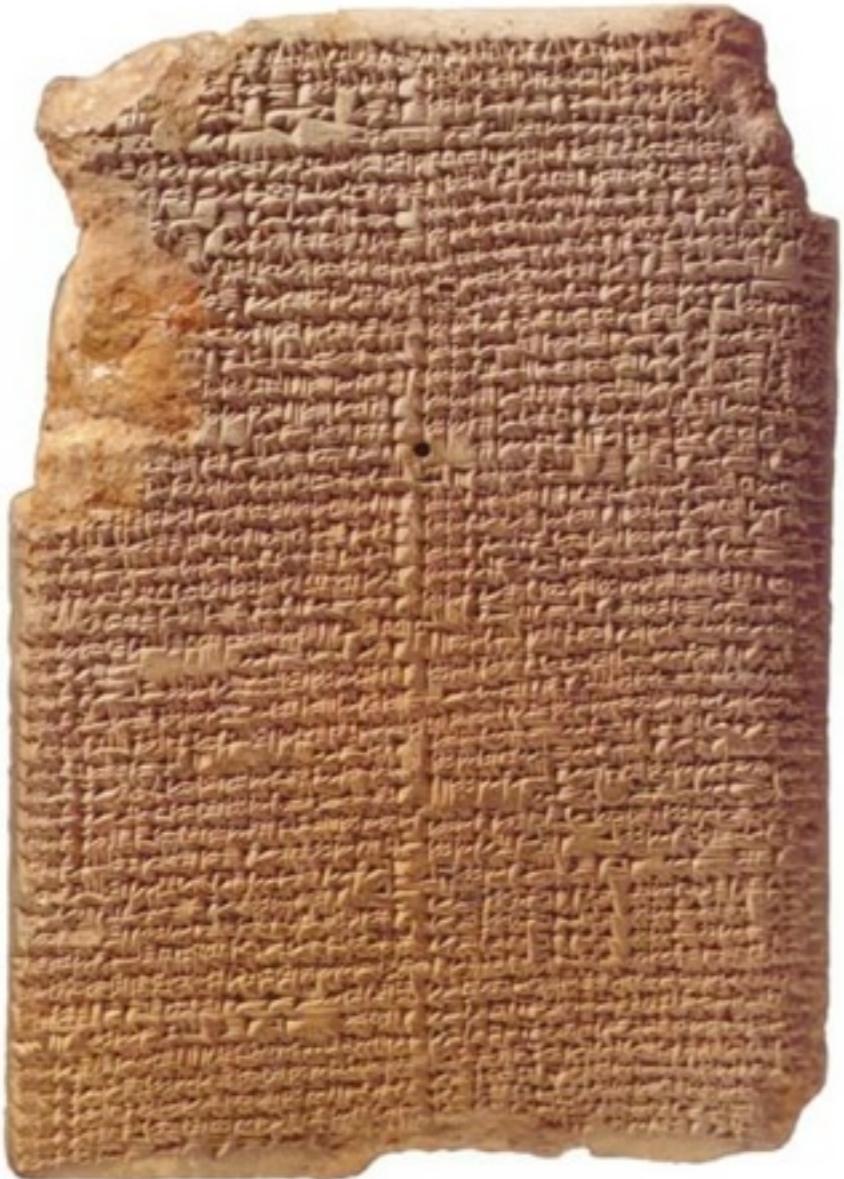


Unsupervised Consonant-Vowel Prediction over Hundreds of Languages

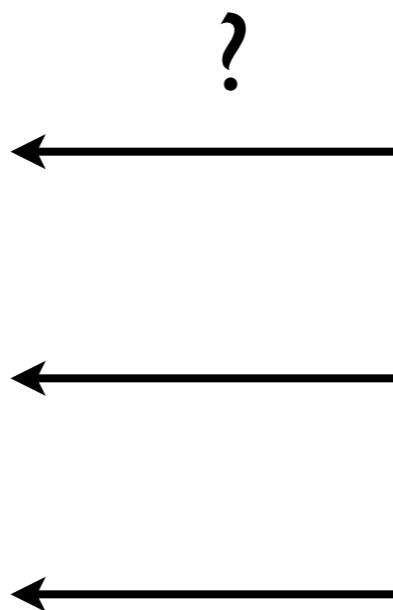


Young-Bum Kim and Benjamin Snyder
Wisconsin-Madison

lost language



known languages



[ACL 2010] Computational Decipherment of *Ugaritic*

⇒ Leverages *related* languages

Remaining Challenges

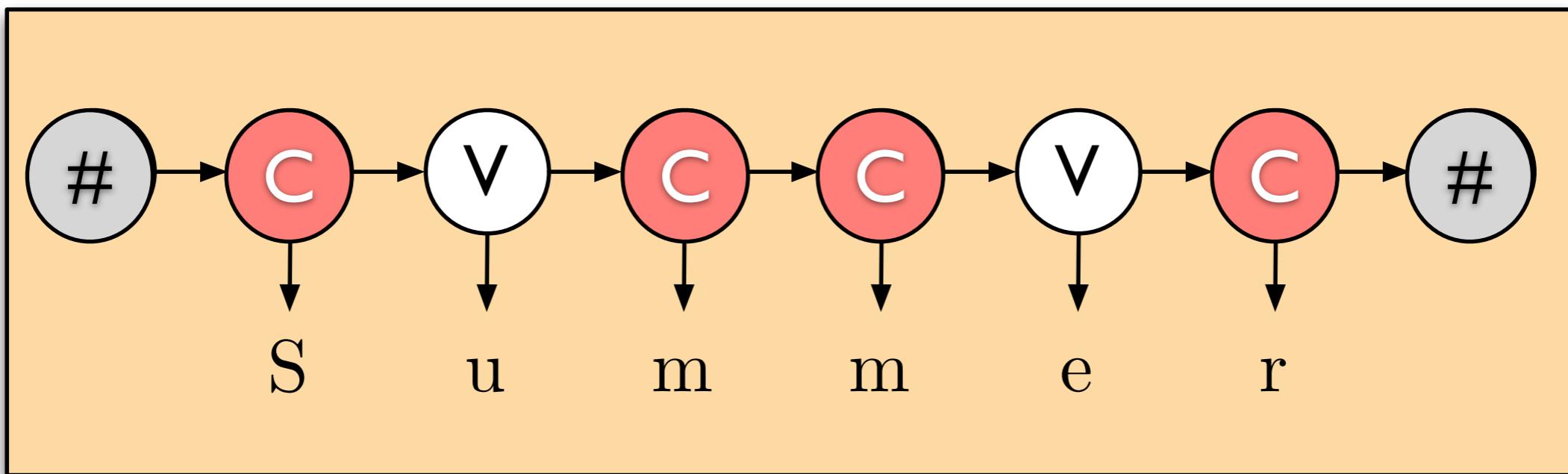
- Language family unknown
- No related living languages
 - ~ *Linear A*
 - ~ *Linear Elamite*
 - ~ *Byblos*
 - ~ *Cypro-Minoan*
 - ~ *Isthmian*
- Syllabic writing systems

Undeciphered:
Linear A



Knight et al 2006

- Trigram HMM over character sequences
 - ~ Two hidden states
 - ~ Estimated using EM
- Distinguishes Spanish consonants/vowels



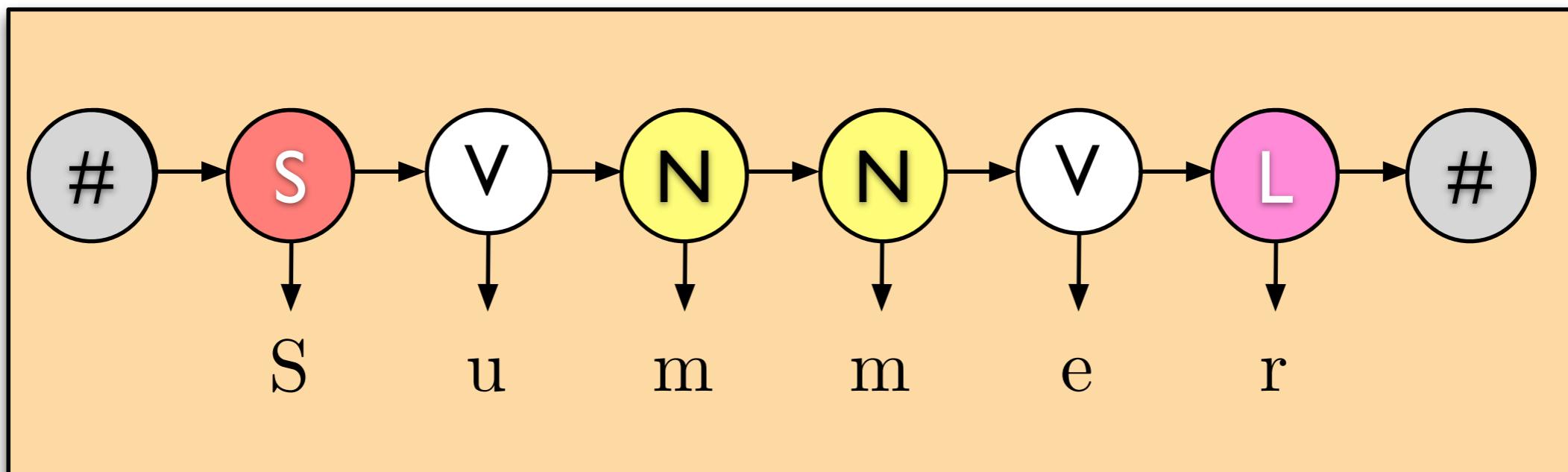
Lingering Questions

1. Just Spanish?

- Georgian: prtskvna (*peeling*) : CCCCCCCCV

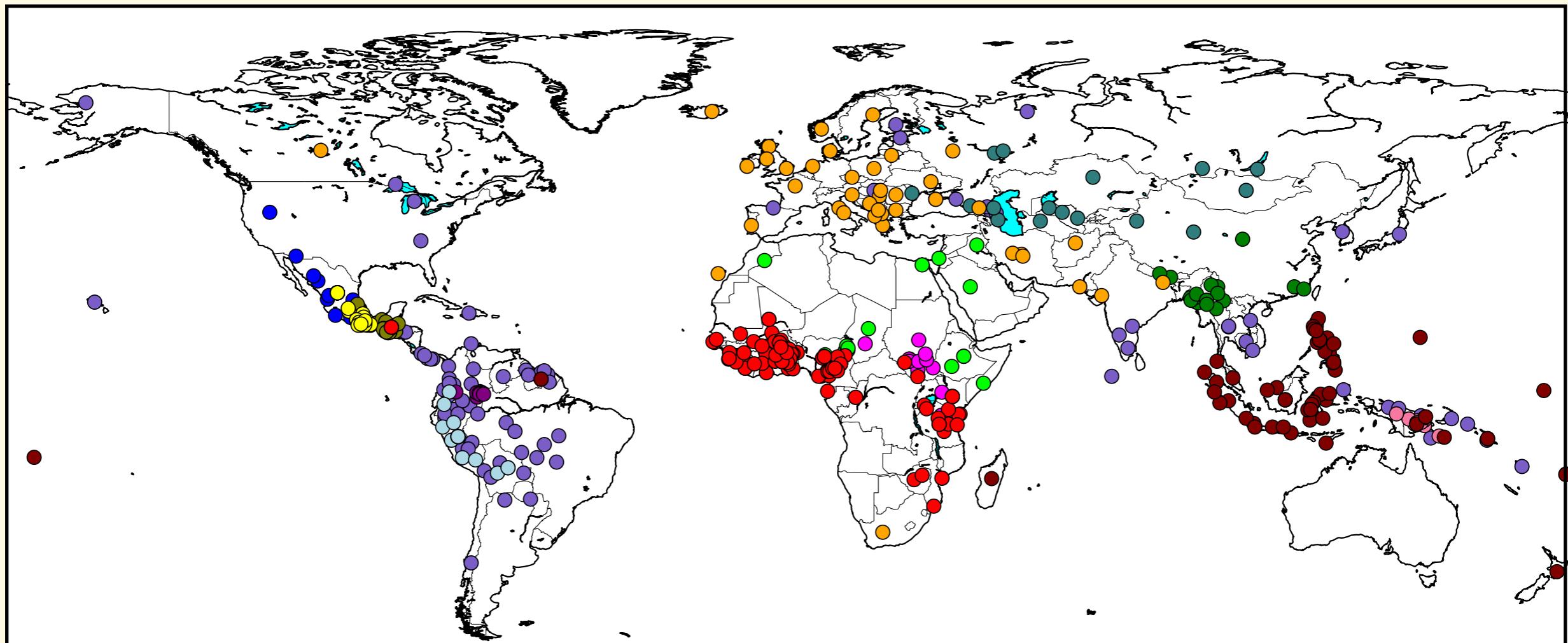
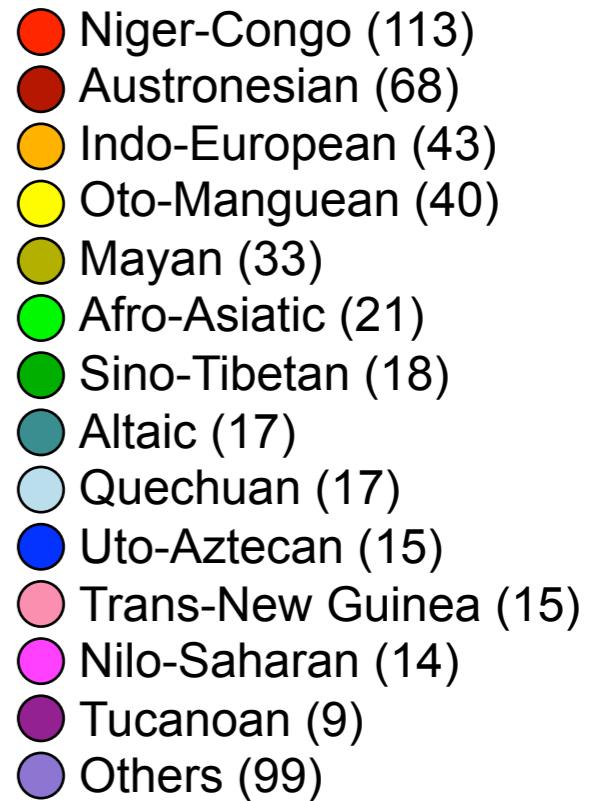
2. Leverage *known* languages?

3. sibilant, nasal, liquid, etc?



Bible Data for 503 alphabetic languages

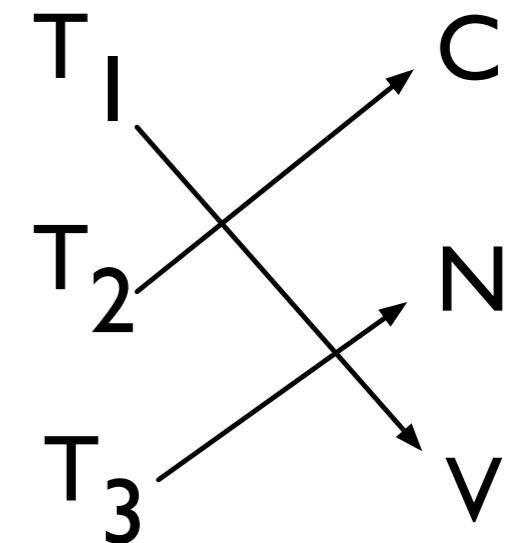
1. Cons vs Vowel
2. Nasal vs Cons vs Vowel



HMM Results

C / V	N / C / V
93.4%	74.6%

- 503 alphabetic languages
- 10 EM restarts
- Oracle tag mappings:



⇒ [C / V] dominant natural pattern

⇒ [N / C / V] less so...

Model Design Goals

1. Cluster letters:

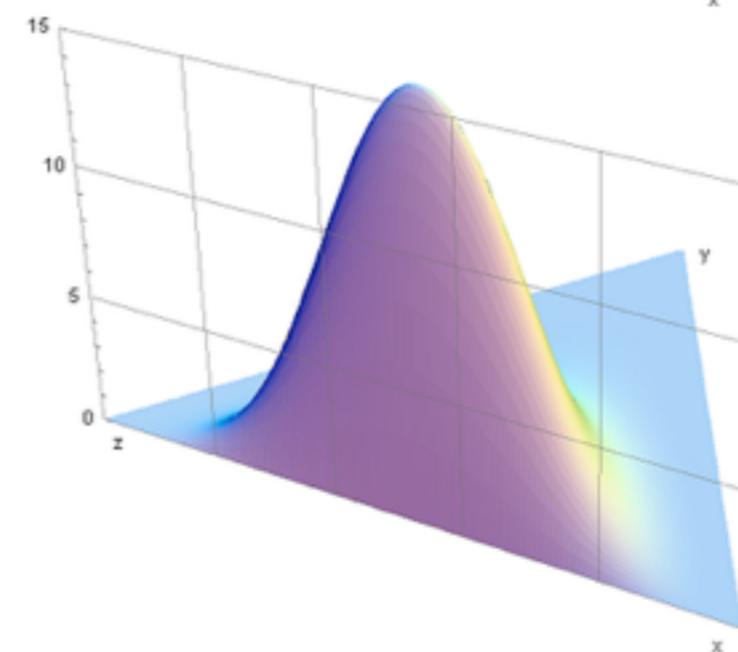
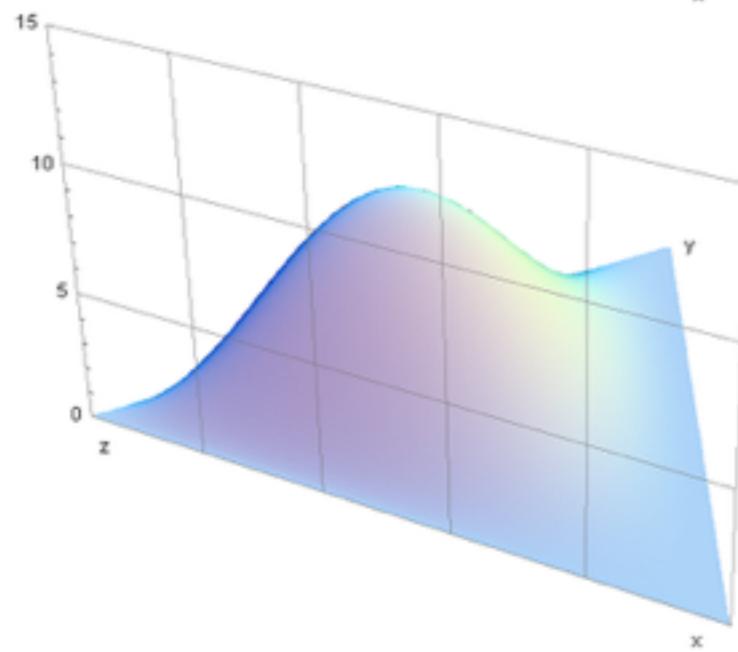
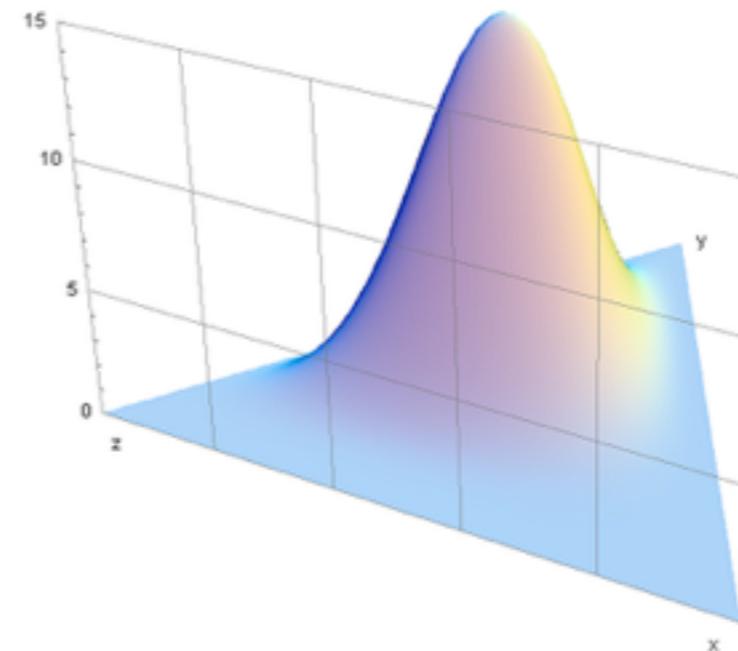
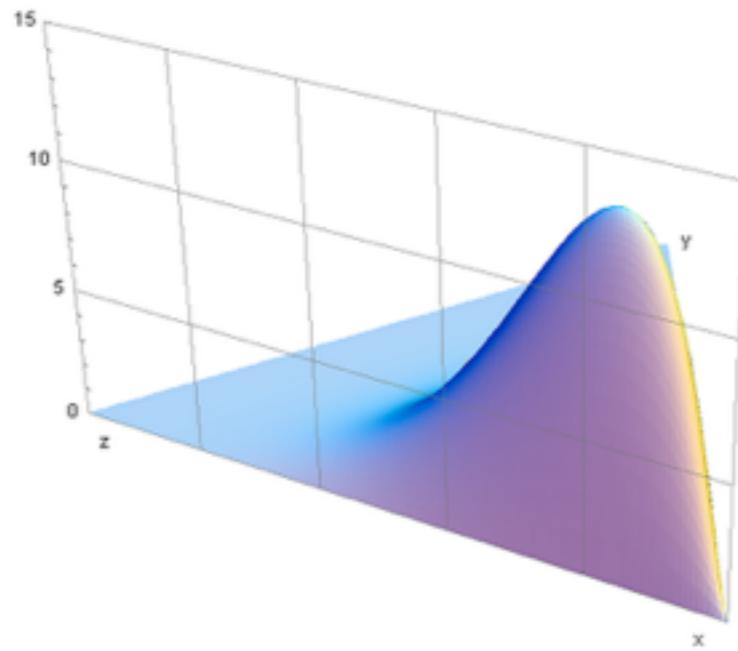
- HMM does *soft* clustering
- letters typically in single phonetic class

2. Learn from hundreds of languages

- Avoid oracle tag mappings

3. Model typologically coherent groupings

Bayesian HMM: Dirichlet Priors



$$\alpha_1, \alpha_2, \alpha_3 > 0$$

Hard Bayesian HMM

- One tag per character type
(stochastic *deterministic* FSA)
- *Dimensionality* of tag emissions unknown:

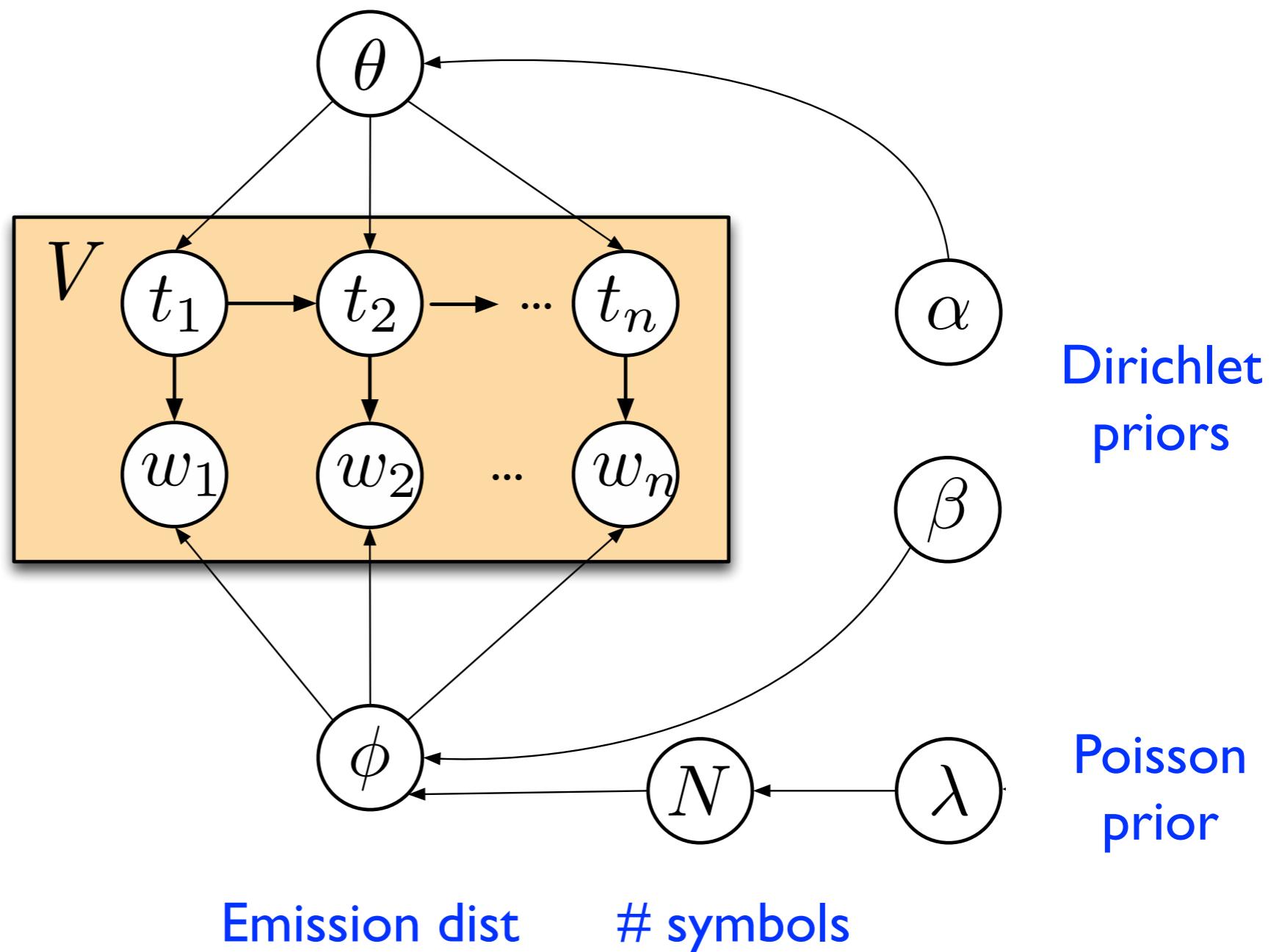
Emission distribution $\phi|N \sim \text{Dirichlet}(\beta_1 \dots \beta_N)$

symbols $N|\lambda \sim \text{Poisson}(\lambda)$

Model I: SYMM

Assume fixed, symmetric priors

transition dist



Inference: Gibbs Sampling

- Integrate out parameters
- Iteratively sample tag assignment t_w of symbol types W :

$$f(t_w | \mathbf{t}_{-w}, w_1 \dots w_n) \propto$$

Inference: Gibbs Sampling

- Integrate out parameters
- Iteratively sample tag assignment t_w of symbol types W :

$$f(t_w | \mathbf{t}_{-w}, w_1 \dots w_n) \propto \int f(N|\lambda) f(\lambda) d\lambda \quad (\text{type counts}) \quad (1)$$

$$\int f(t_1 \dots t_n | \theta) f(\theta) d\theta \quad (\text{tags}) \quad (2)$$

$$\int f(w_1 \dots w_n | t_1 \dots t_n, \phi) f(\phi|N) d\phi \quad (\text{symbols}) \quad (3)$$

Tag Predictive Distribution

$$\frac{\prod_{t,t'} (\alpha_{t,t'} + n(t,t'))^{[\delta(t,t')]} \text{tag bigram counts excluding symbol } W}{\prod_t \dots \dots \dots}$$

Dirichlet hyperparameter

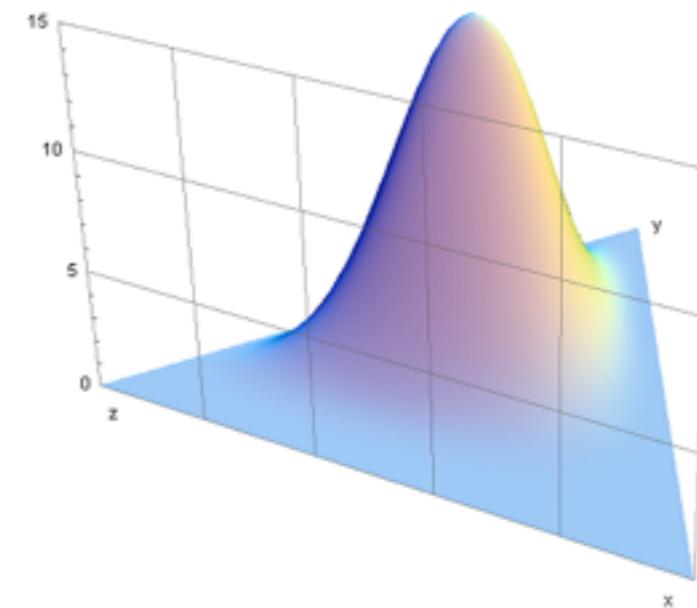
tag bigram counts excluding symbol W

tag bigram counts only with symbol W

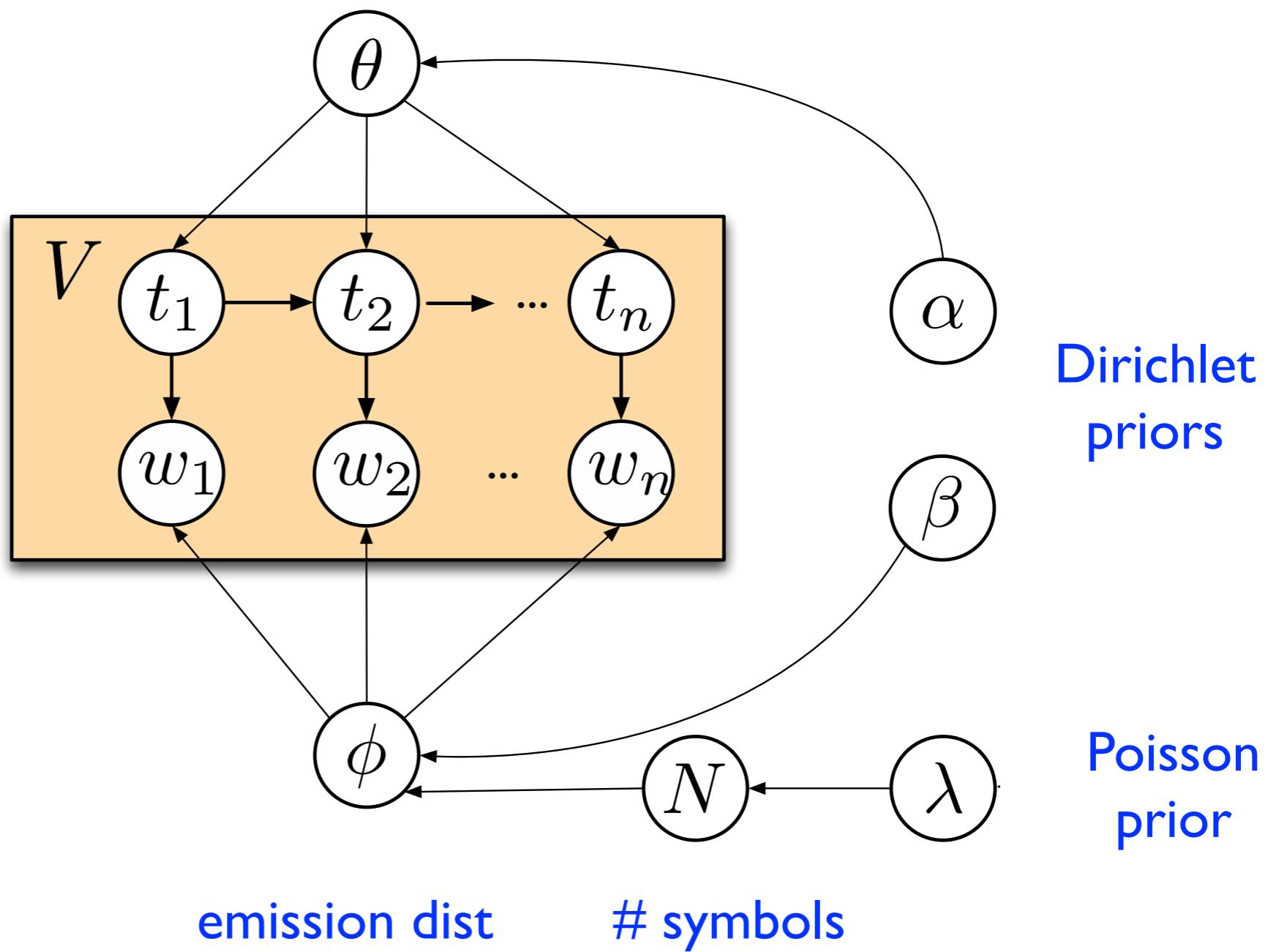
Next Idea

- Trying to decipher an unknown language
- Learn from our 502 other languages where the tags are *known*

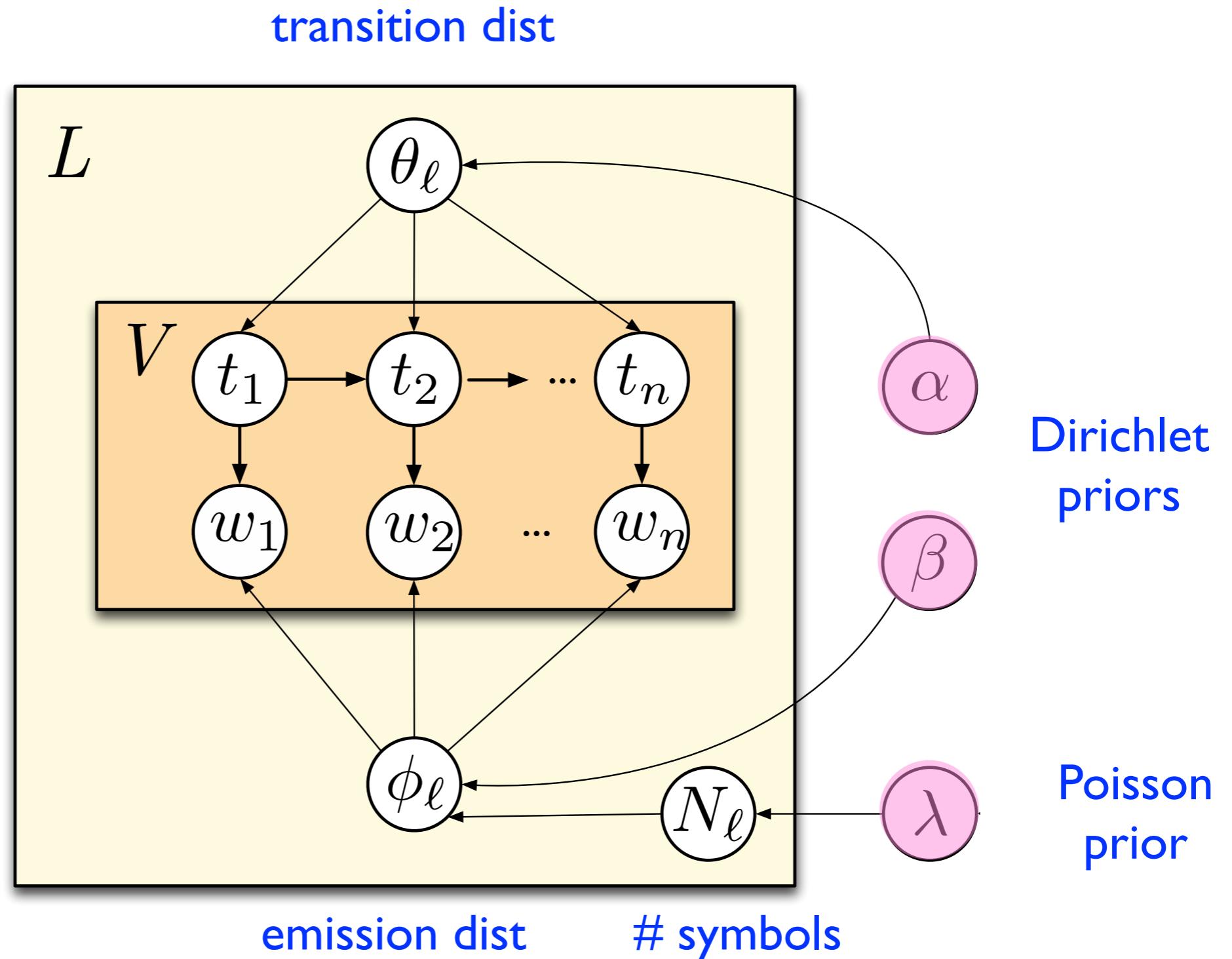
⇒ Assume shared,
unknown priors
across languages



transition dist



Model 2: Infer universal language priors MERGED



Inference: Gibbs Sampling

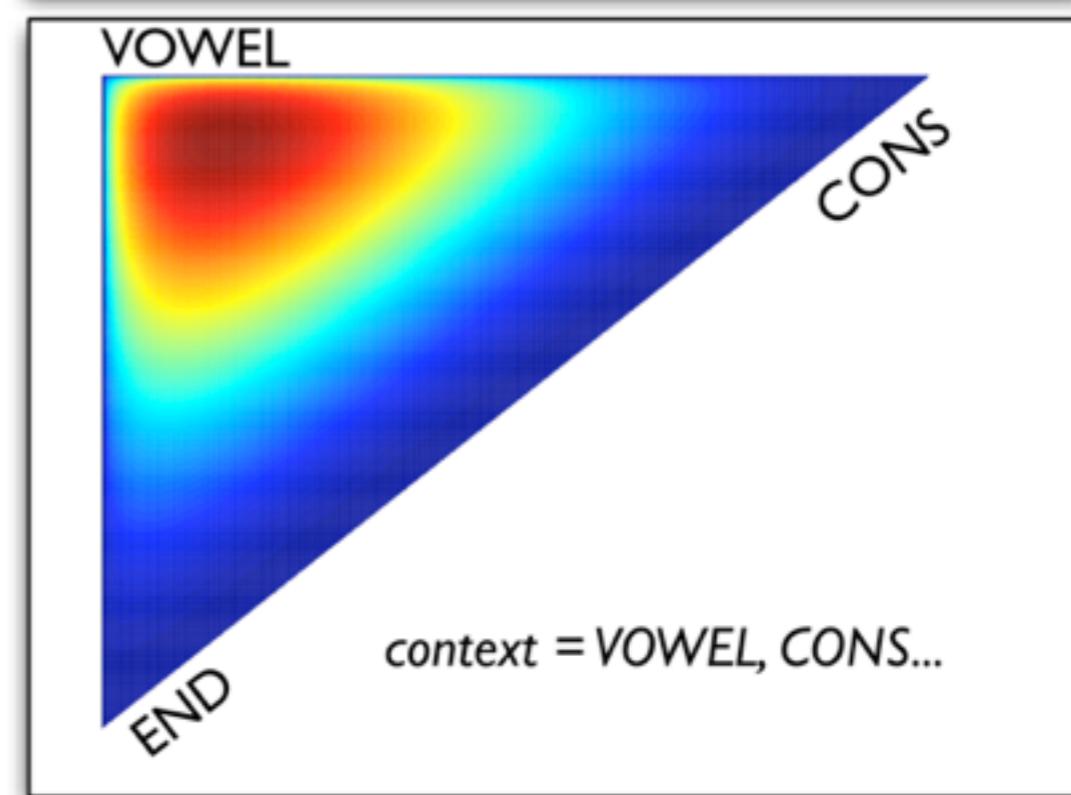
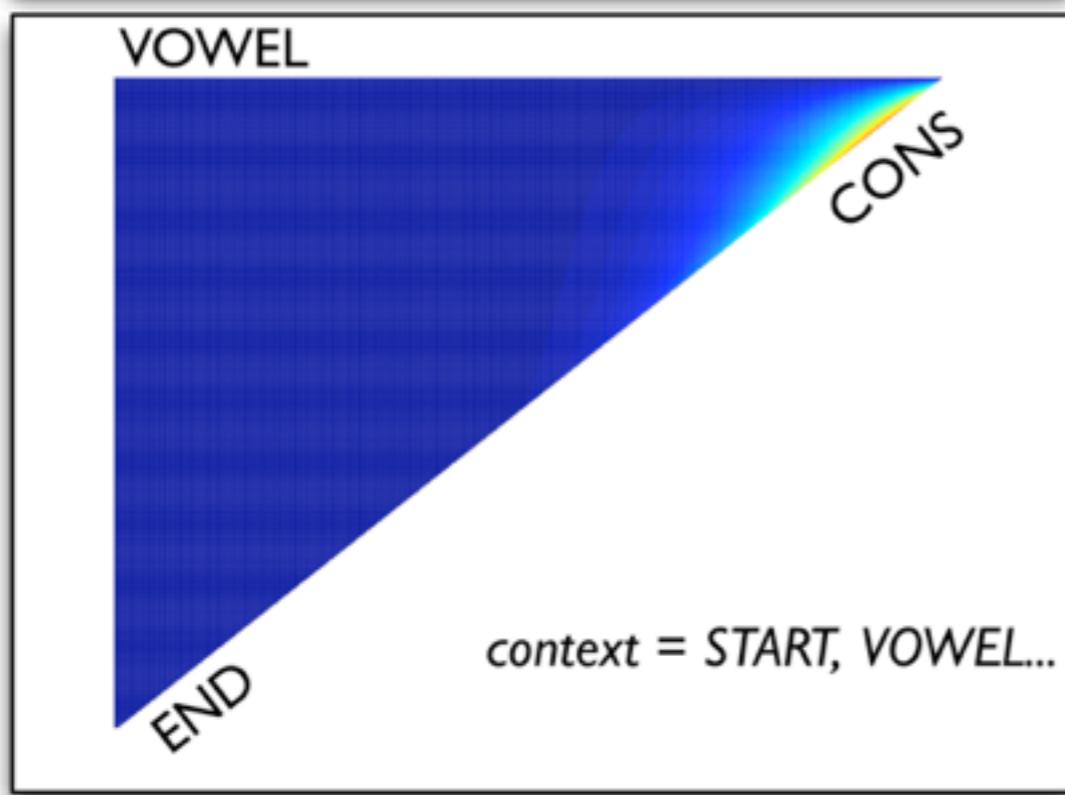
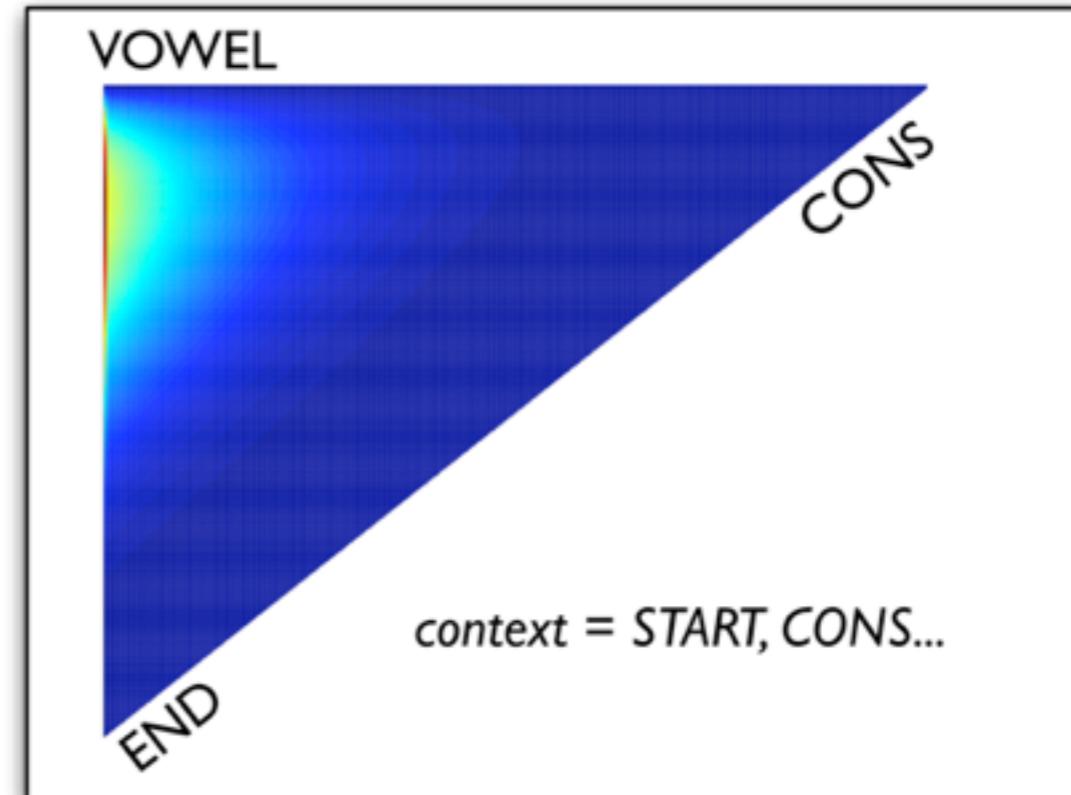
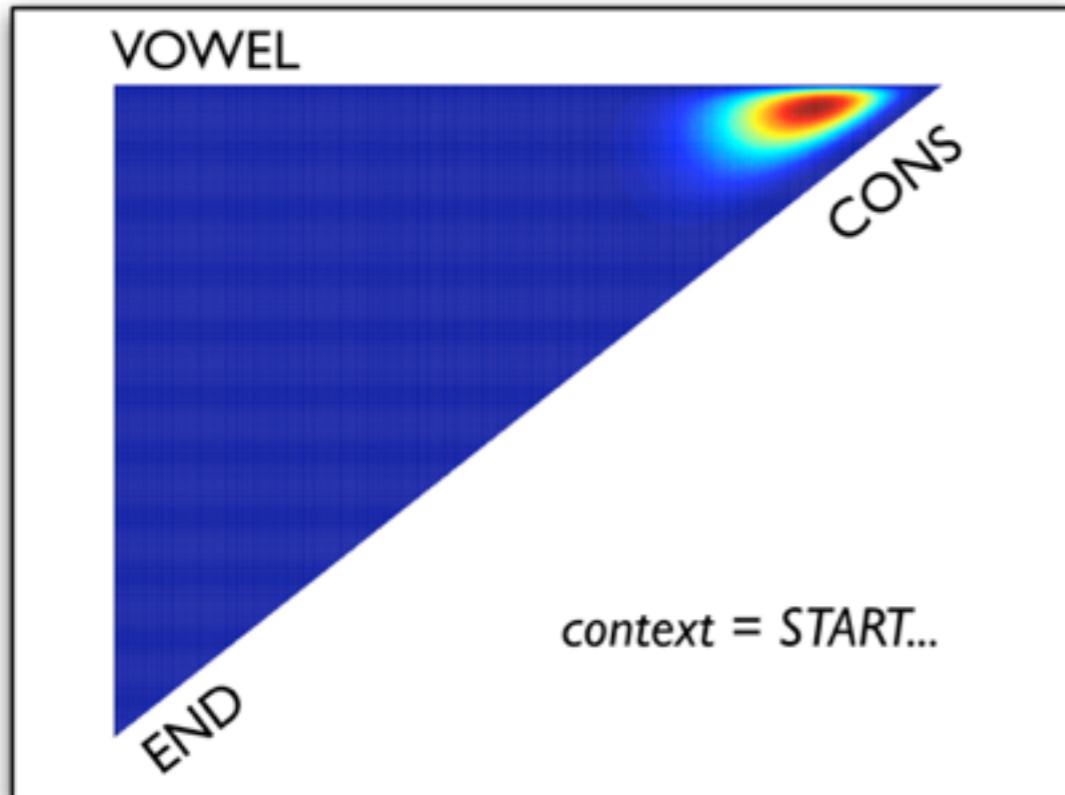
- Slice sample hyperparameters
 - ~ *product of predictive distributions with common prior*
 - ~ **combine terms across languages for efficiency**

$$f(\alpha_{t,t'} | \dots) \propto \prod_k \frac{(\alpha_{t,t'} + k)^{\#(k,t,t')}}{\dots \dots}$$

number of languages which have $\geq k$ occurrences of bigram (t, t')

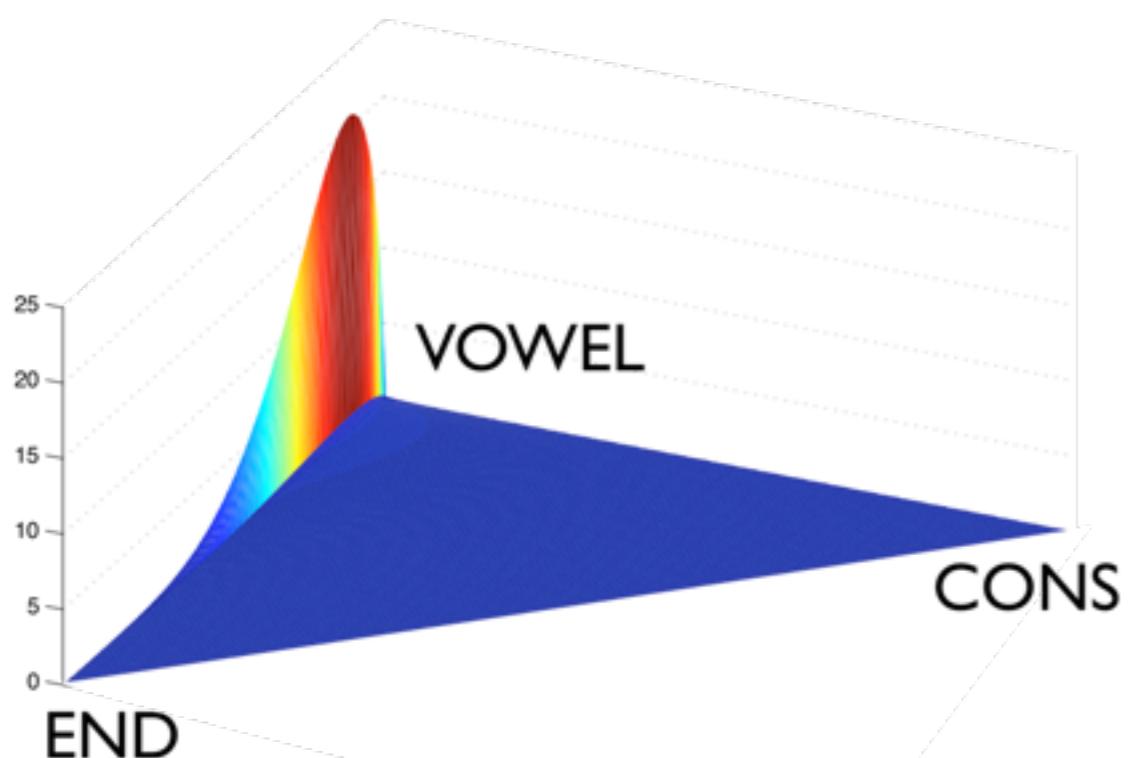


MAP Transition Priors: 503 languages

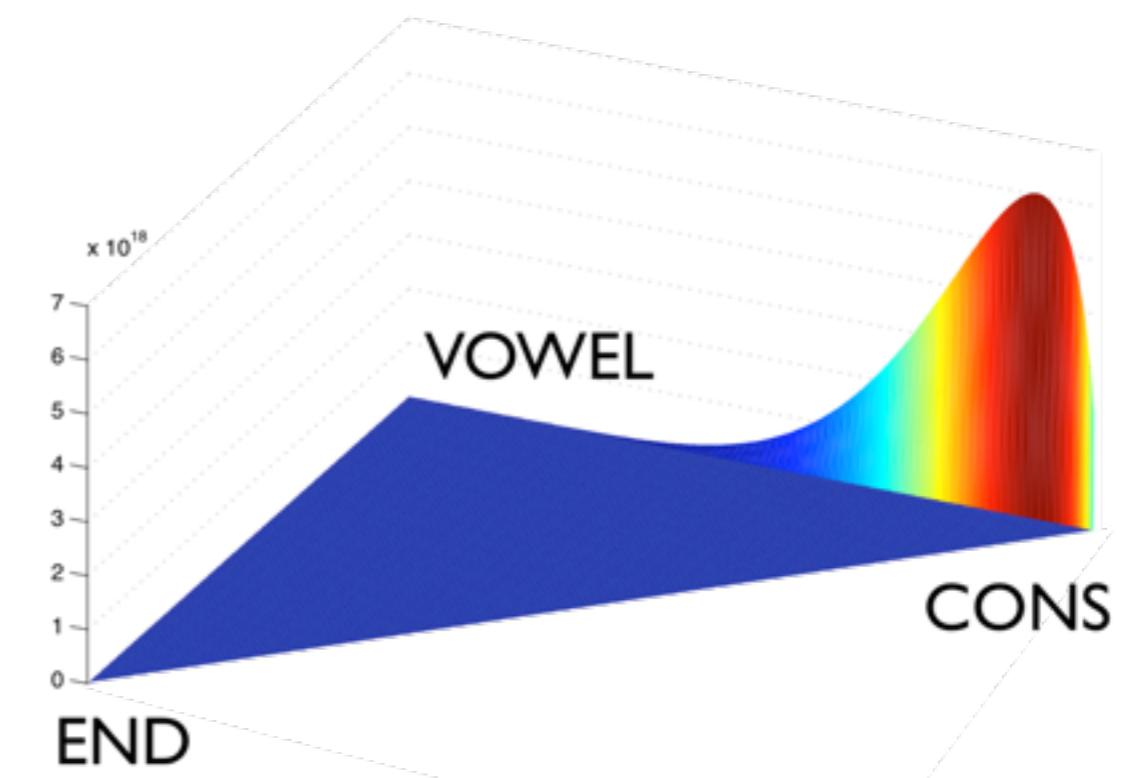


MAP Transition Priors: 503 languages

context = CONS, CONS...



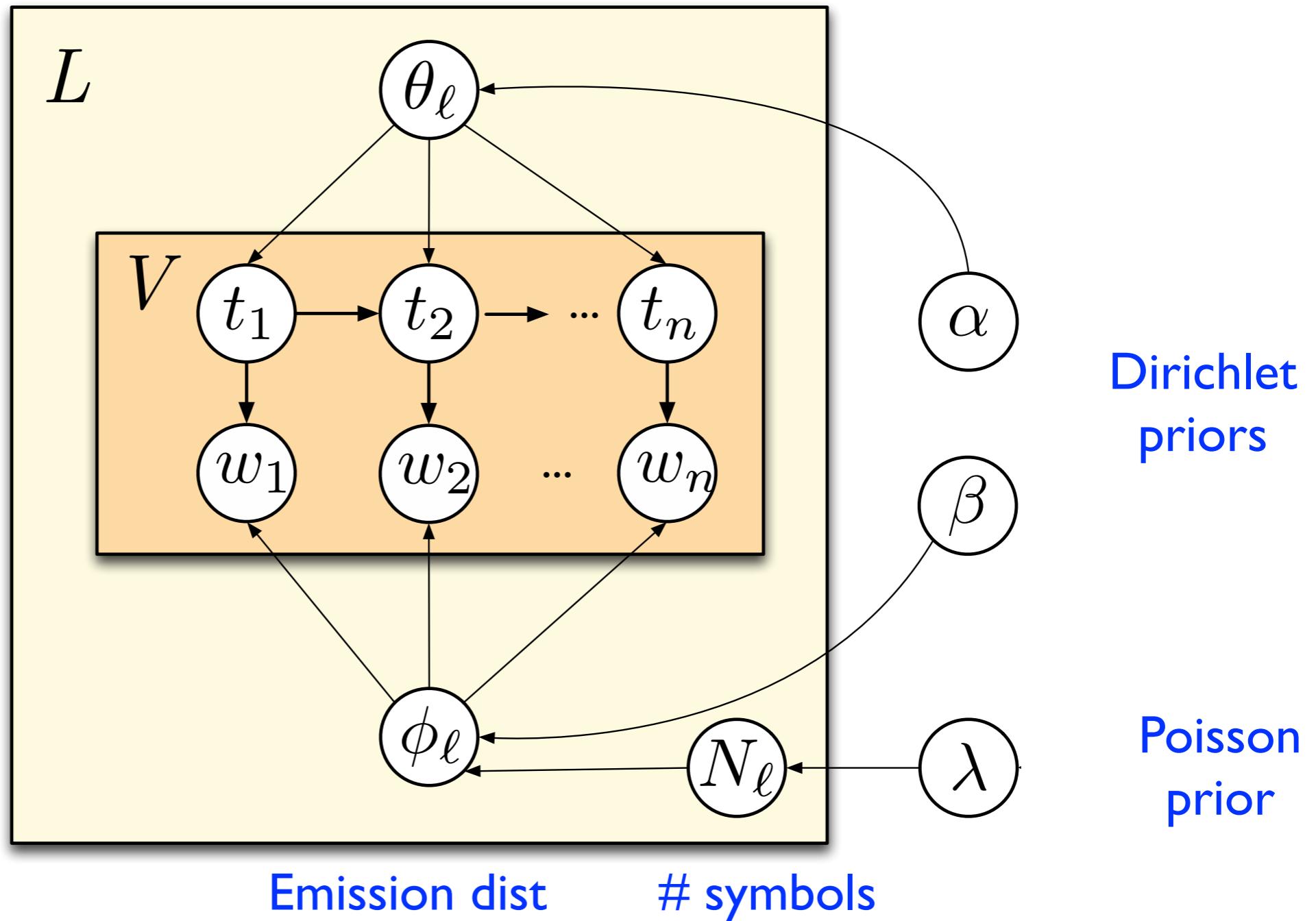
context = CONS, VOWEL...



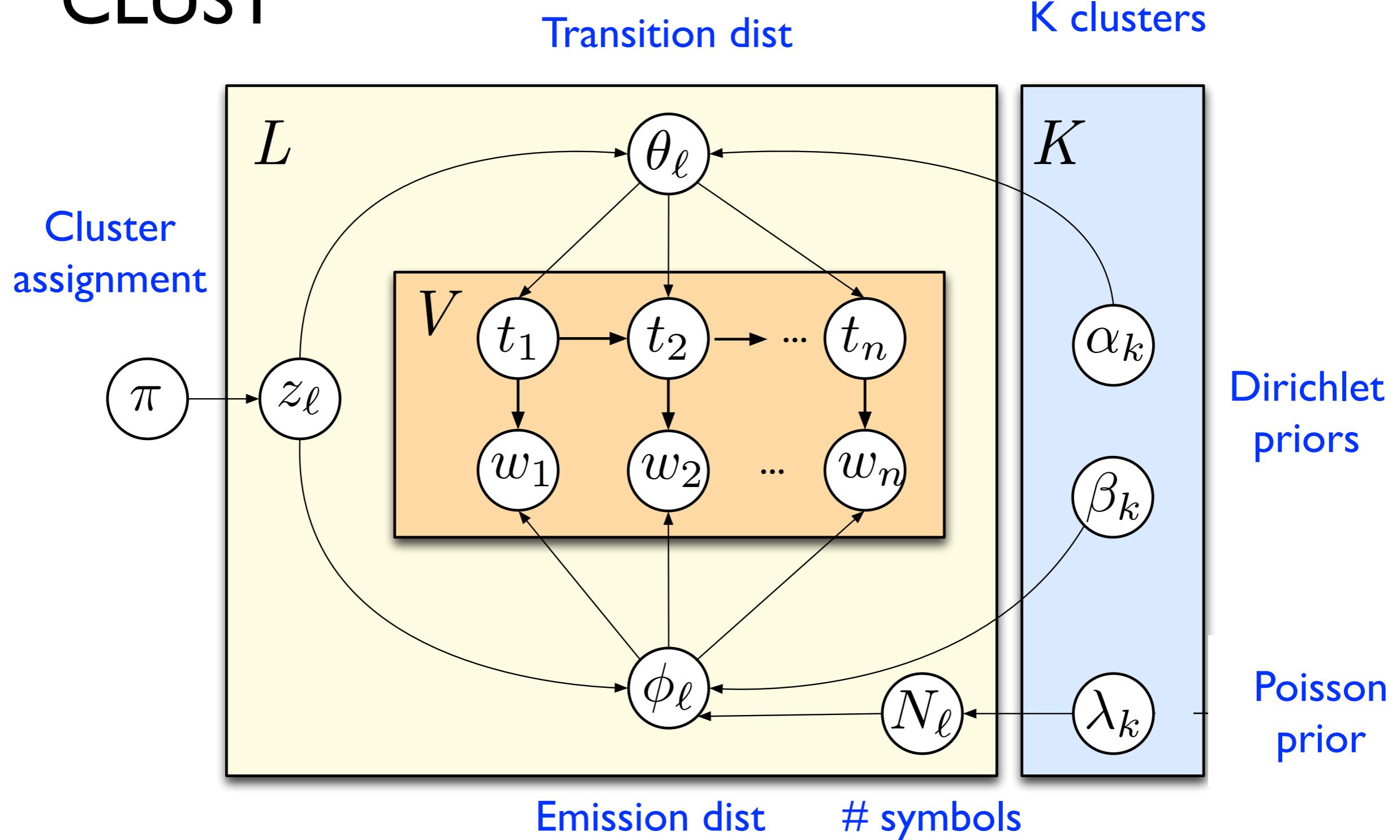
Final Idea

- Languages vary typologically
 - ⇒ Assume latent clustering of languages
 - ⇒ Hyperparameters shared at the cluster level

Transition dist



Model 3: CLUST



Evaluation

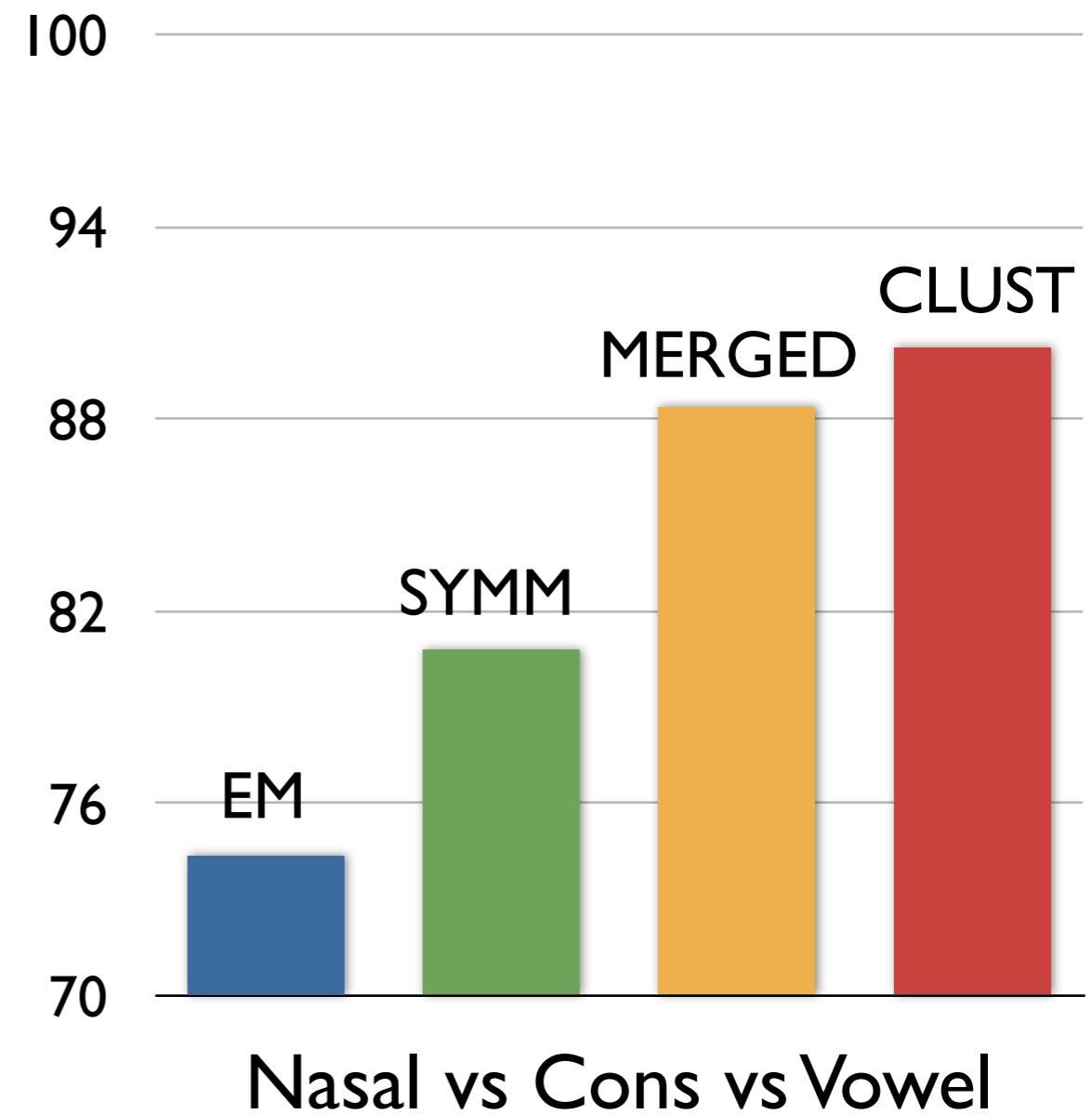
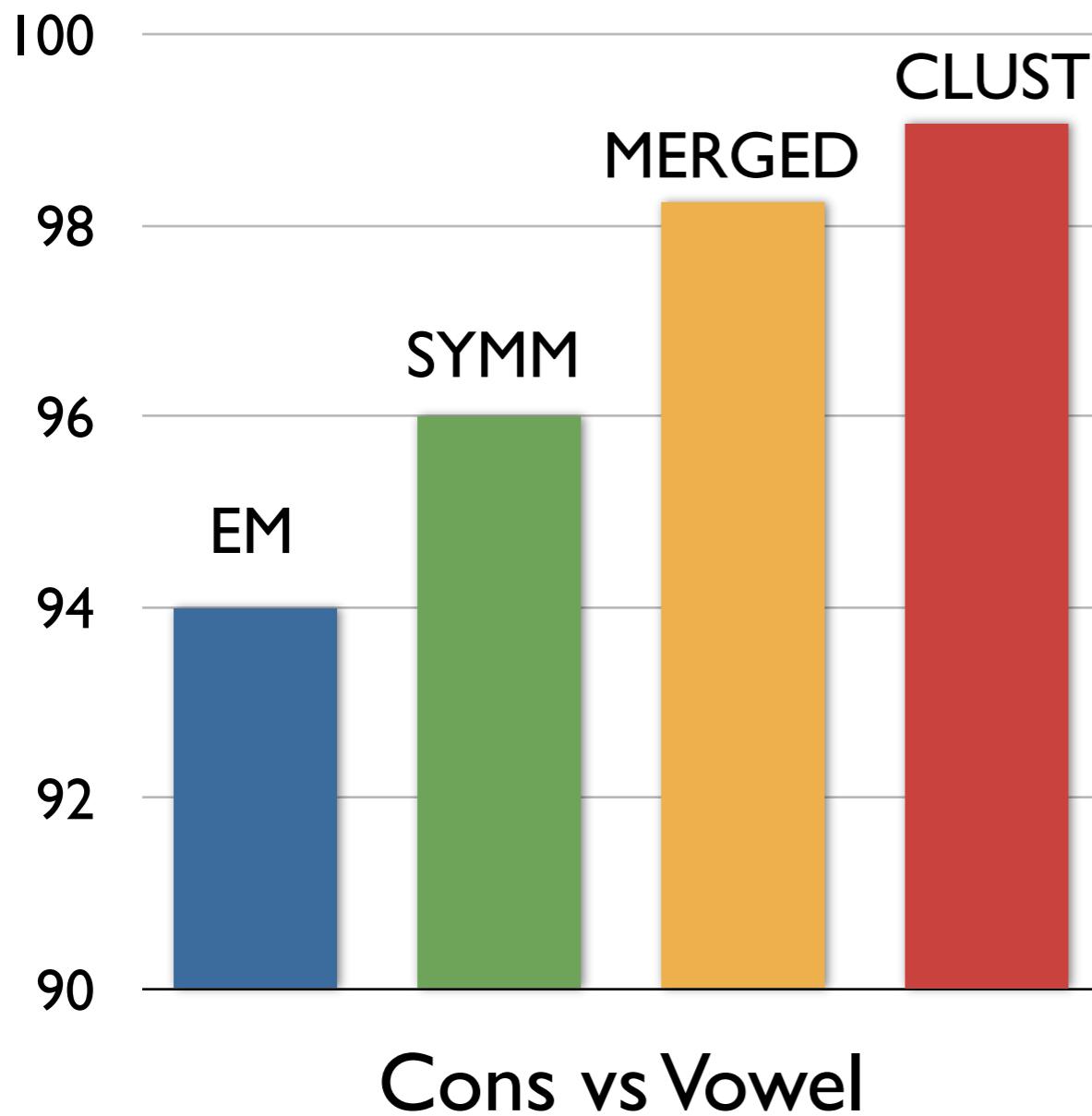
- Token-level accuracy
 - ~ 503 languages
 - ~ leave-one-out average
 - Baseline:
 - ~ HMM (EM) [Knight 2006] ← oracle mappings
 - Hard Bayesian HMM
 - ~ SYMM: symmetric fixed priors ← oracle mappings
 - ~ MERGE: uses 502 observed languages
 - ~ CLUST: 20 latent language clusters
- Two Tasks

Cons / Vowel

Nasal / Cons / Vowel
- 

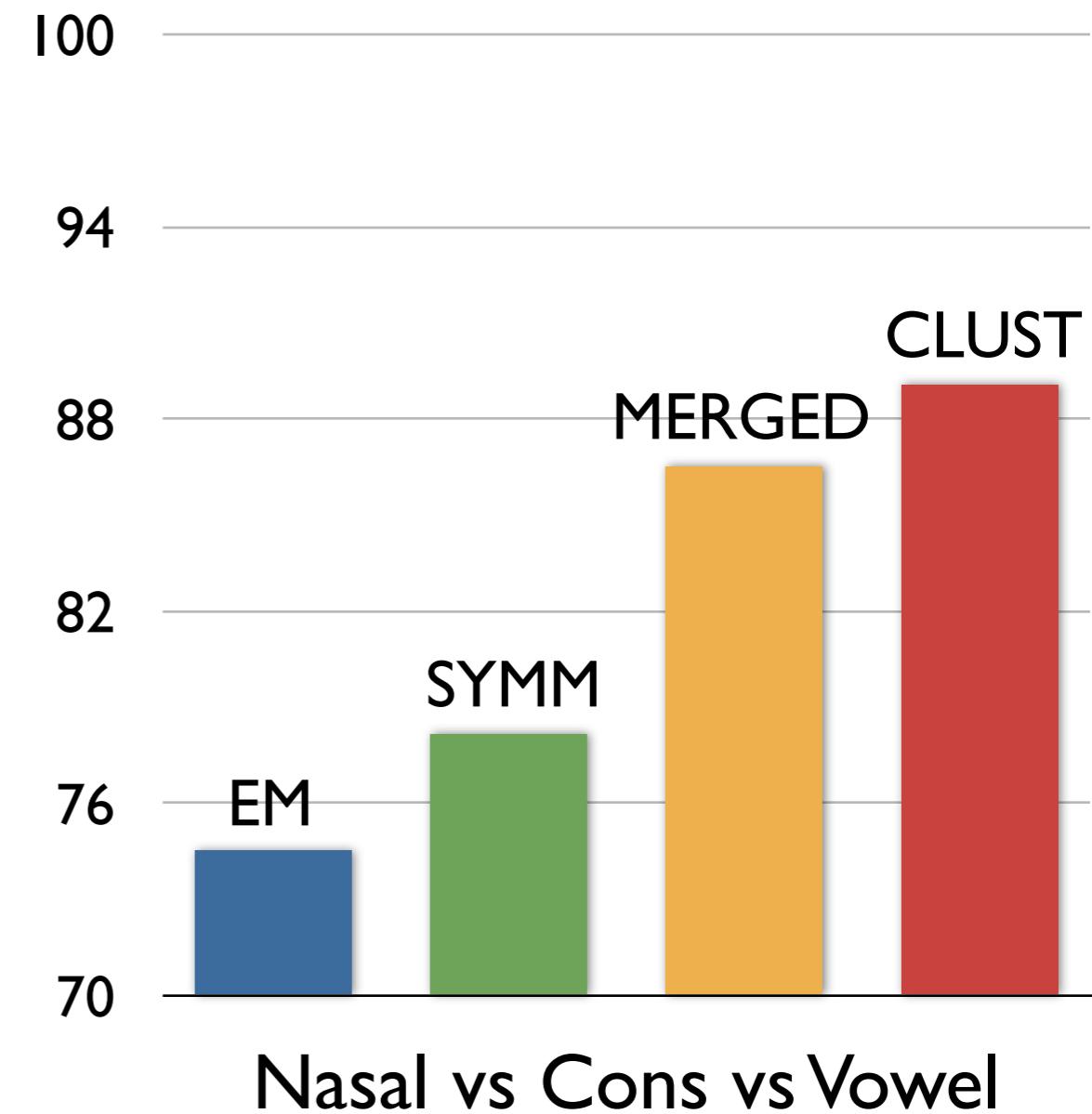
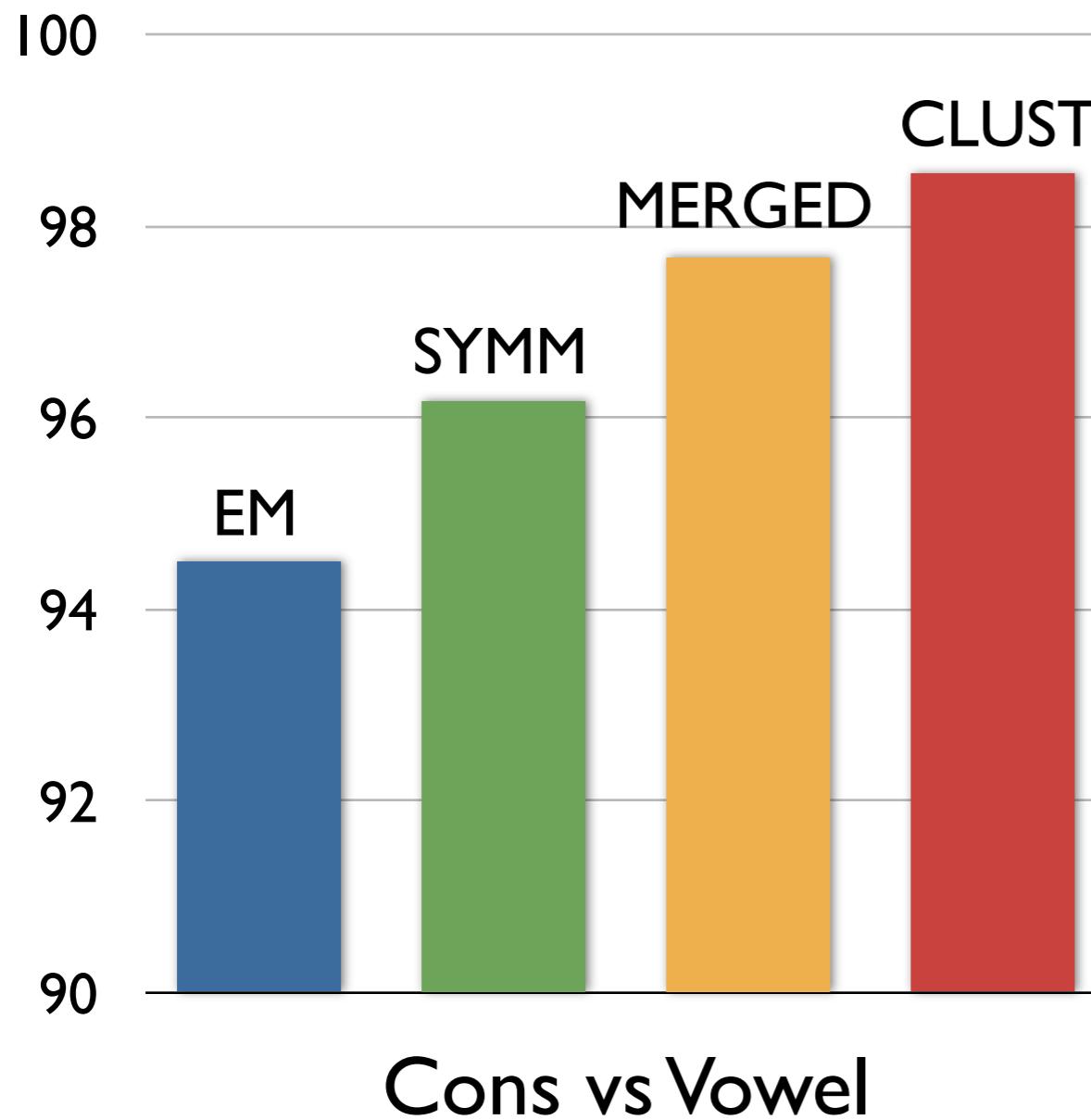
Results

- 503 alphabetic languages

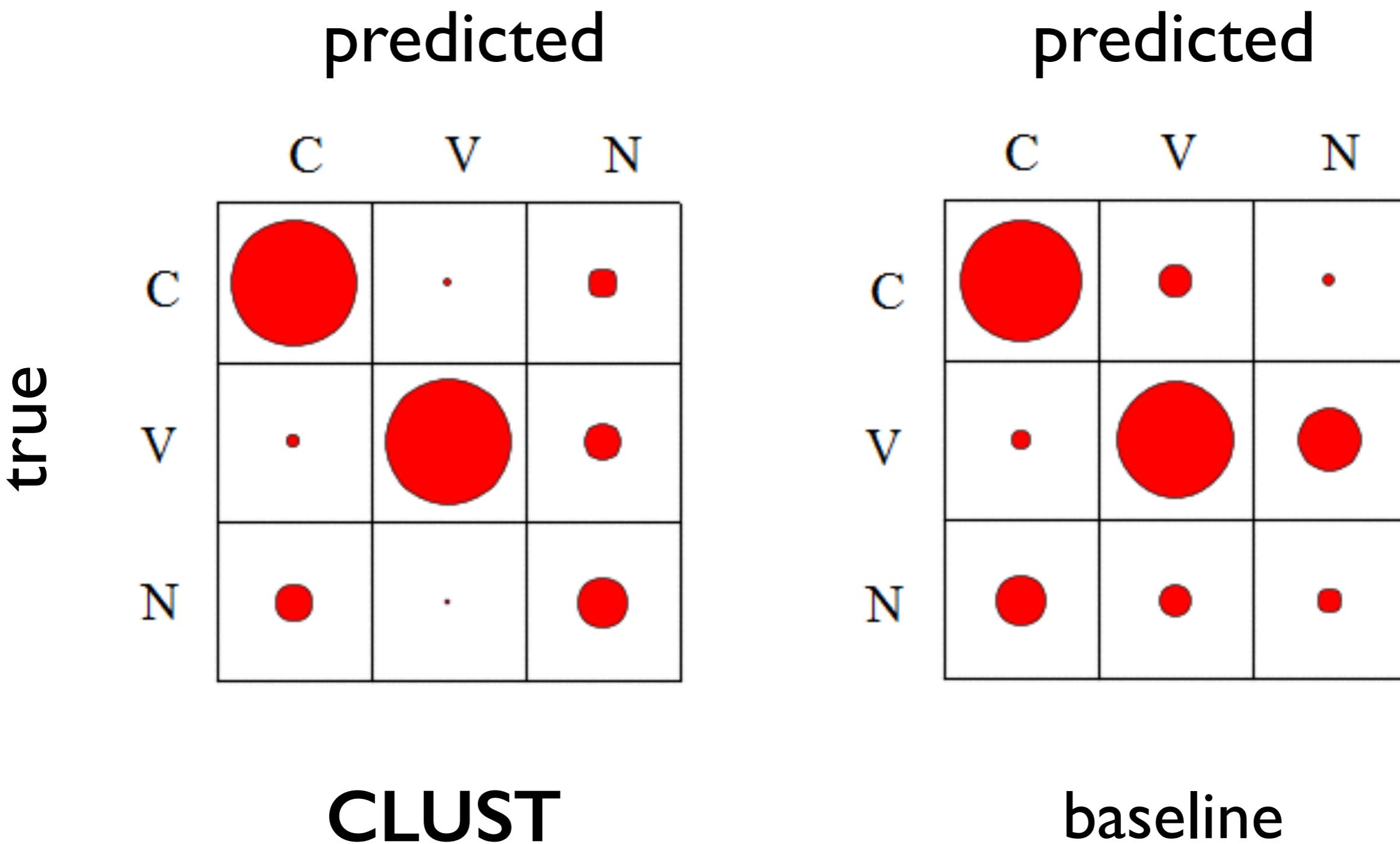


Results

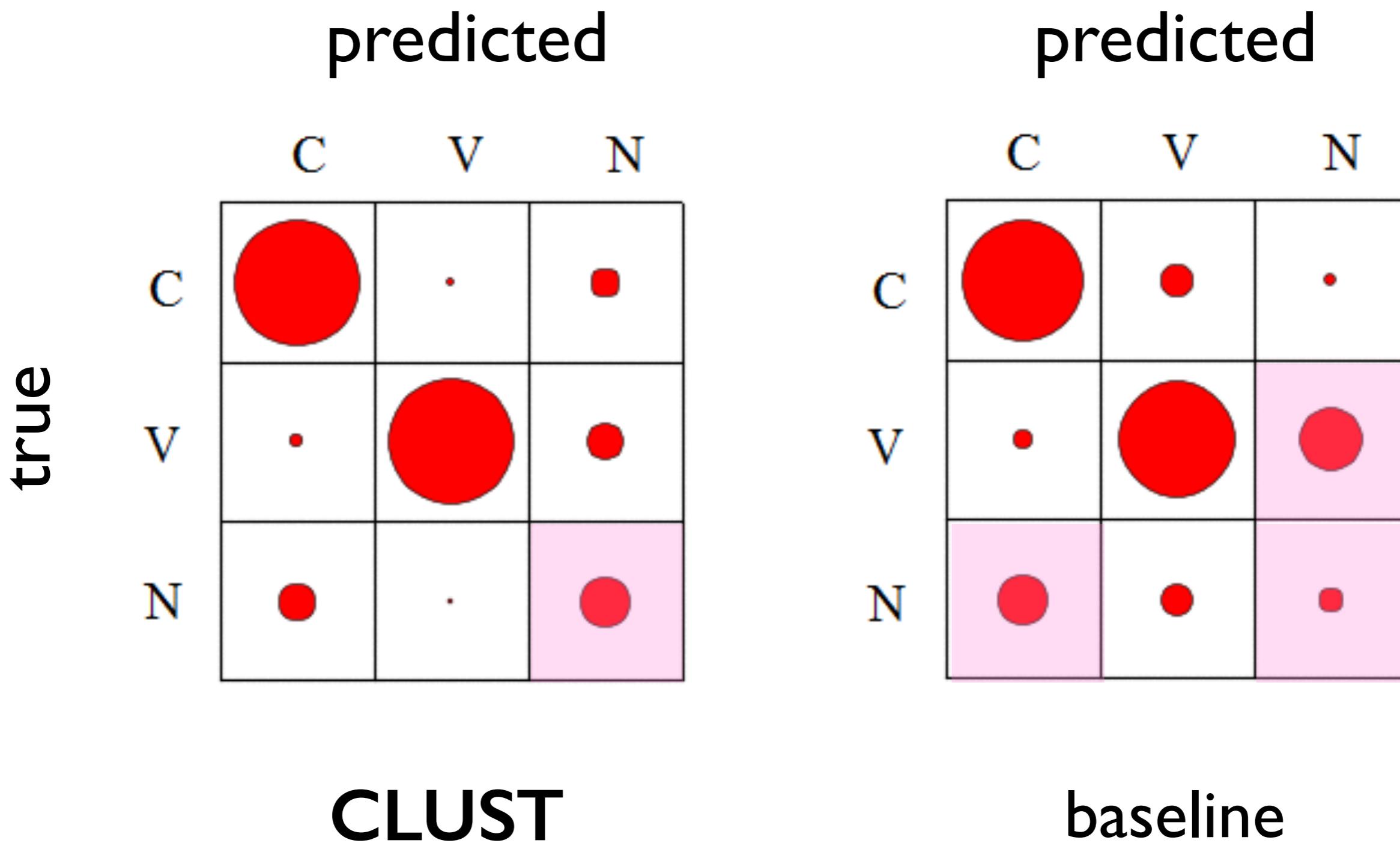
- 30 language isolates



Confusion matrix



Confusion matrix



Cluster Analysis

<i>Plurality Family</i>	<i>Proportion (%)</i>
Indo-European	38
	24
	21
Quechuan	89
Mayan	64
Oto-Manguean	55
Maipurean	25
Tucanoan	20
Uto-Aztecian	40
Altaic	44

<i>Plurality Family</i>	<i>Proportion (%)</i>
Austronesian	91
	71
	24
Niger-Congo	100
	78
	74
	68
	67
	50
	24

Future Work

- Finer-grained phonetic categories
- Syllable writing systems
 - ~ Syllables as hidden variables in our data
- Apply to Linear A

Thank you!