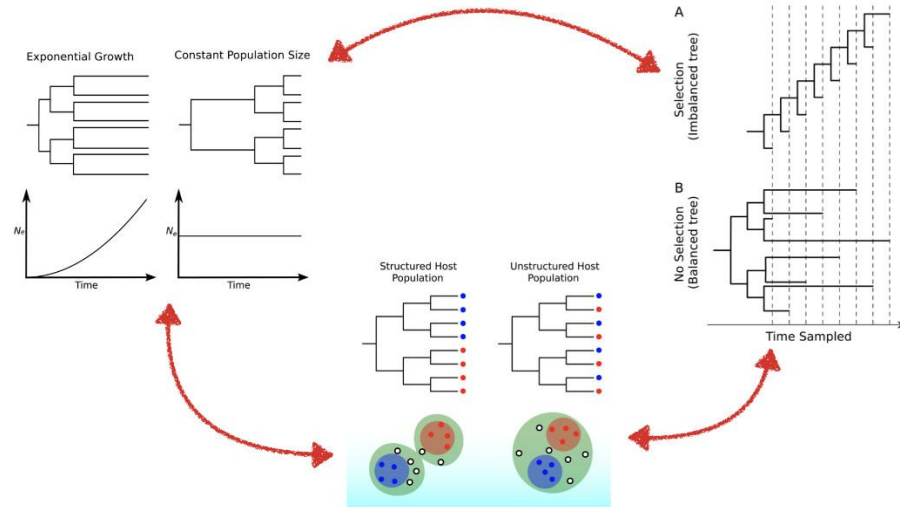
Molecular clock, strict/relaxed/random



Source: Valenzuela (2021), TTB Online first steps

Model - V

Different population dynamics give *very* different trees!



Volz *et al.* **PLoS Comp Biol** 2013
Grenfell *et al.* **Science** 2004

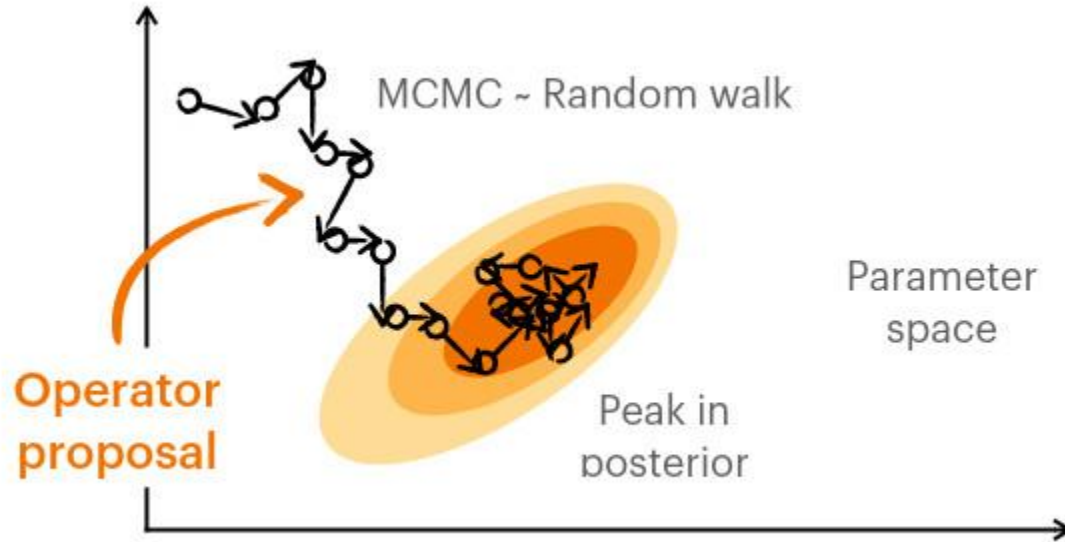
Source: Valenzuela (2021), TTB Online first steps

Model - VI

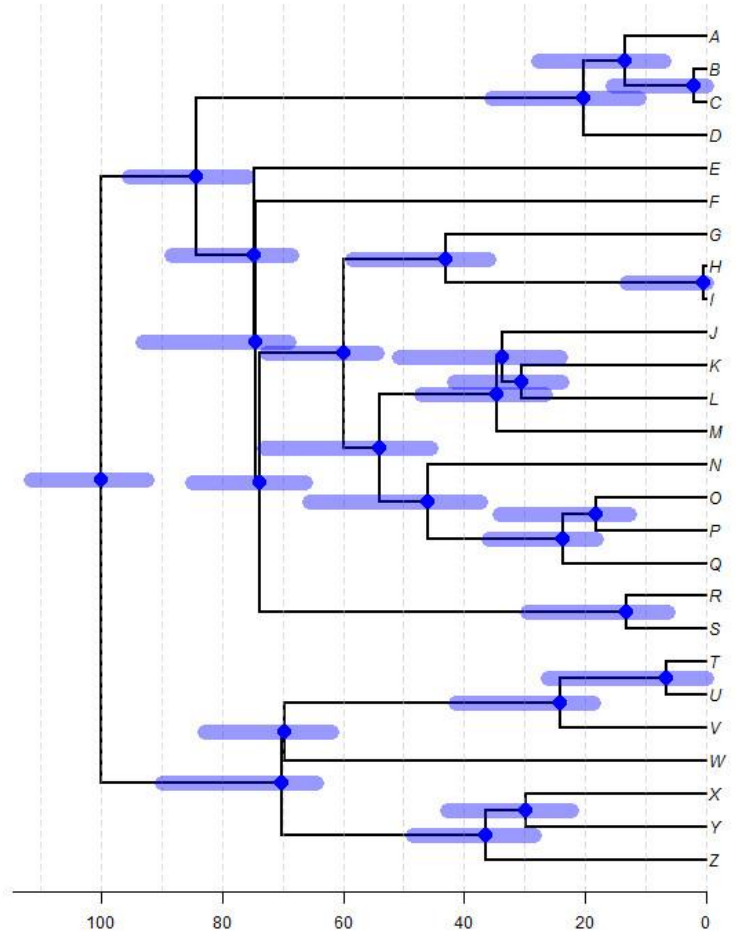
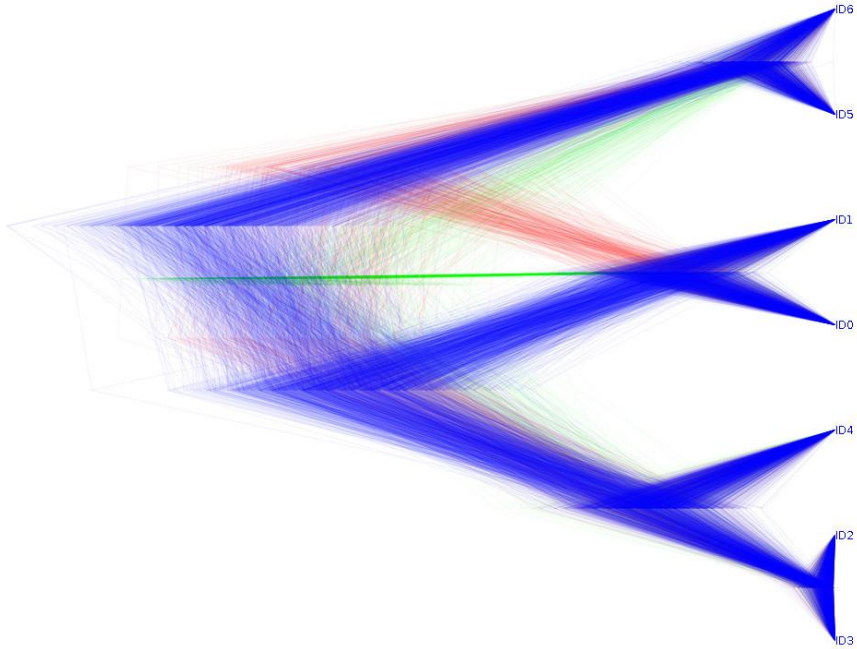
$$P(\text{Diagram 1} | \text{ACAC... TCAC... ACAG...}) = \frac{P(\text{Diagram 2} | \text{Diagram 1}) P(\text{Diagram 1})}{P(\text{ACAC... TCAC... ACAG...})}$$

The diagram in the numerator consists of four components: a tree diagram, a circular diagram with two arrows, a 4x4 grid of colored dots (blue and orange), and a circular diagram with a red dot and a blue arrow.

MCMC



Tree summary



A	Eminentiora	prolixarum	arborum	culmina	perindeque	distenta	acuto	sonitu	resultabant
B	Eminentiora	prolixarum	abietum	cacumina	perindeque	distantia	acuto	sonitu	resultabant
C	Altiora	prolixarum	arborum	culmina	perindeque	distenta	acuto	sonitu	resonabant
D	Altiora	promissarum	arborum	culmina	perindeque	distenta	acutissimo	sonitu	resonabant
E	Altiora	prolixarum	arborum	culmina	perindeque	distenta	acuto	sono	resonabant
F	Altiora	prolixarum	arborum	fulmina	perindeque	et distenta	acuto	tinnitu	resonabant
G	Altiora		arborum	culmina	perindeque	discreta	acuto	sono	resonabant
H	Altiora	prolixarum	arborum	culmina	proptereaque	distenta	acuto	sono	resonabant
I	Eminentiora	promissarum	abietum	culmina	perindeque	distenta	acutissimo	sonitu	resultabant
J	Eminentiora	prolixarum	abietum	cacumina	per insignem	distantiam	acuto	sonitu	resultabant

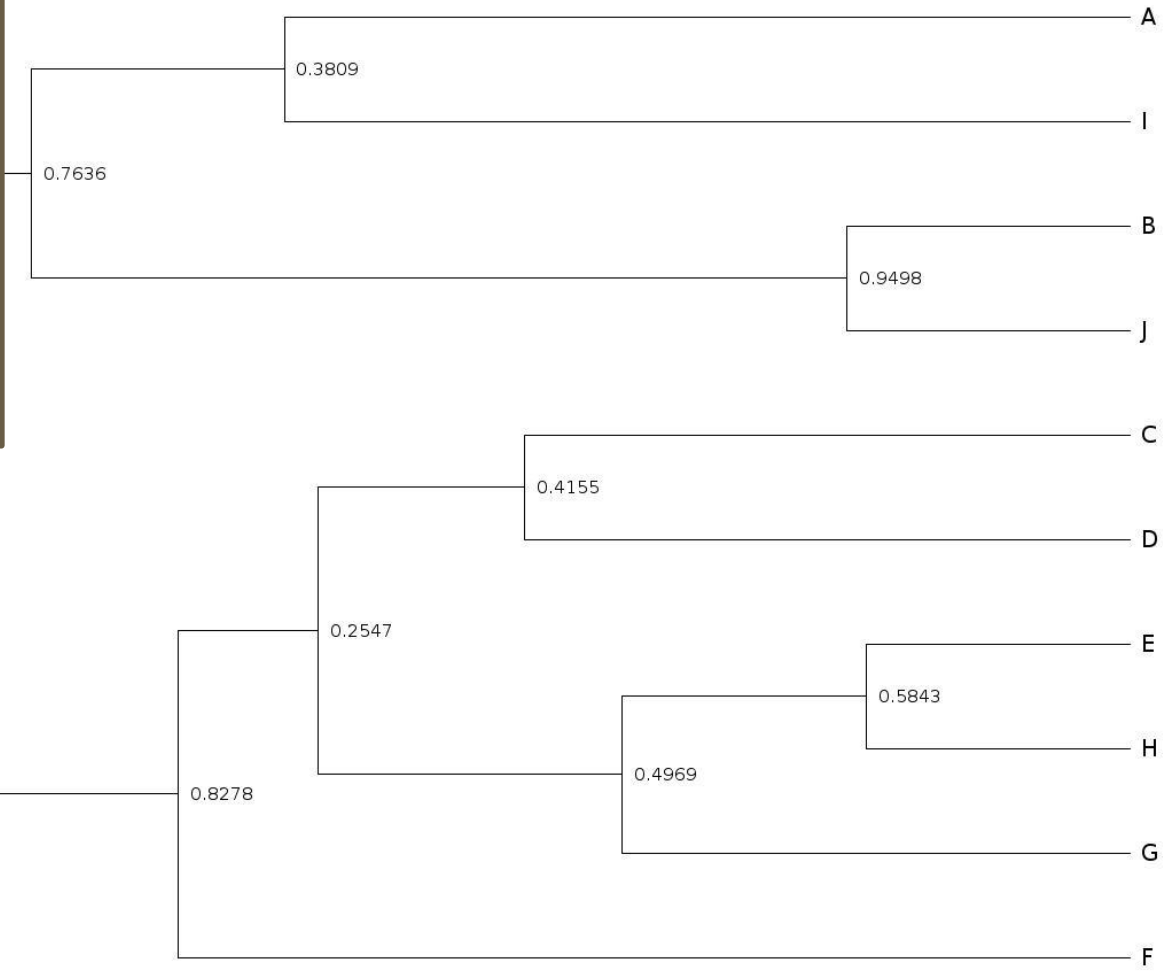
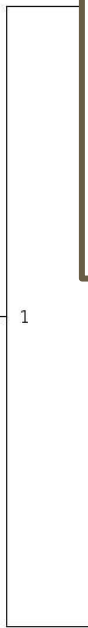
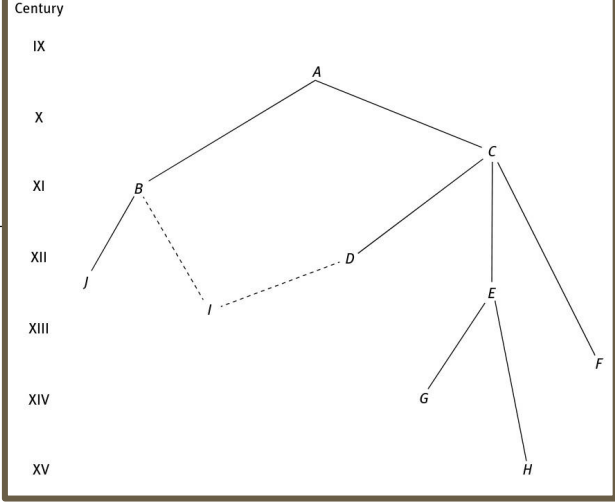
A	Eminentiora	prolixarum	arborum	culmina	perindeque	distenta	acuto	sonitu	resultabant
B	Eminentiora	prolixarum	abietum	cacumina	perindeque	distantia	acuto	sonitu	resultabant
C	Altiora	prolixarum	arborum	culmina	perindeque	distenta	acuto	sonitu	resonabant
D	Altiora	promissarum	arborum	culmina	perindeque	distenta	acutissimo	sonitu	resonabant
E	Altiora	prolixarum	arborum	culmina	perindeque	distenta	acuto	sono	resonabant
F	Altiora	prolixarum	arborum	fulmina	perindeque	et distenta	acuto	tinnitu	resonabant
G	Altiora		arborum	culmina	perindeque	discreta	acuto	sono	resonabant
H	Altiora	prolixarum	arborum	culmina	proptereaque	distenta	acuto	sono	resonabant
I	Eminentiora	promissarum	abietum	culmina	perindeque	distenta	acutissimo	sonitu	resultabant
J	Eminentiora	prolixarum	abietum	cacumina	per insignem	distantiam	acuto	sonitu	resultabant

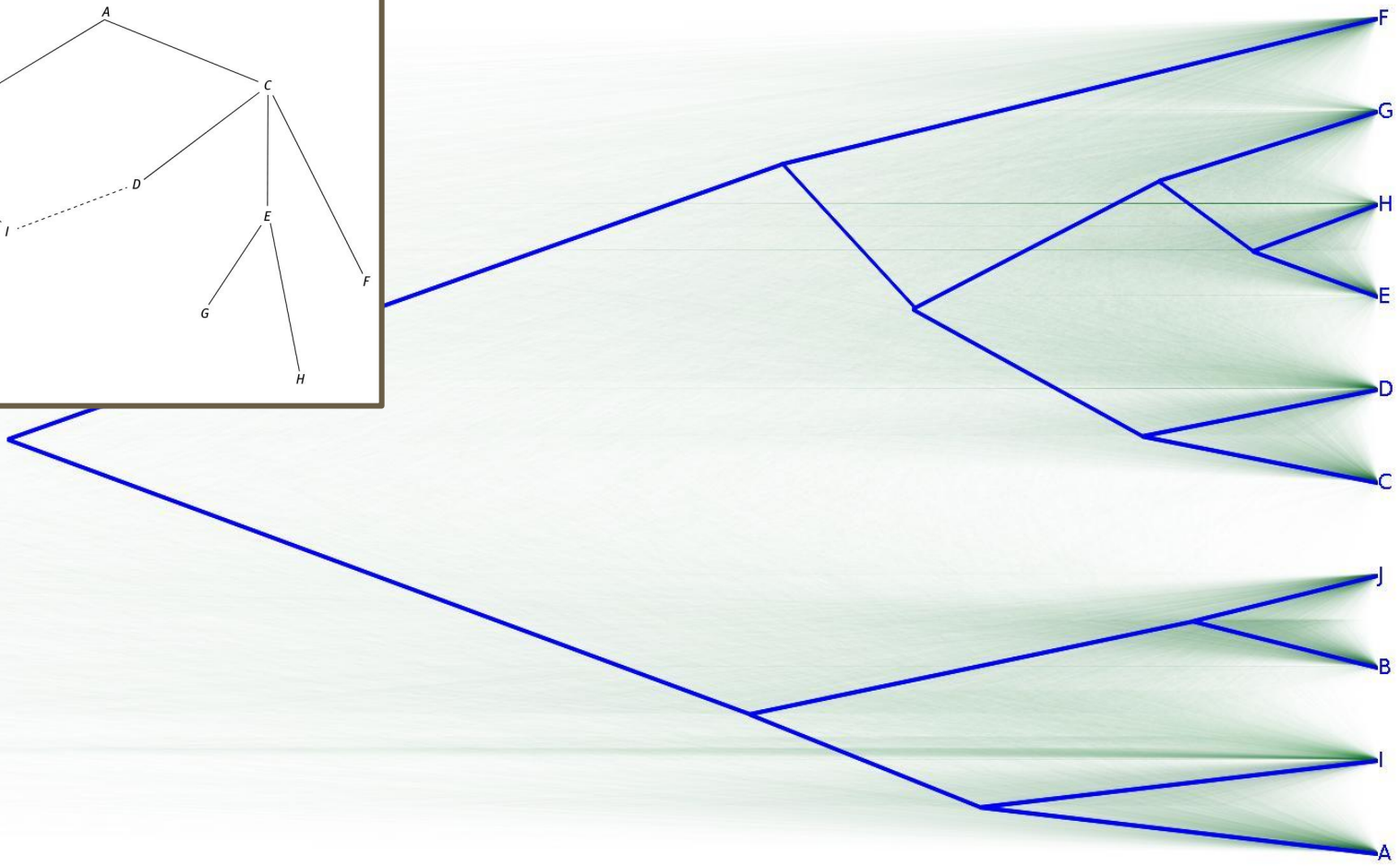
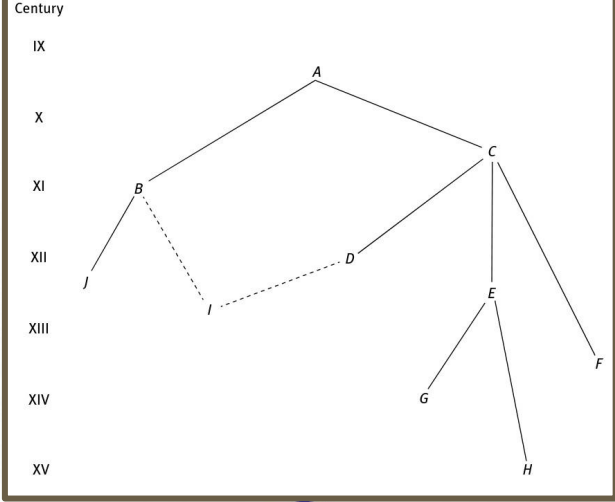
A	A	C	F	H	K	N	S	U	X
B	A	C	G	I	K	O	S	U	X
C	B	C	F	H	K	N	S	U	Y
D	B	D	F	H	K	N	T	U	Y
E	B	C	F	H	K	N	S	V	Y
F	B	C	F	J	K	P	S	W	Y
G	B	E	F	H	K	Q	S	V	Y
H	B	C	F	H	L	N	S	V	Y
I	A	D	G	H	K	N	T	U	X
J	A	C	G	I	M	R	S	U	X

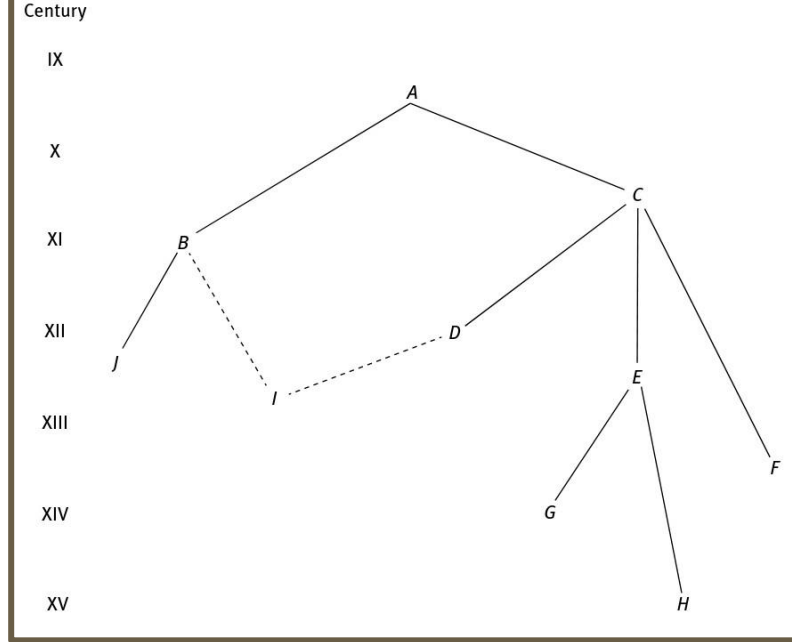
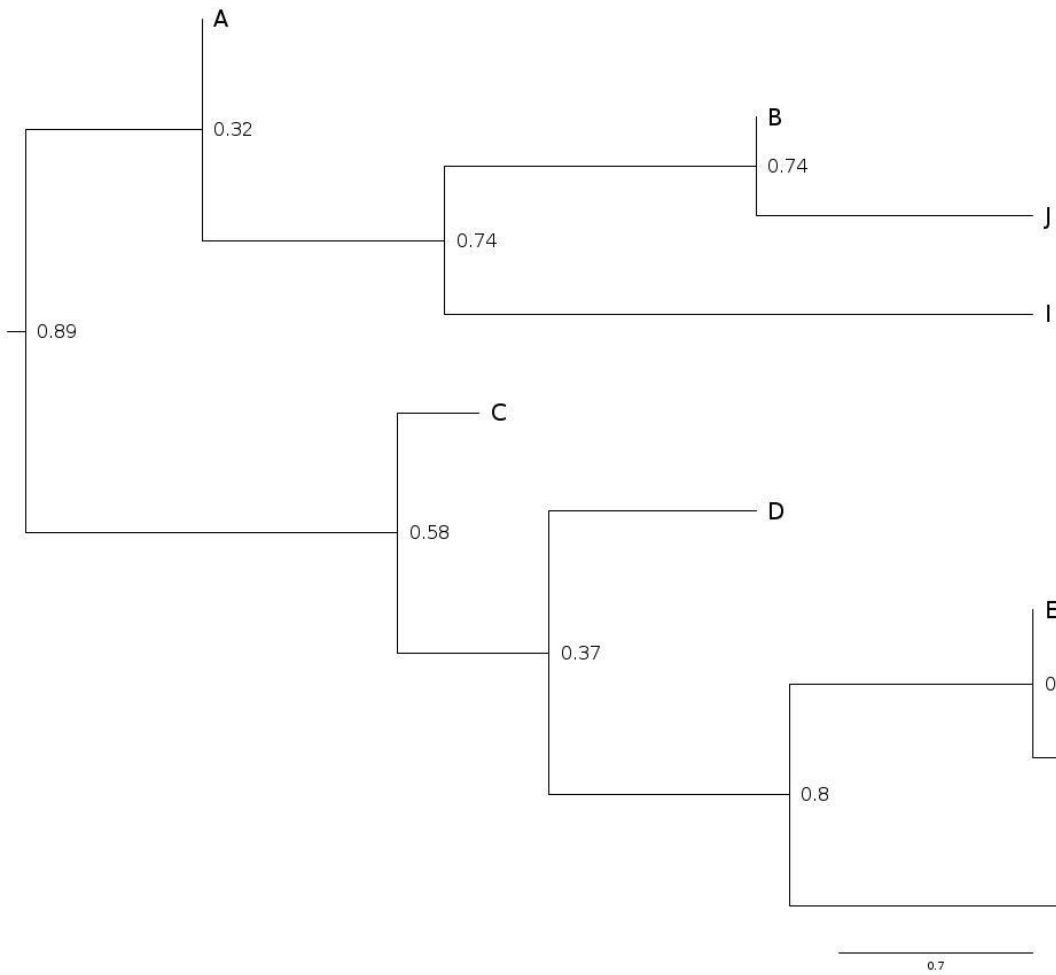
```

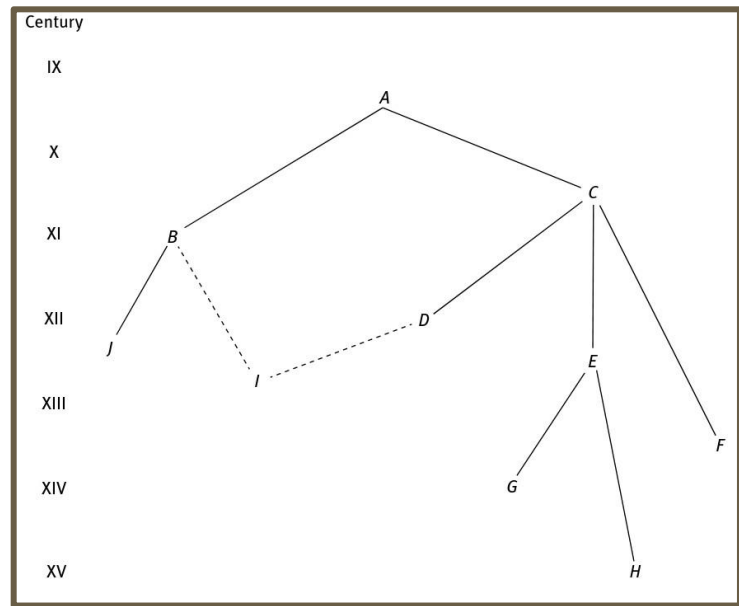
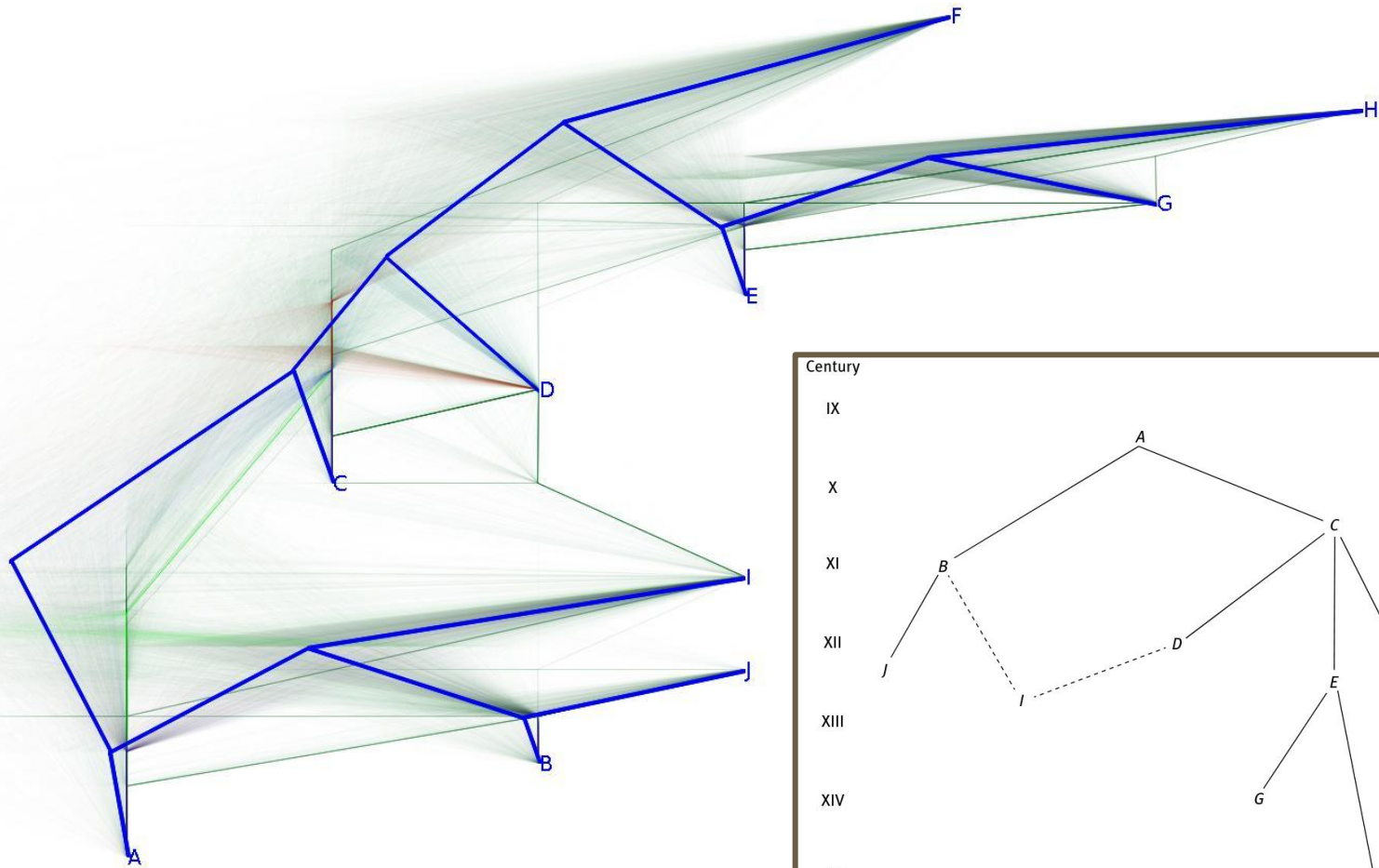
1 #NEXUS
2
3 BEGIN DATA;
4     DIMENSIONS NTAX=10 NCHAR=34;
5     FORMAT DATATYPE=STANDARD MISSING=? GAP=- SYMBOLS="01";
6     CHARSTATELABELS
7         1 CHAR_1_ASCERT,
8         2 CHAR_1_Altiora,
9         3 CHAR_1_Eminentiora,
10        (...)
11        34 CHAR_9_resultabant
12;
13 MATRIX
14 A 0010010001001000100000100010100001
15 B 0010010010010000100010000010100001
16 C 0100010001001000100000100010100010
17 D 0100001001001000100000100100100010
18 E 0100010001001000100000100010010010
19 F 010001000100100100000010010001010
20 G 0100100001001000100100000010010010
21 H 0100010001001000010000100010010010
22 I 0010001010001000100000100100100001
23 J 0010010010010001000001000010100001
24;
25 begin assumptions;
26     charset CHAR_1 = 1-3;
27     charset CHAR_2 = 4-7;
28     charset CHAR_3 = 8-10;
29     charset CHAR_4 = 11-14;
30     charset CHAR_5 = 15-18;

```







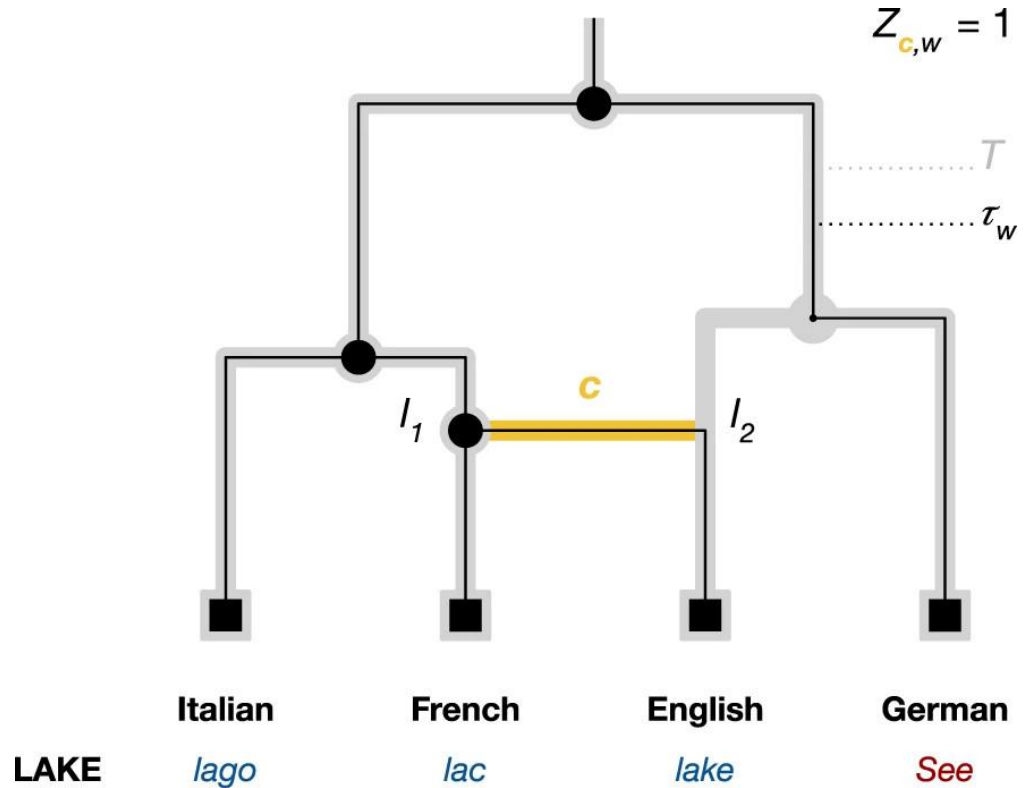


Tools are not ready: some issues

- Tree prior: the biological models don't reflect manuscript evolution (e.g. binary division)
- Molecular clock: time works differently when making copies from texts
- Site models: transitions tend to be very asymmetric, it is hard to encode expert knowledge
- We don't really know what most parameters should be like

Where we are

- BEAST2 and MrBayes are limited in their offerings
- BEAST2 can be expanded, and RevBayes could be a solution
- We are currently experimenting with the detection on synthetic data (of which we know the “real” stemma) and the Divine Comedy (with data from Shaw 2011)



Source: Neureiter et al. (2022)



UPPSALA
UNIVERSITET

Thank you!

tiago.tresoldi@lingfil.uu.se