

Diachrony-aware Induction of Binary Latent Representations from Typological Features

MURAWAKI Yūgo
Kyoto University



The source code is available at
<https://github.com/murawaki/latent-typology/>

IJCNLP? No Text Processing!

And I abandoned a *neural network* model I proposed 2 years ago, and switched back to a *Bayesian* model

IJCNLP? No Text Processing!

And I abandoned a *neural network* model I proposed 2 years ago, and switched back to a *Bayesian* model

- What I am playing with is matrix X
 - $L = 2,679$, $N = 104$

N categorical features

L languages	1	1	...	2
	2	4	...	3
	1	1	...	3
	\vdots	\vdots	\vdots	\vdots
	3	1	...	1

X

IJCNLP? No Text Processing!

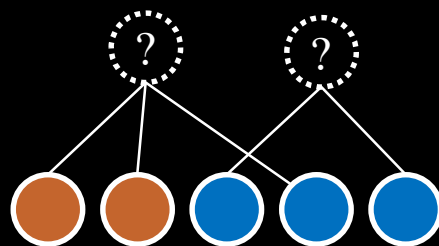
And I abandoned a *neural network* model I proposed 2 years ago, and switched back to a *Bayesian* model

- What I am playing with is matrix X
 - $L = 2,679$, $N = 104$
- And make use of two additional resources

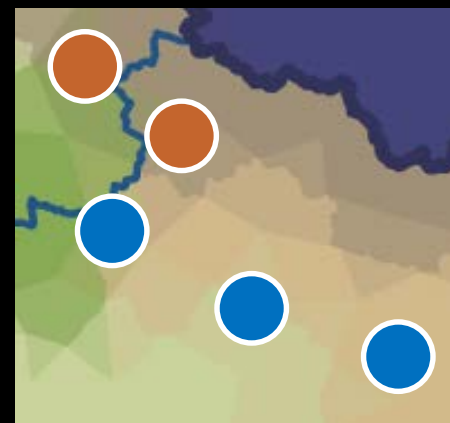
N categorical features

L languages	1	1	...	2
	2	4	...	3
	1	1	...	3
	\vdots	\vdots	\vdots	\vdots
	3	1	...	1

X



Phylogenetic groupings
of languages



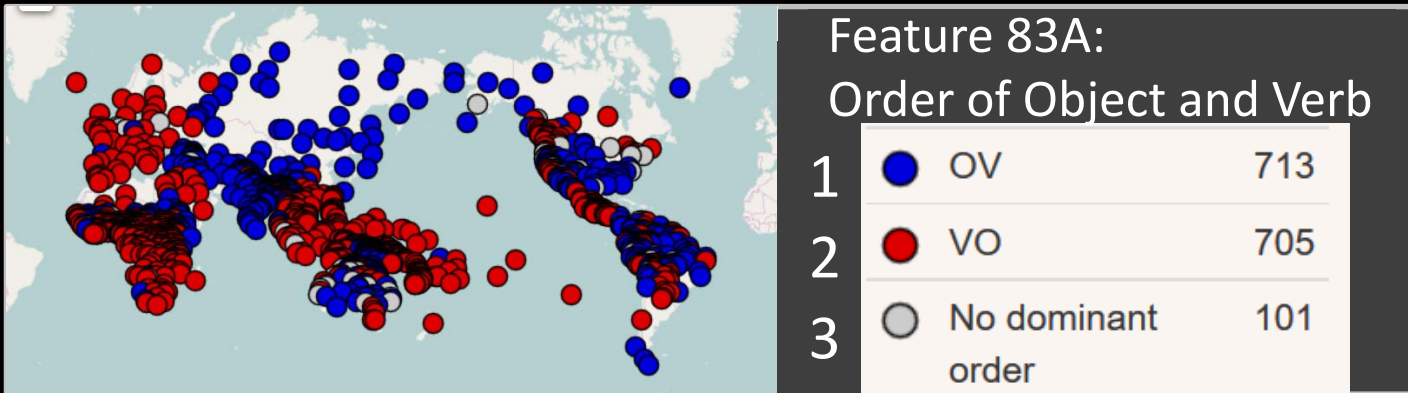
Spatial locations of
languages

Features of Linguistic Typology

N categorical features

L languages	1	1	...	2
	2	4	...	3
	1	1	...	3
	⋮	⋮	⋮	⋮
	3	1	...	1
X				

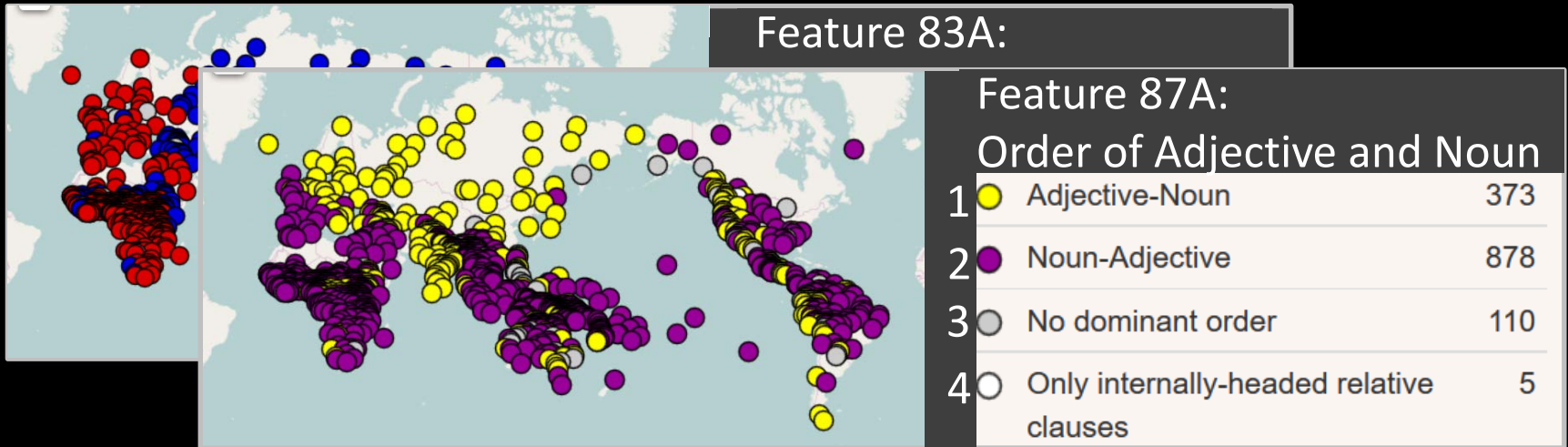
Features of Linguistic Typology



N categorical features

L languages	1	1	...	2
	2	4	...	3
	1	1	...	3
	⋮	⋮	⋮	⋮
	3	1	...	1
	X			

Features of Linguistic Typology



N categorical features

L languages	1	1	...	2
	2	4	...	3
	1	1	...	3
	⋮	⋮	⋮	⋮
	3	1	...	1

X

Questions

(Not Addressed in This Talk)

1. How have these features changed over time?
2. Can we use these features to trace the deep history of languages?

N categorical features

L languages	1	1	...	2
	2	4	...	3
	1	1	...	3
	\vdots	\vdots	\vdots	\vdots
	3	1	...	1
X				

Missing Values are a Real Problem

- 73.1% of items are missing

N categorical features

L languages	1	1	...	2
	?	4	...	3
	1	?	...	?
	:	:	:	:
	3	1	...	?
X				

Missing Values are a Real Problem

- 73.1% of items are missing
- My model imputes missing values by combining
 1. **Synchronic** clues (inter-feature dependencies)
 2. **Diachronic** clues (inter-language dependencies)

N categorical features

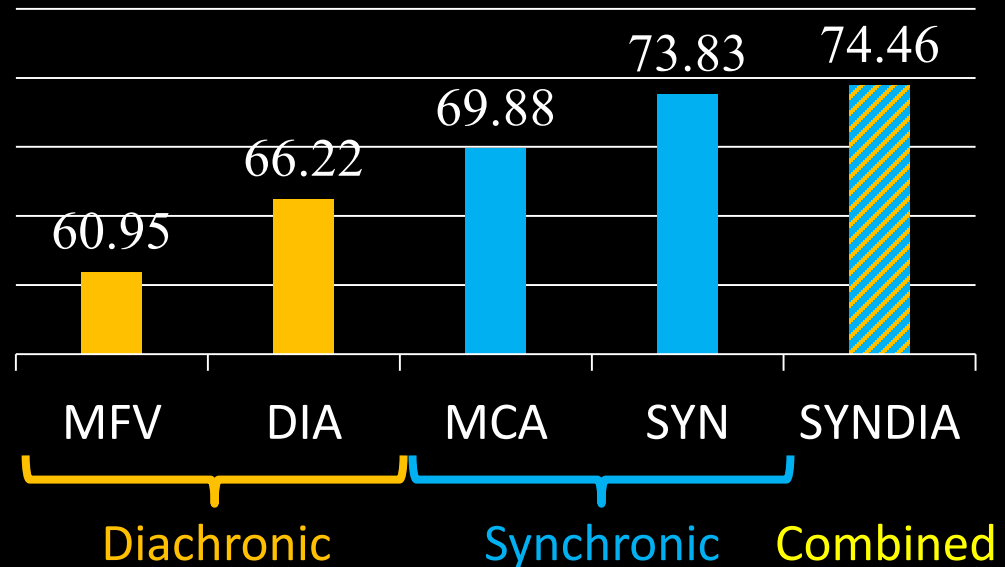
L languages	1	1	...	2
	?	4	...	3
	1	?	...	?
	\vdots	\vdots	\vdots	\vdots
	3	1	...	?
X				

Missing Values are a Real Problem

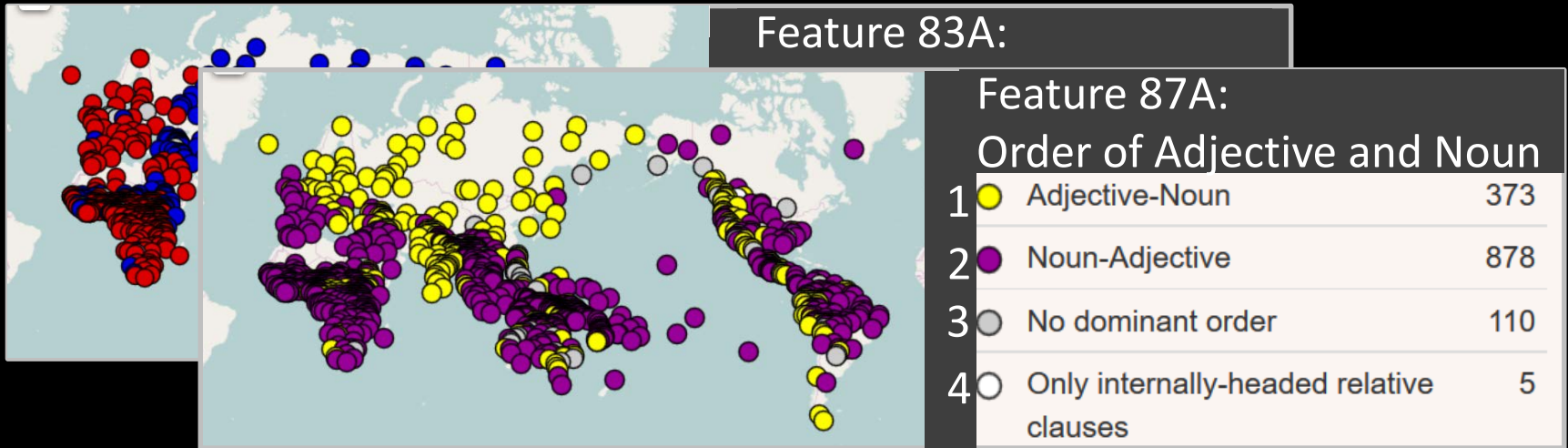
- 73.1% of items are missing
- My model imputes missing values by combining
 1. **Synchronic** clues (inter-feature dependencies)
 2. **Diachronic** clues (inter-language dependencies)

N categorical features

L languages	1	1	...	2
	?	4	...	3
	1	?	...	?
	:	:	:	:
	3	1	...	?
X				



Synchronic Clues: Inter-feature Dependencies

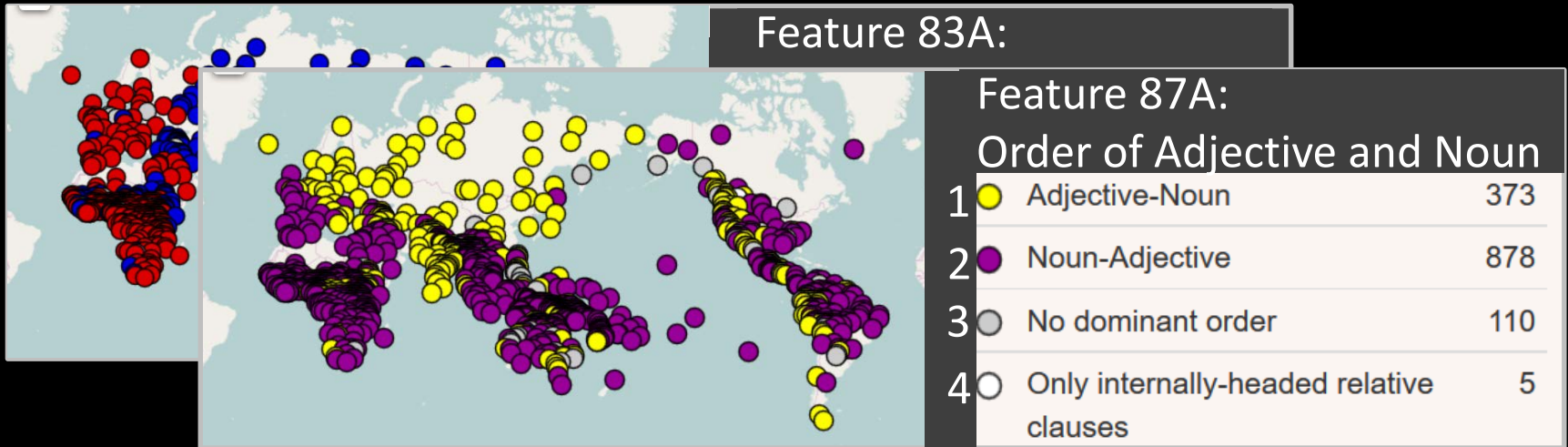


N categorical features

L languages	1	1	...	2
	?	4	...	3
	1	?	...	?
	:	:	:	:
	3	1	...	?

X

Synchronic Clues: Inter-feature Dependencies

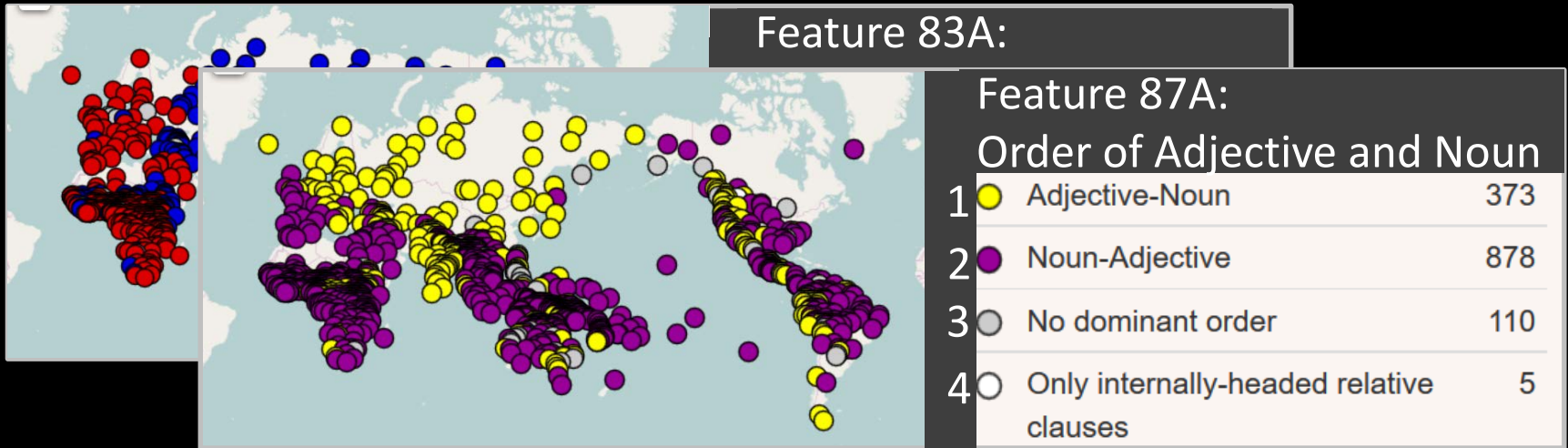


N categorical features

L languages	1	1	...	2
	?	4	...	3
	1	?	...	?
	:	:	:	:
	3	1	...	?
X				

If Feature 83A is 1 (OV),
then Feature 87A is likely
to be 1 (Adjective-Noun)

Synchronic Clues: Inter-feature Dependencies

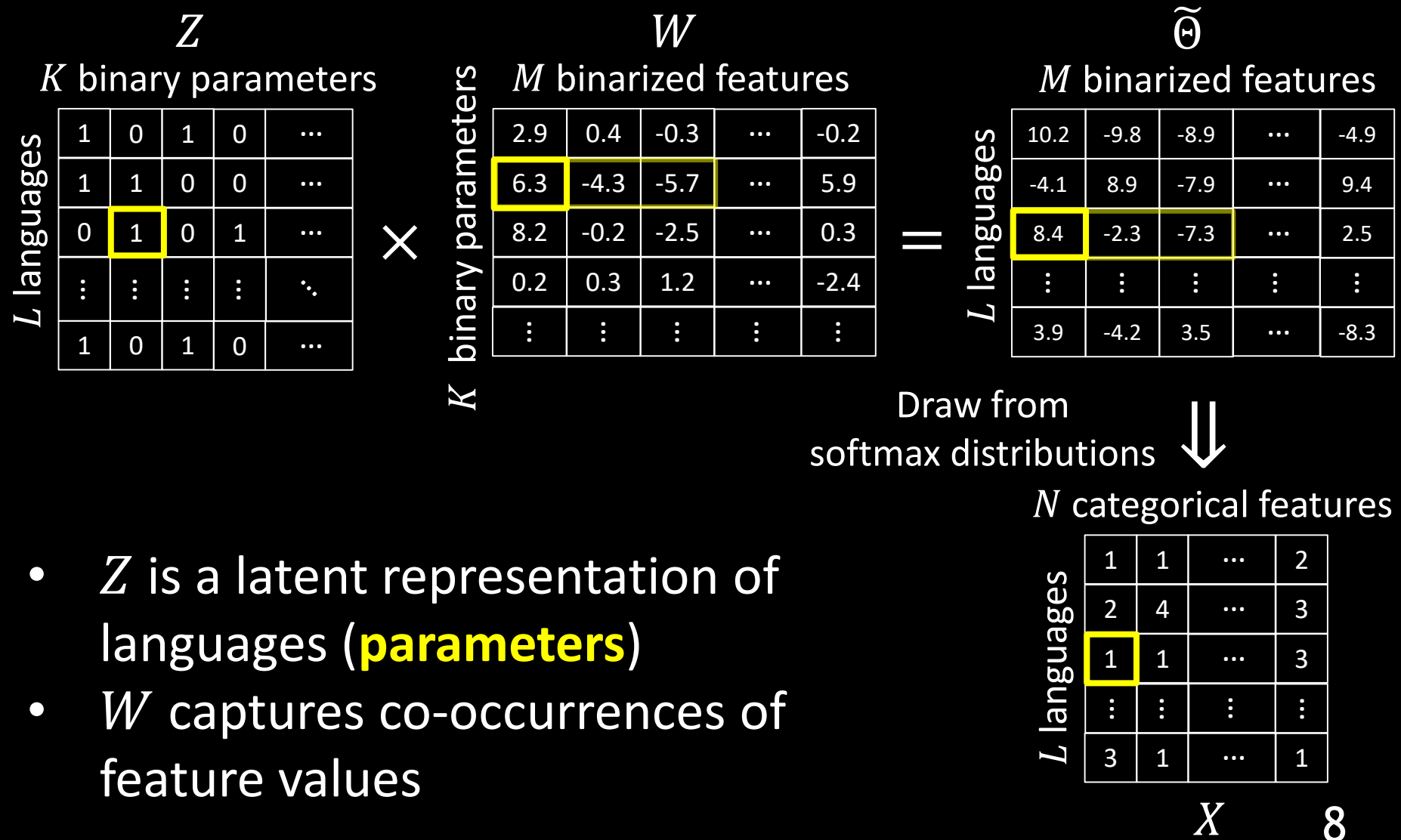


N categorical features

L languages	1	1	...	2
	?	4	...	3
	1	1	...	?
	:	:	:	:
	3	1	...	?
X				

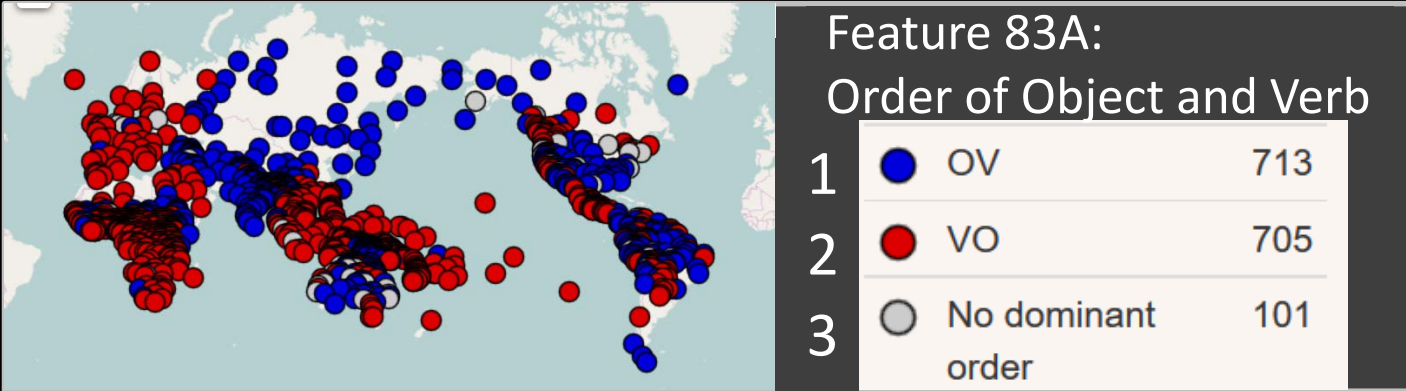
If Feature 83A is 1 (OV),
then Feature 87A is likely
to be 1 (Adjective-Noun)

Matrix Decomposition

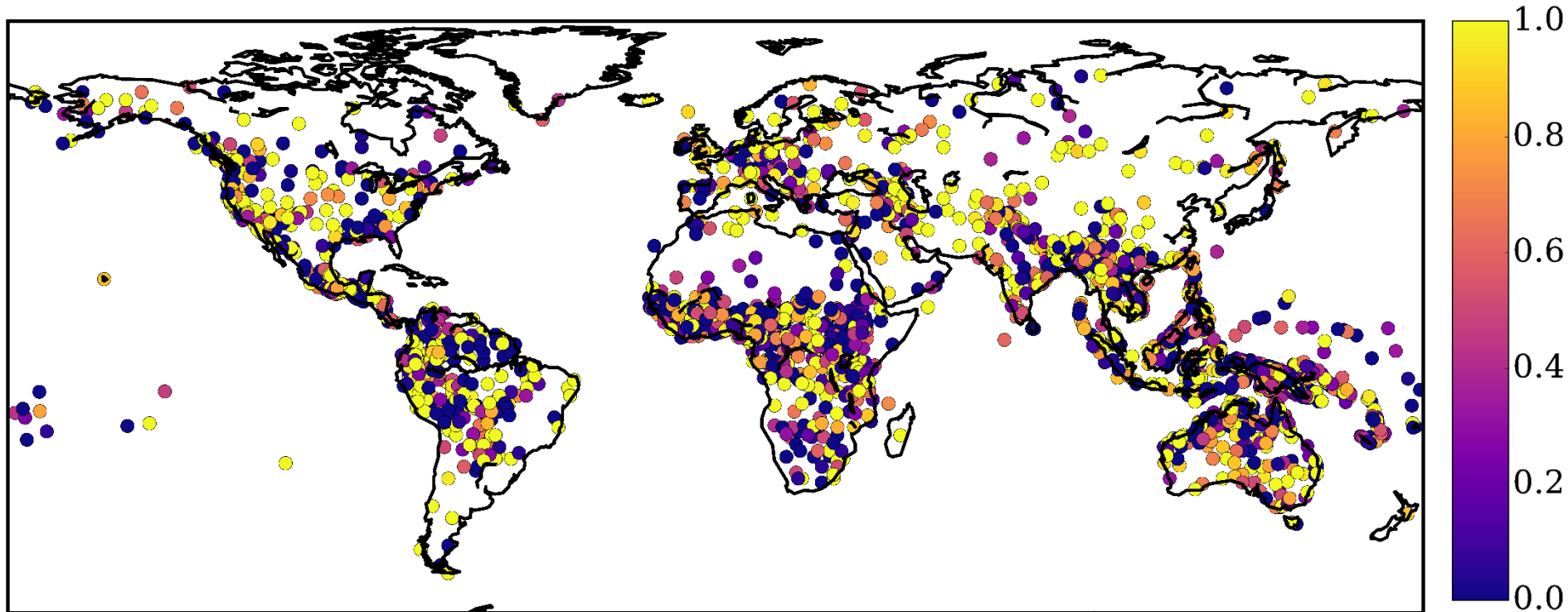


- Z is a latent representation of languages (**parameters**)
- W captures co-occurrences of feature values

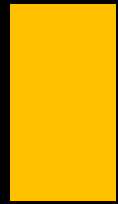
Geographic Distribution of an Induced Parameter



Geographic Distribution of an Induced Parameter



Geographic signals observed in surface features
are **somehow lost** during the induction

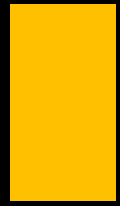


Diachronic Clues: Inter-language Dependencies

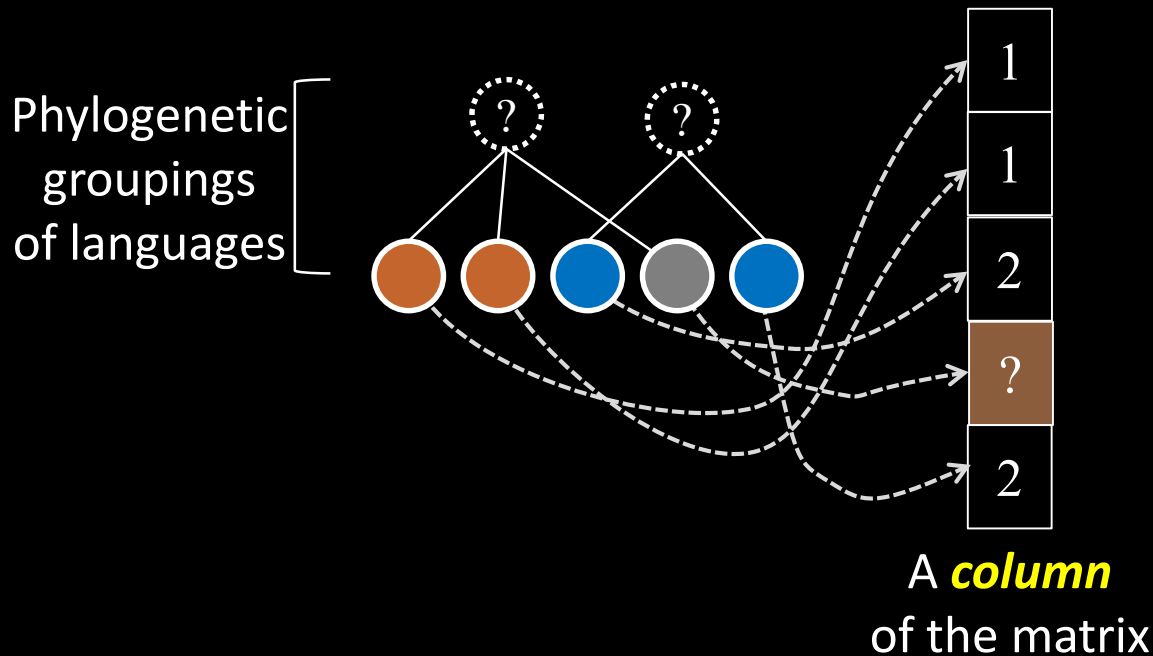
1
1
2
?
2

A **column**
of the matrix

Languages may share the same value because they ...



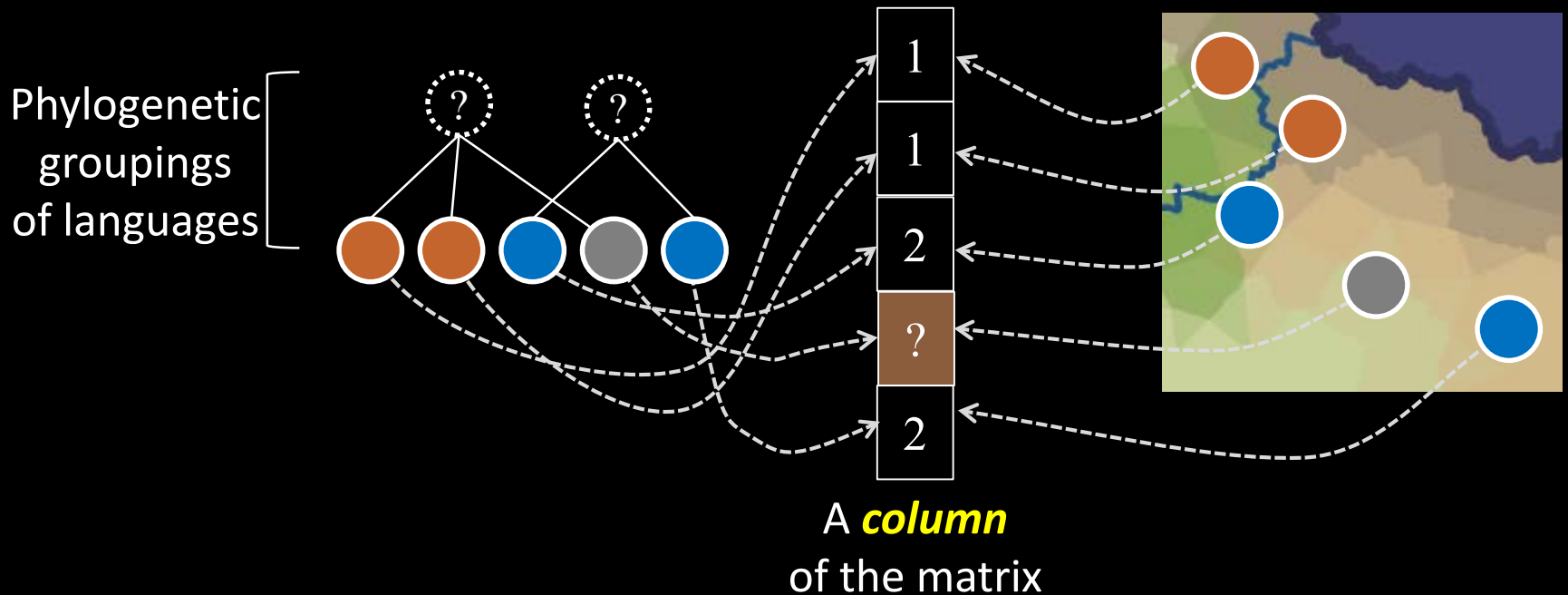
Diachronic Clues: Inter-language Dependencies



Languages may share the same value because they ...

1. are **phylogenetically** related,

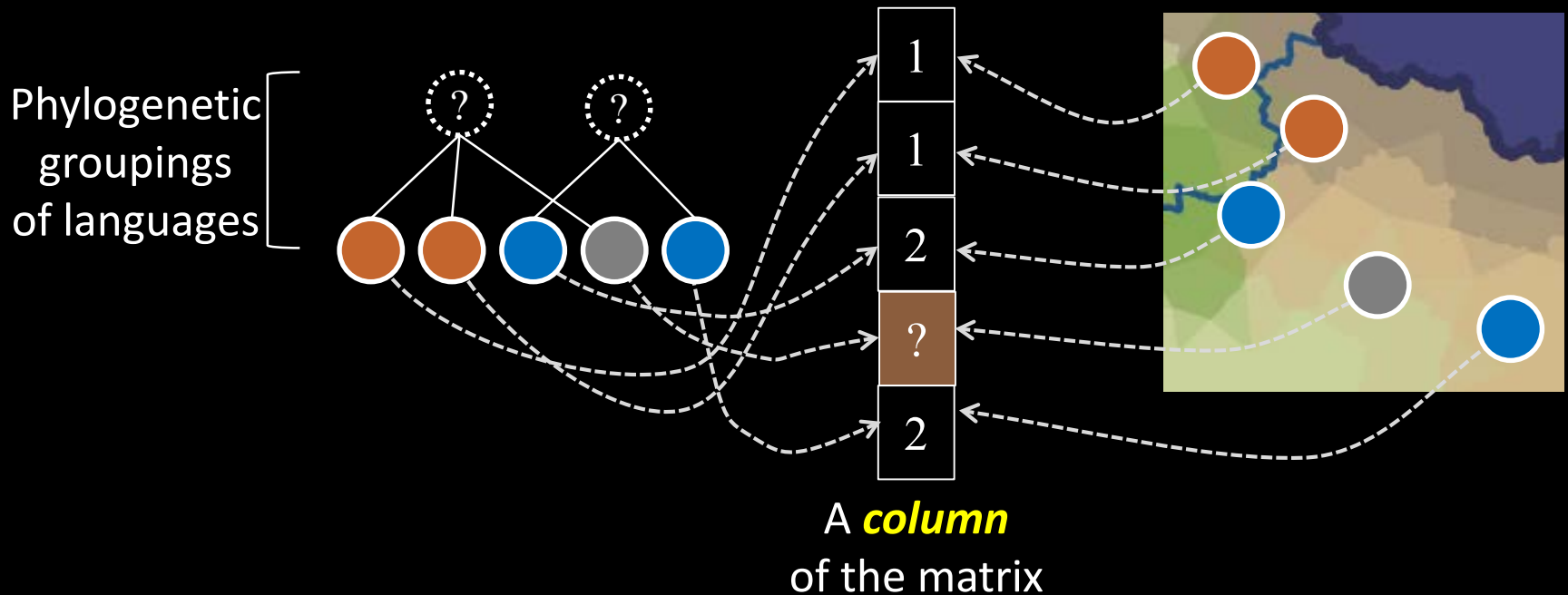
Diachronic Clues: Inter-language Dependencies



Languages may share the same value because they ...

1. are **phylogenetically** related,
2. are **spatially** close (in contact), and/or

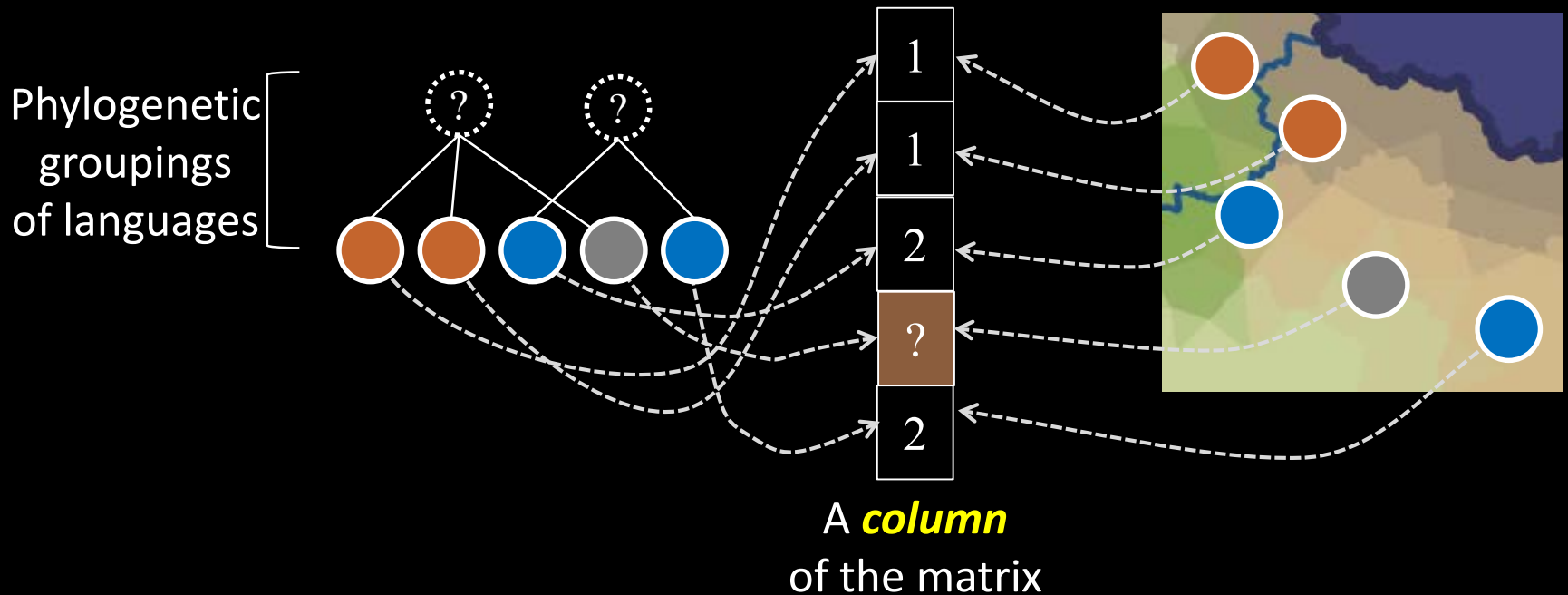
Diachronic Clues: Inter-language Dependencies



Languages may share the same value because they ...

1. are **phylogenetically** related,
2. are **spatially** close (in contact), and/or
3. just take a **universally** (globally) frequent value

Diachronic Clues: Inter-language Dependencies

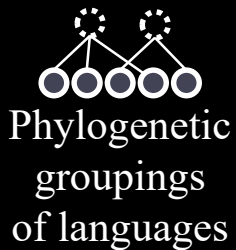


Languages may share the same value because they ...

1. are **phylogenetically** related,
2. are **spatially** close (in contact), and/or
3. just take a **universally** (globally) frequent value

Strength controlled by
vertical stability
horizontal diffusibility
universality

Autologistic Model



The prob. of language l taking value j for feature i is

$$P(x_{l,i} = j | \mathbf{x}_{-l,i}, v_i, h_i, u_i) \propto \exp(v_i V_{l,i,j} + h_i H_{l,i,j} + u_{i,j})$$

vertical stability

of l 's phylogenetic neighbors taking value j

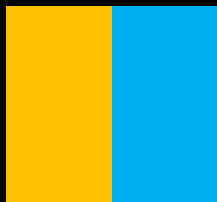
horizontal diffusibility

of l 's spatial neighbors taking value j

universality

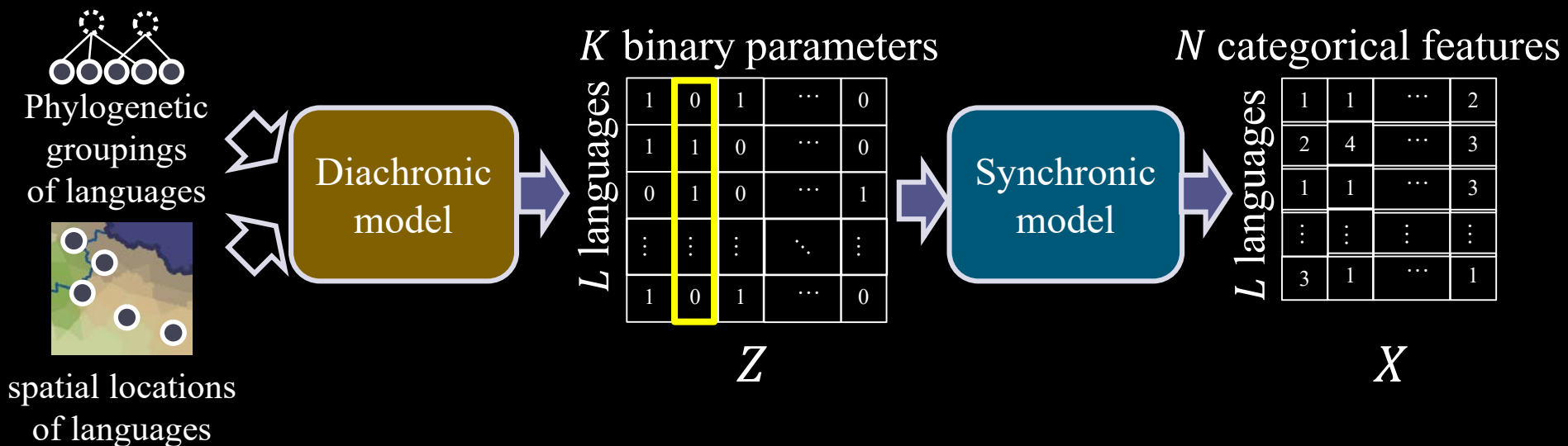
[Yamauchi+, COLING2016]

A Bayesian extension will appear in the Journal of Language Evolution



Combined Model:

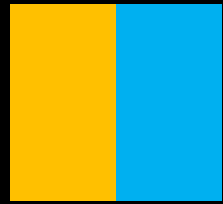
Diachronic model works in the latent space



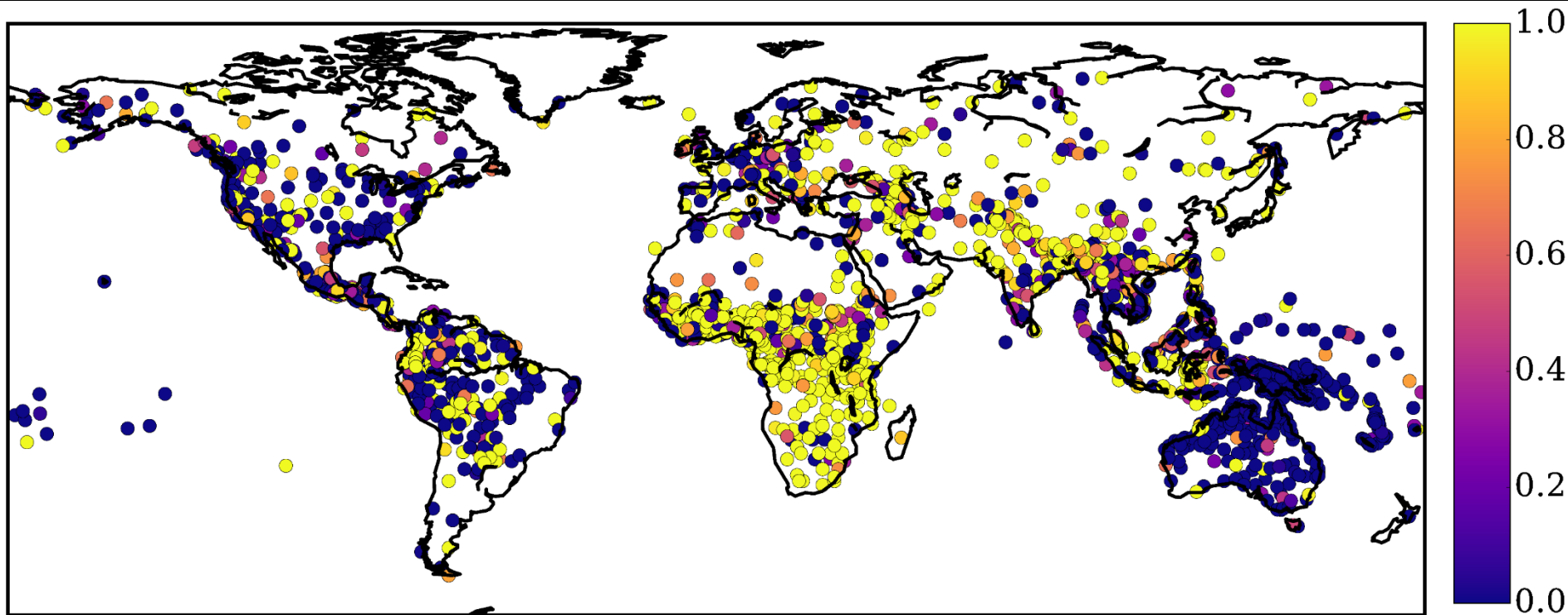
$$P(A, Z, W, X) = P(A)P(Z|A)P(W)P(X|Z, W)$$

Diachronic
model

Synchronic
model



Geographic Distribution of an Induced Parameter, Revised



Some parameters appear to preserve geographic signals observed in surface features

Conclusions

- Proposed a Bayesian model that combines synchronic and diachronic clues
- Induced a binary latent representation that appears to preserve geographic signals
- Remaining question: Can the latent representation help uncovering the deep history of languages?

