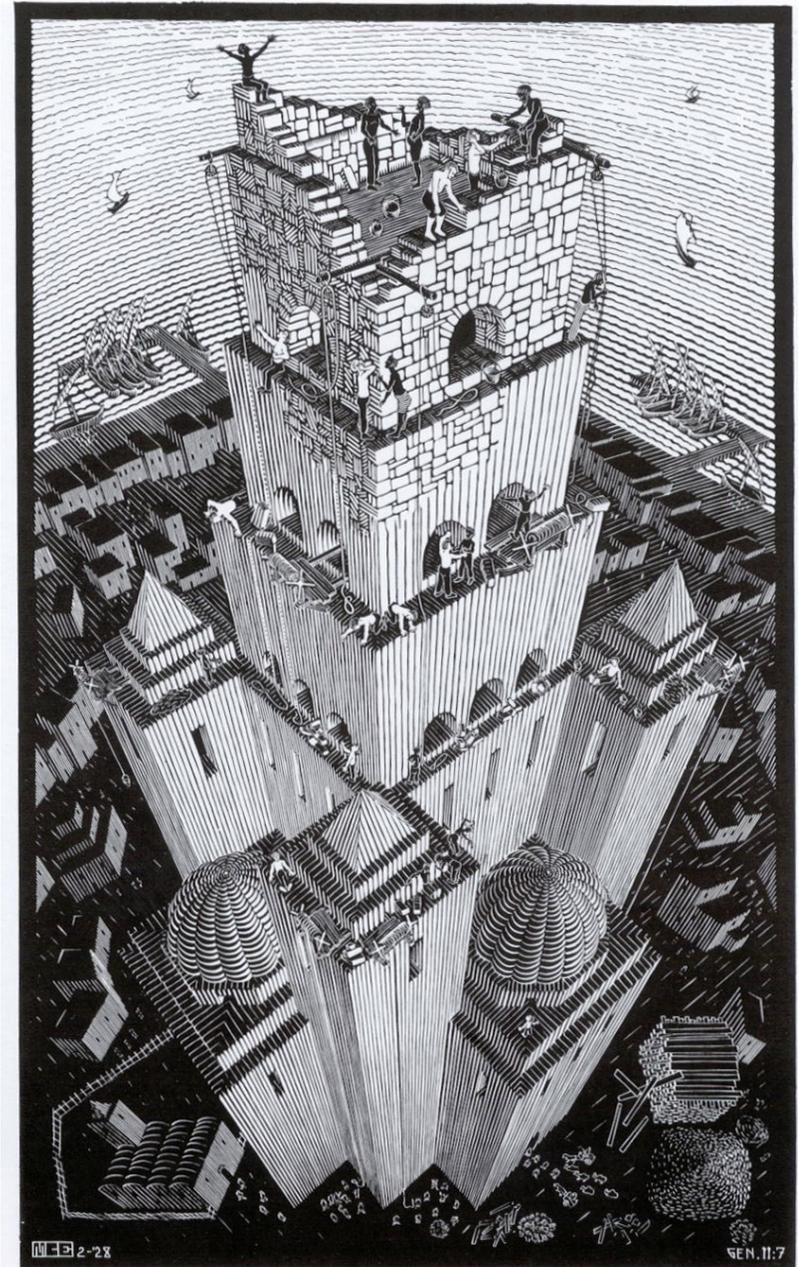


# Network models of sound change

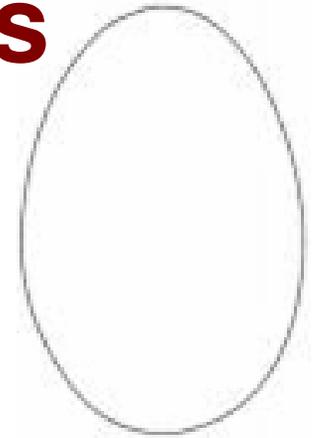
Dan Hruschka  
Santa Fe Institute  
[dhrusch@santafe.edu](mailto:dhrusch@santafe.edu)

in collaboration with Tanmoy Bhattacharya,  
Eric Smith & Jon Wilkins

“Statistical Inference for Complex Networks”  
December 2008



# “Egg” in 24 Turkic languages



- jumurtGa
- simit
- umurxa

j	j	j	j	j	j	j	j	j	n	n	d	j	s	s	h	č		ž	ž	j	j	z
u	u	u	o	u	u	o	u	u	ı	ı	ı	u	ы	ı	ı	u	u	u	u	u	o	u
m	m	m	m	m	m	m	m	m	m	b	m	m	m	m	m		m	m	m	m	m	n
u	u	u	ı	u	u	u	u	u	ı	ı	ı	u	a	ī	ī	u	u	u	ı	ı	o	u
r	r	r	r	r	r		r	r	r	r	r	r	r			r	r	r	r	r	r	r
t	t	t		t	t	t	t	t		t	t		d	t	t			t	t	t	t	t
y	y		q	y	q		G	x	q	q	q					y	x	q	q	q	q	q
a	a	a	a	a	a		a	a	a	a	a	a				a	a	a	a	a	a	a

# The Plan

**Question**—What theories of **sound change** best account for observed diversity in related languages?

- Brief introduction—concepts, units
- Specify one model in a likelihood framework
- Discuss plans for specifying other models and raise issues of inference and model selection



# Many Theories

- No expectation of what changes to what
- Regularity, but no expectation of what changes to what
- Random walk through feature space
- Lenition ( $t \rightarrow ts \rightarrow s \rightarrow h$ )
- Articulatory Reduction
- Other accounts
  - “An eft” becomes “a newt”
  - Sound symbolism—'stamp', 'stomp', 'tamp', 'tromp', 'tramp'
  - Constraints on language design (e.g. optimality theory)

# An Opportunity

- “Although no comprehensive study of sound change that would allow us to distinguish common from uncommon innovations has ever been undertaken, historical linguists have acquired a sense of what kinds of change are likely to occur” (Blust 2004)

# Common Methods

- **Sociolinguistics**—Study dialectical variation and change over short periods of time
- **Psycholinguistics**—Study variation in production and perception

## **Advantage**

- Close view of change in action

## **Disadvantage**

- Limited to frequent or fast changes

# Historical Linguistics

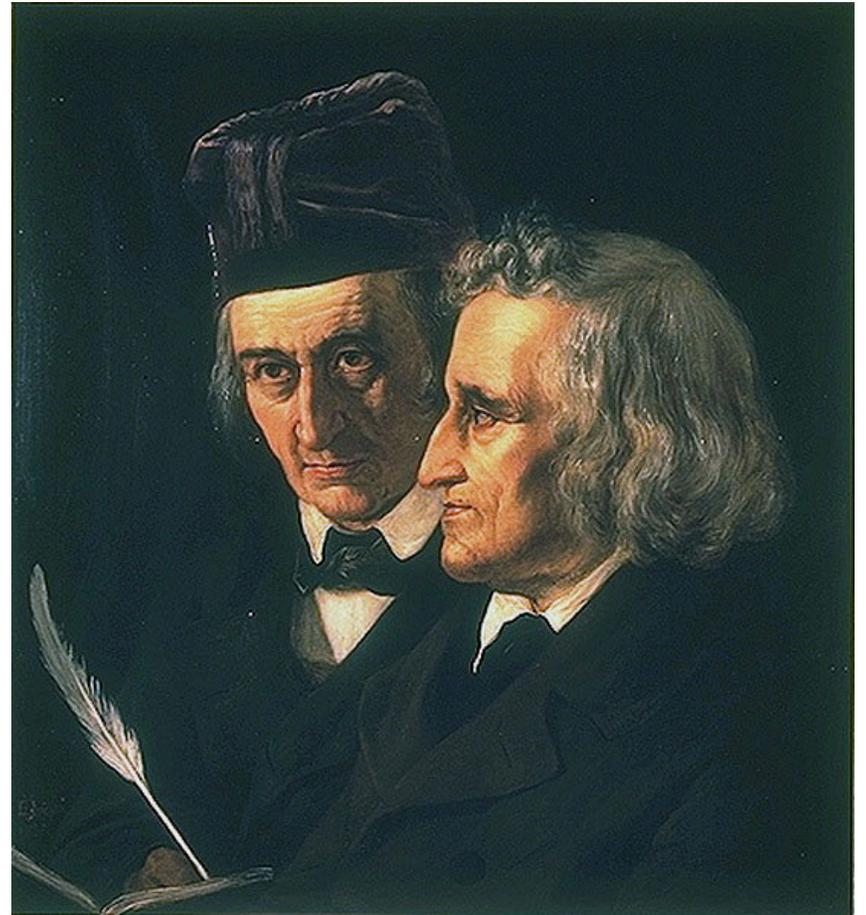
Study extant between-language variation and ancient recorded languages (Greek, Sanskrit)

## Advantages

- Longer time scales

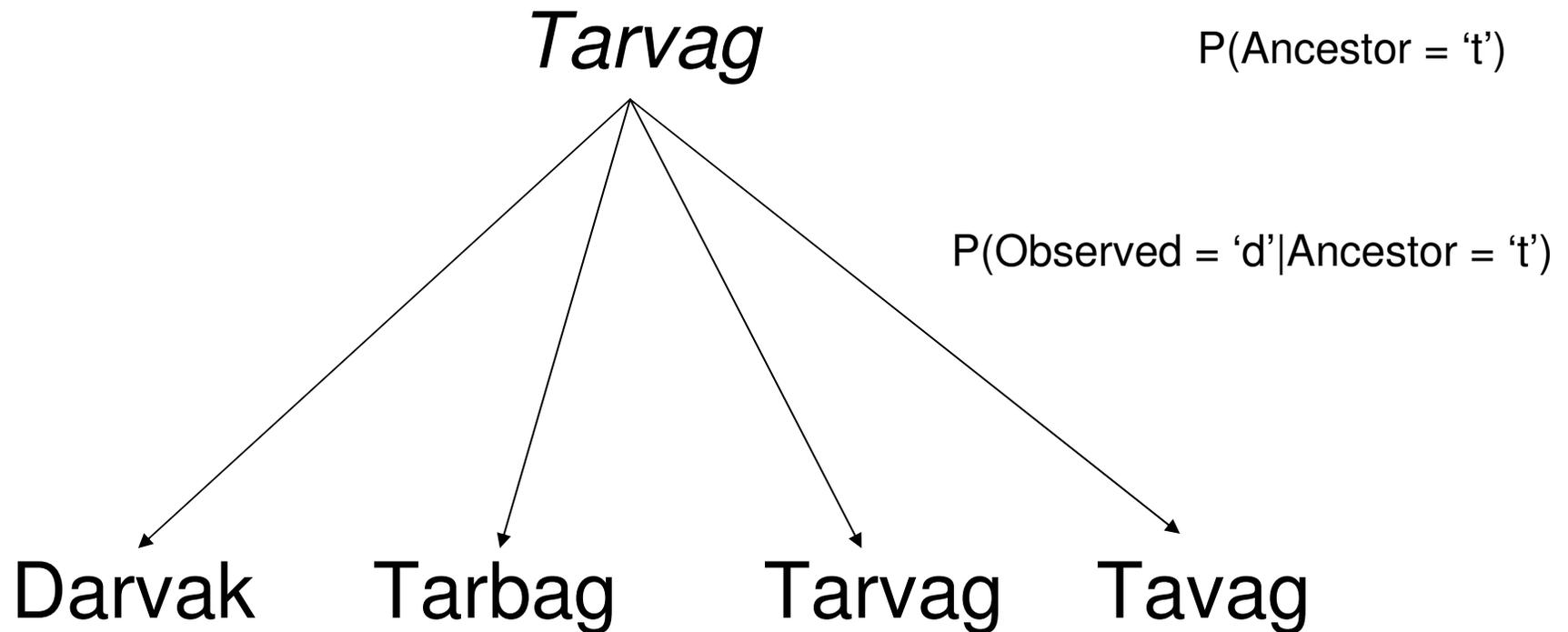
## Disadvantages

- Can't see process
- **Requires some way of making inferences about generative models**



Jacob and Wilhelm Grimm (1855)

# A Model



1. **Align words**
2. **Estimate probabilities of sound change**
3. **Estimate probabilities of ancestral sounds**
4. Identify cognates
5. Infer historical relationships

# Likelihood Function

$$p(L_1, L_2, \dots, L_N \mid p(A), p(s \mid A), \text{alignment}) =$$
$$\prod_{\text{CognateClasses}} \prod_{\text{positions}} \sum_{\text{Ancestors}} p(A) \prod_L p(s \mid A)$$

$P(s|A)$  = probability of observed sound given an ancestor

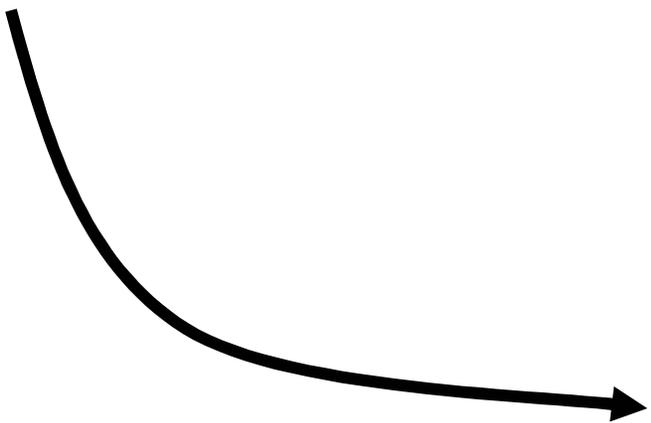
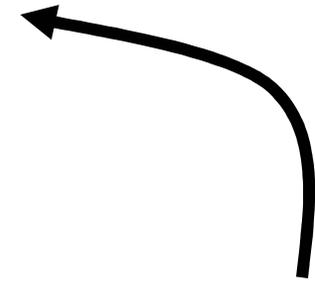
$P(A)$  = probability of ancestral sound

Some assumptions were required to get here

# Estimation Strategy

$P(\text{Observed sound} | \text{Ancestral sound})$

$P(\text{Ancestral sound})$

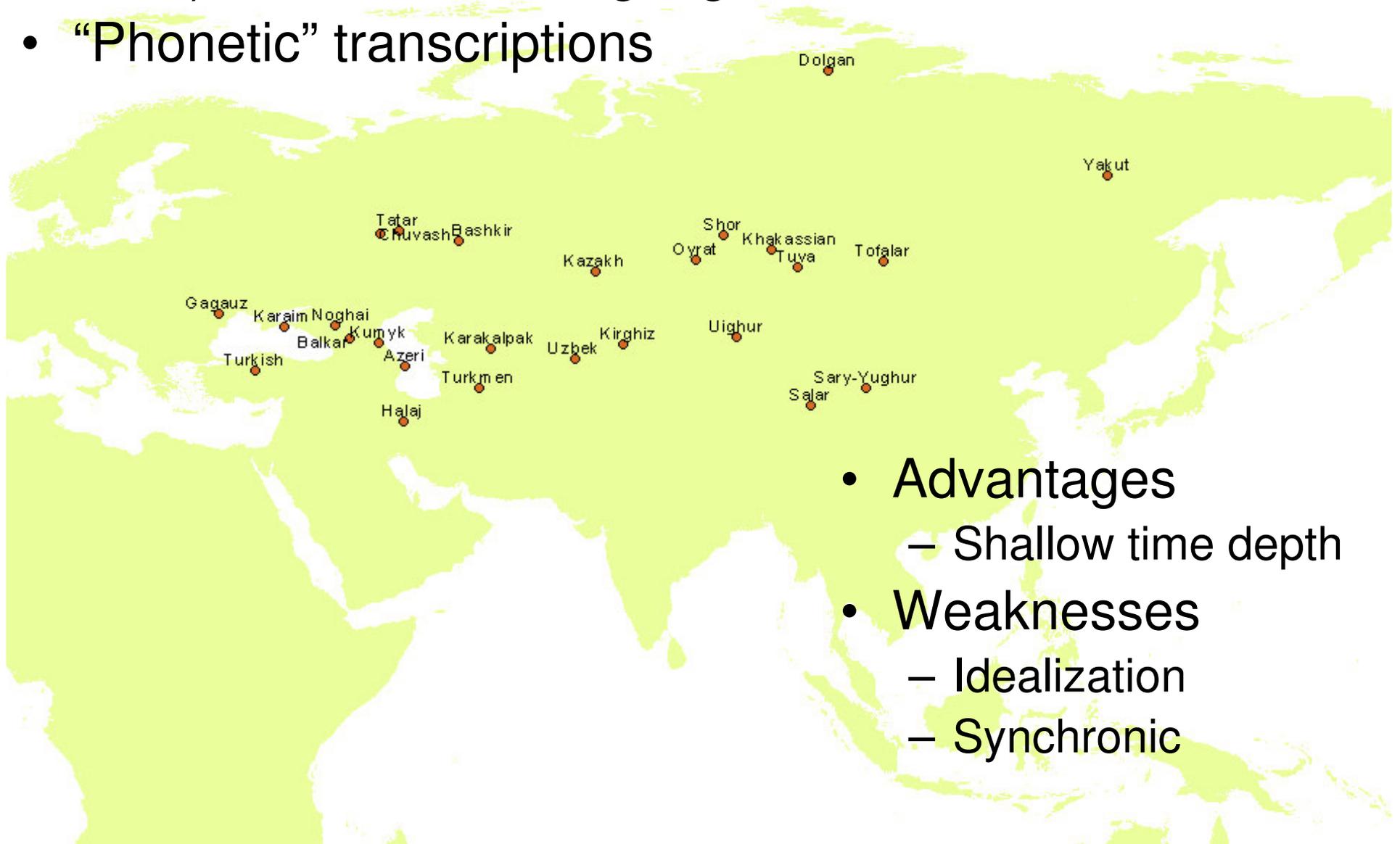


t	č	t
i	ə	i
r	r	r
n	n	
a	e	a
q		x

Find alignment that maximizes L

- Word lists (350 to 1400 words each) for 29 Turkic languages
- “Phonetic” transcriptions

# Data

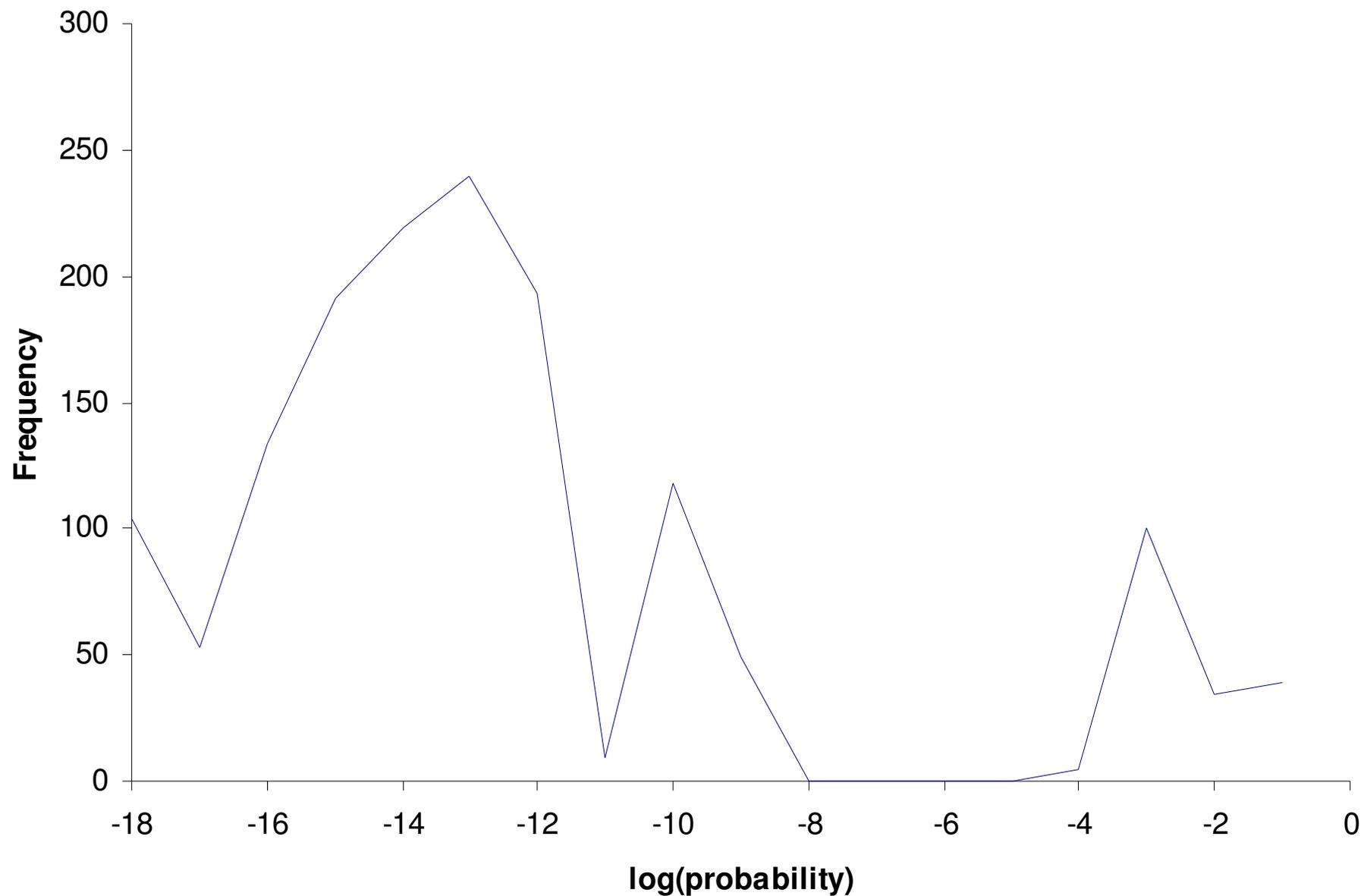


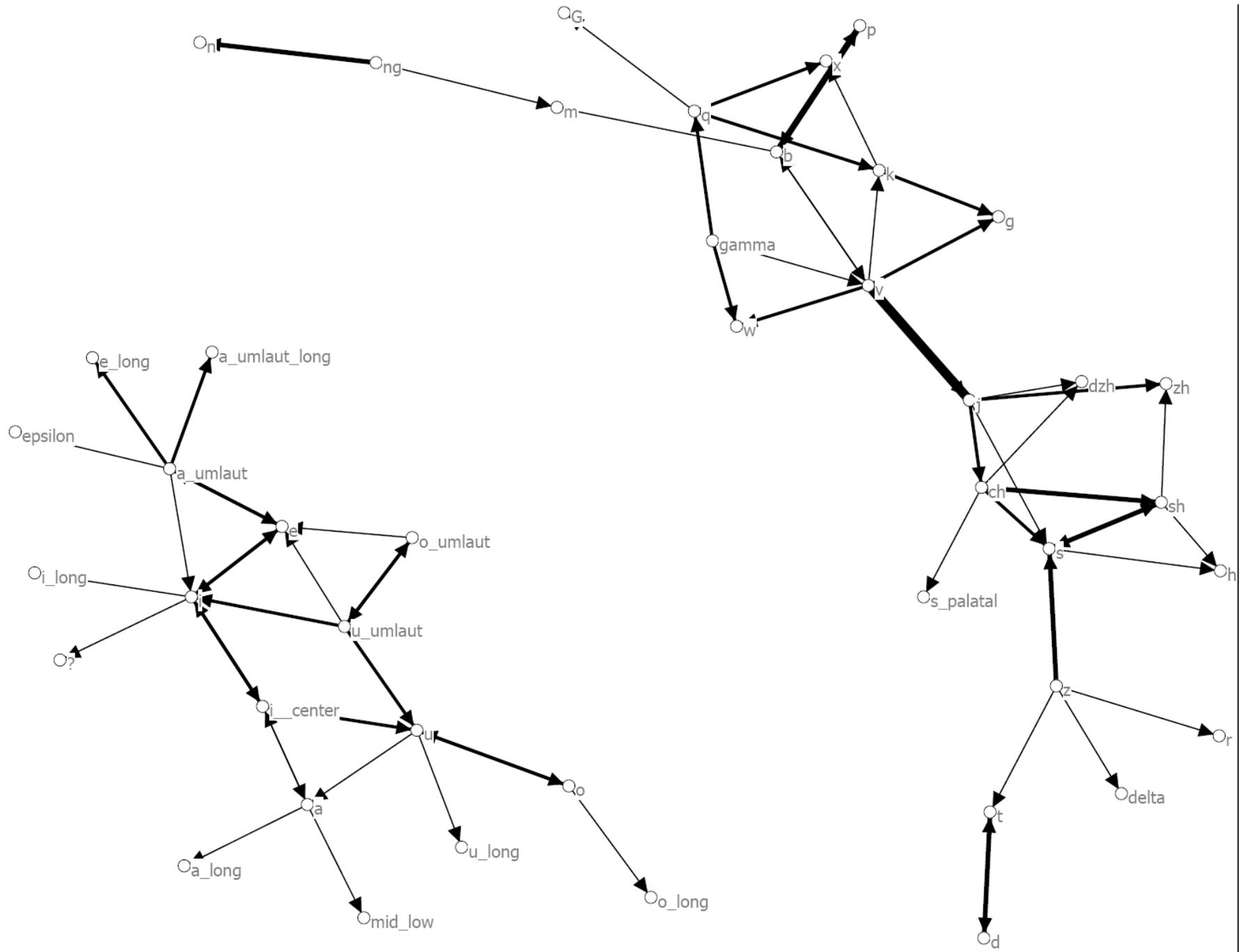
- Advantages
  - Shallow time depth
- Weaknesses
  - Idealization
  - Synchronic

# Finding ML estimates

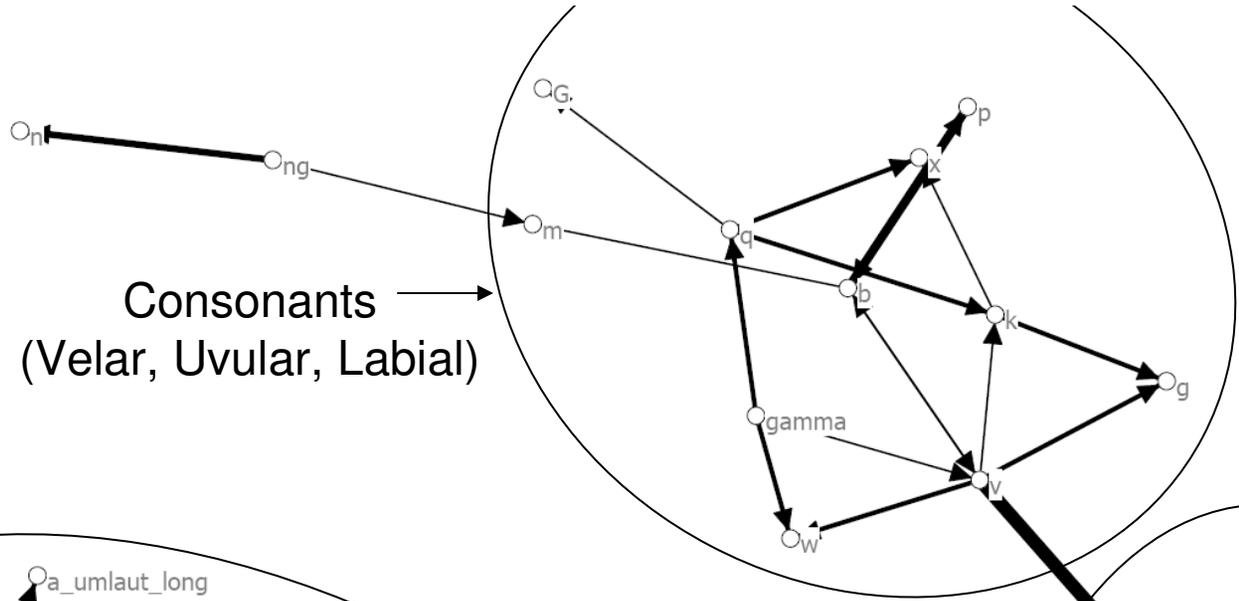
- 100 words
  - 9635 observations
  - 24 hours
  - loglikelihood = -8239
  - Ancestral sounds—15 consonants and 9 vowels with probability  $> 10^{-17}$
- 200 words
  - 16127 Observations
  - 118 hours
  - loglikelihood = -16903
  - Ancestral sounds— 18 consonants and 9 vowels with probability  $> 10^{-17}$

# Distribution of Probabilities in the Sound Change Matrix

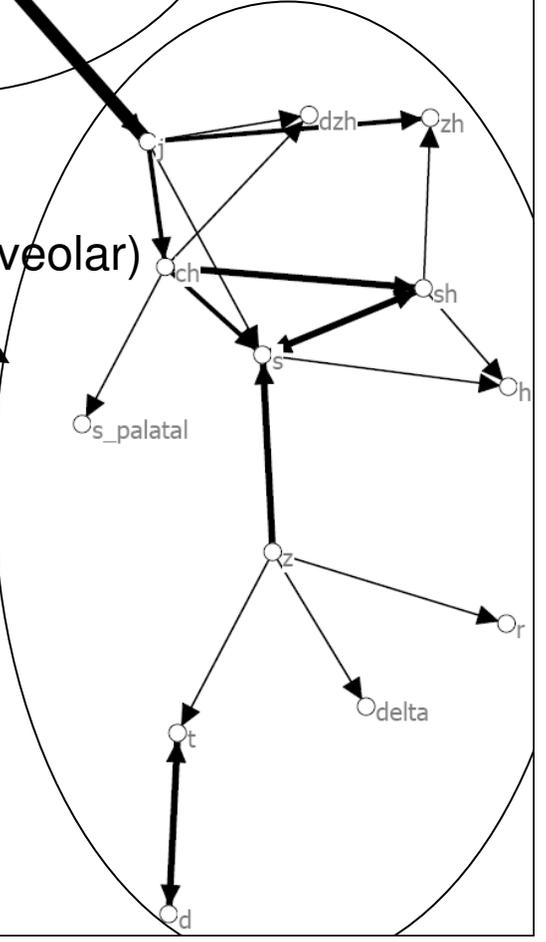




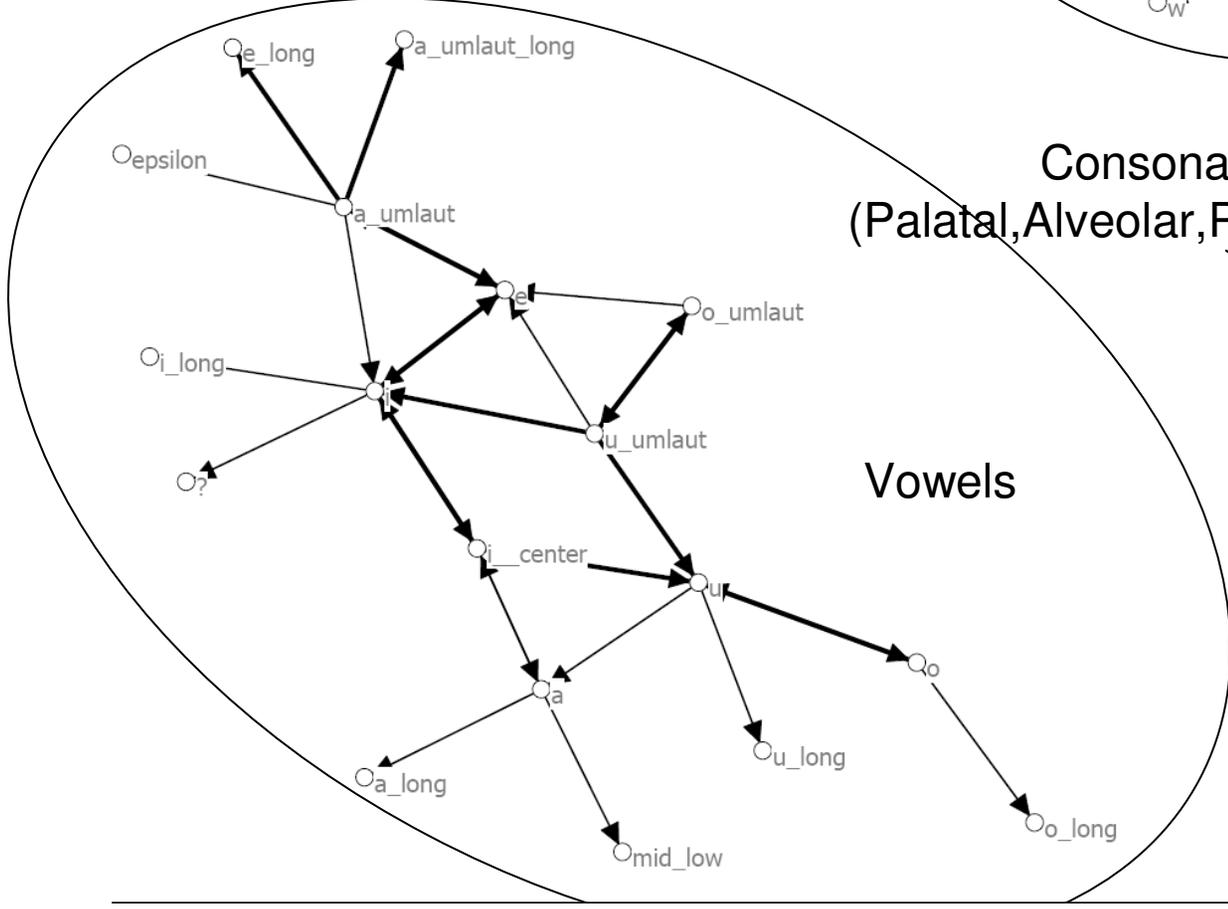
Consonants  
(Velar, Uvular, Labial)



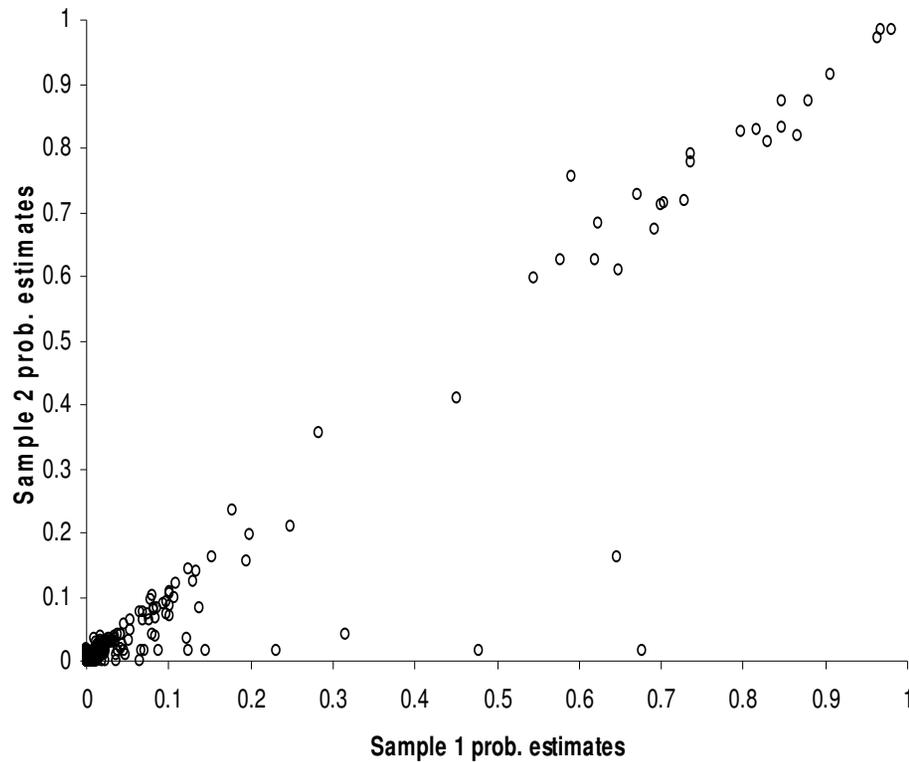
Consonants  
(Palatal, Alveolar, Postalveolar)



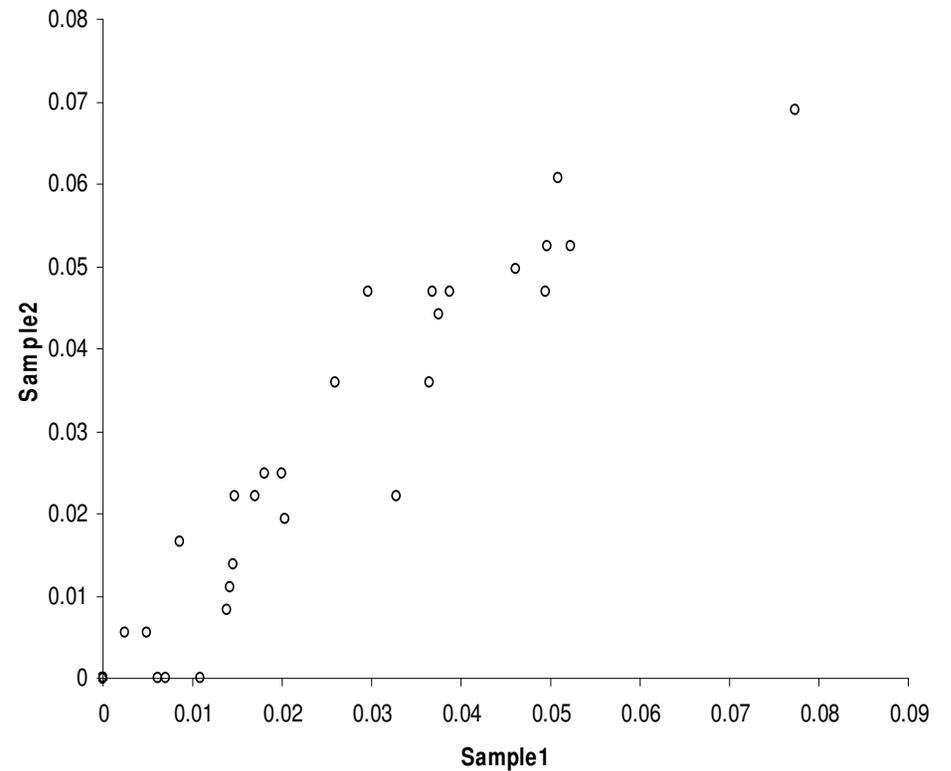
Vowels



# Comparing 2 samples' estimates



Estimated Sound Change Probabilities



Estimated Ancestral Probabilities

# Recap

- We have ML estimates for data given model
- They fit qualitative observations made by many linguists
- They are consistent for different samples

# Questions

- How to compare with other models?
  - Feature-based changes
  - Natural classes (Bouchard-Cote et al. 2007)
- We will always have misspecification. What to do?
  - Historical relationships
  - Regularity of sound change (Foot-Pied, Five-Pent)
  - Idiosyncratic, but common kinds of change (“A newt” & “an eft”)
- Is there a way to estimate faster?