



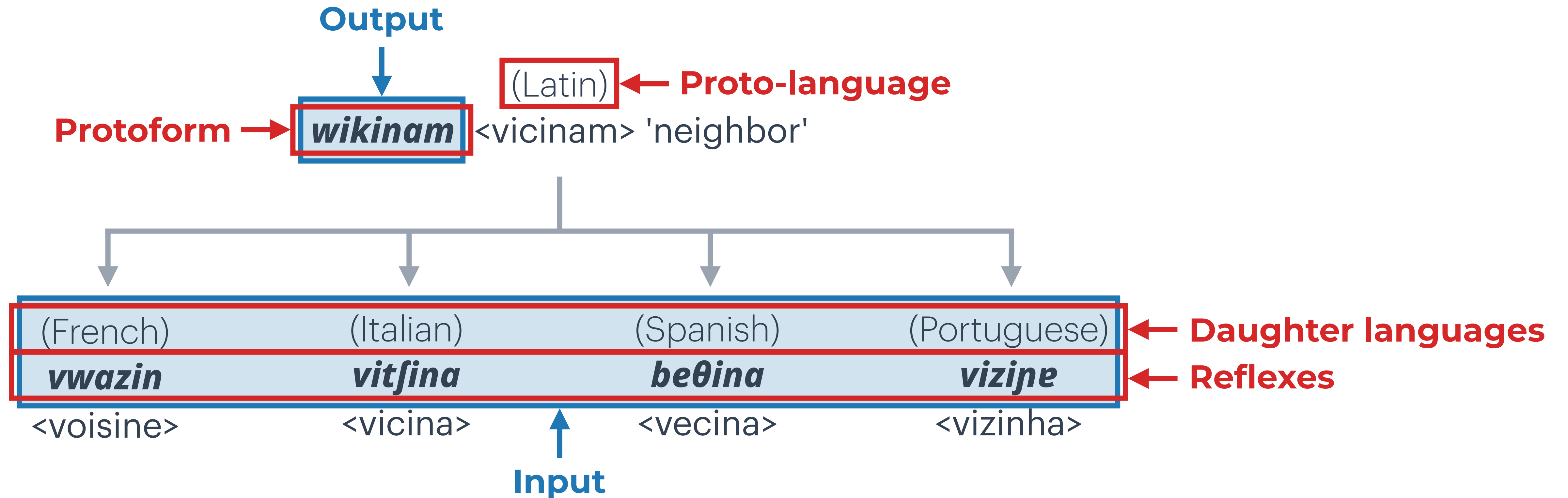
Semisupervised Neural Proto-Language Reconstruction

Liang Lu¹, Peirong Xie², David R. Mortensen¹

¹Carnegie Mellon University, ²University of Southern California

lianglu@cs.cmu.edu, louisxie@usc.edu, dmortens@cs.cmu.edu

Protoform Reconstruction



Shown: an example from the Romance dataset (Meloni et al., 2021; Ciobanu and Dinu, 2018)

A 19th-Century Discovery

Languages change in systematic ways, and it is possible to reproducibly reconstruct proto-languages using these systematic patterns, even when no record of the proto-language survives.

Historical linguists use the **comparative method** to reconstruct proto-languages.

The comparative method's **regularity principle**:

- ▶ Sound changes are regular
- ▶ Reflexes should be derivable deterministically from reconstructions using a single set of sound change rules

Ancestor language



Systematic changes

Descendent language

Supervised Neural Reconstruction

- ▶ **RNN with language embedding** (Meloni et al., 2021)
- ▶ **Transformer** (Kim et al., 2023)
- ▶ **VAE** (Cui et al., 2022)



Shown: the 媚 cognate set from the Chinese dataset (Chang et al., 2022)

A Hypothetical Example: Tangkhulic Languages

Gloss	'grandchild'	'bone'	'breast'	'laugh'
Labeled?	Yes	Yes	No	No
Kachai	ð e	r e	n e	n i
Huishu	r u k	r u k	n u k	n u k
Ukhrul	r u	r u	n u	n u
Protoform Label	d u	r u	n u	n ï

A Hypothetical Example: Tangkhulic Languages

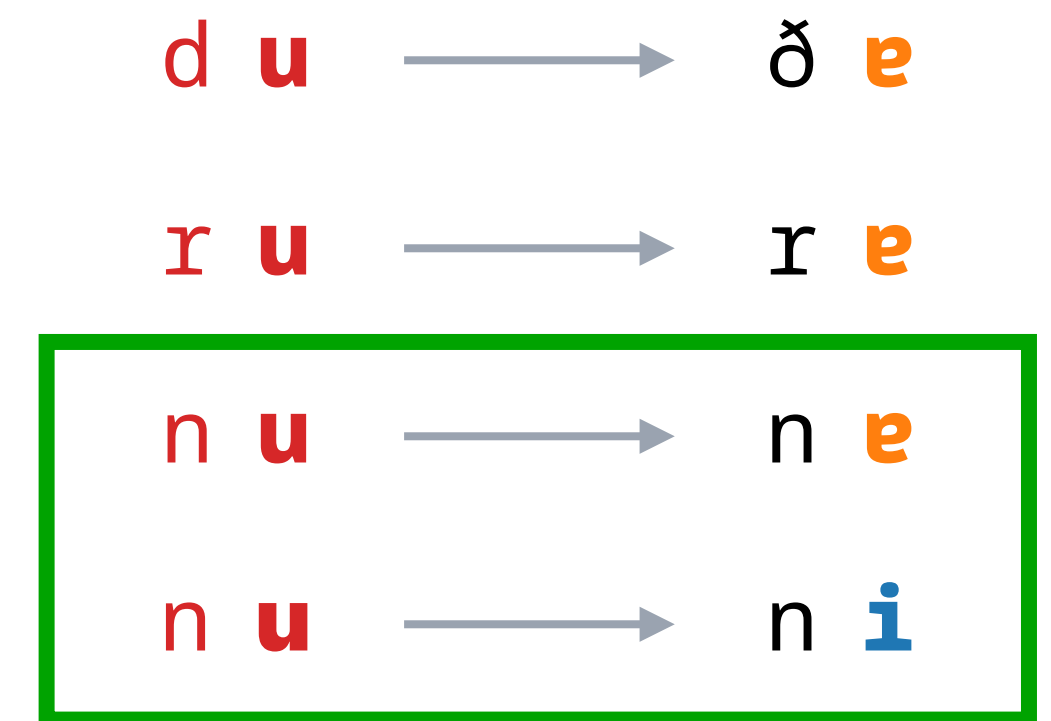
Gloss	'grandchild'	'bone'	'breast'	'laugh'
Labeled?	Yes	Yes	No	No
Kachai	ǎ e	r e	n e	n i
Huishu	r u k	r u k	n u k	n u k
Ukhrul	r u	r u	n u	n u
Protoform Label	d u	r u	(hidden)	(hidden)

A Hypothetical Example: Tangkhulic Languages

Gloss	'grandchild'	'bone'	'breast'	'laugh'
Labeled?	Yes	Yes	No	No
Kachai	ǎ e	r e	n e	n i
Huishu	r u k	r u k	n u k	n u k
Ukhrul	r u	r u	n u	n u
Protoform Label	d u	r u	(hidden)	(hidden)
Supervised Only	d u	r u	n u	n u

A Hypothetical Example: Tangkhulic Languages

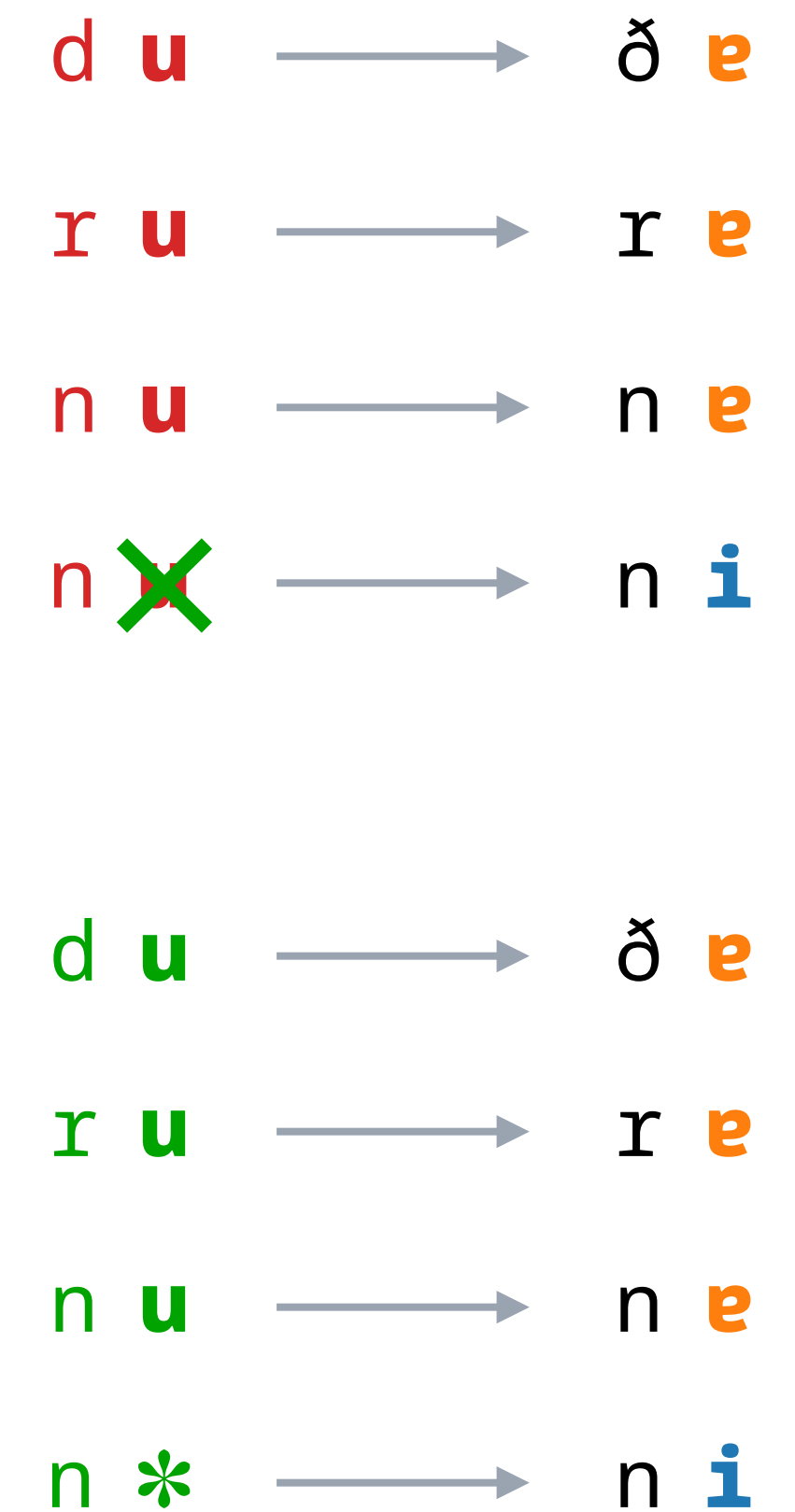
Gloss	'grandchild'	'bone'	'breast'	'laugh'
Labeled?	Yes	Yes	No	No
Kachai	ð e	r e	n e	n i
Huishu	r u k	r u k	n u k	n u k
Ukhrul	r u	r u	n u	n u
Protoform Label	d u	r u	(hidden)	(hidden)
Supervised Only	d u	r u	n u	n u



↑
 Trouble: Cannot deterministically derive the reflexes!

A Hypothetical Example: Tangkhulic Languages

Gloss	'grandchild'	'bone'	'breast'	'laugh'
Labeled?	Yes	Yes	No	No
Kachai	ð e	r e	n e	n i
Huishu	r u k	r u k	n u k	n u k
Ukhrul	r u	r u	n u	n u
Protoform Label	d u	r u	(hidden)	(hidden)
Supervised Only	d u	r u	n u	n u
Our Model	d u	r u	n u	n *
				↑ Something other than u



A Hypothetical Example: Tangkhulic Languages

Gloss	'grandchild'	'bone'	'breast'	'laugh'
Labeled?	Yes	Yes	No	No
Kachai	ǎ e	r e	n e	n i
Huishu	r u k	r u k	n u k	n u k
Ukhrul	r u	r u	n u	n u
Protoform Label	d u	r u	n u	n i ← Indeed not u
Supervised Only	d u	r u	n u	n u
Our Model	d u	r u	n u	n * ↑ Something other than u

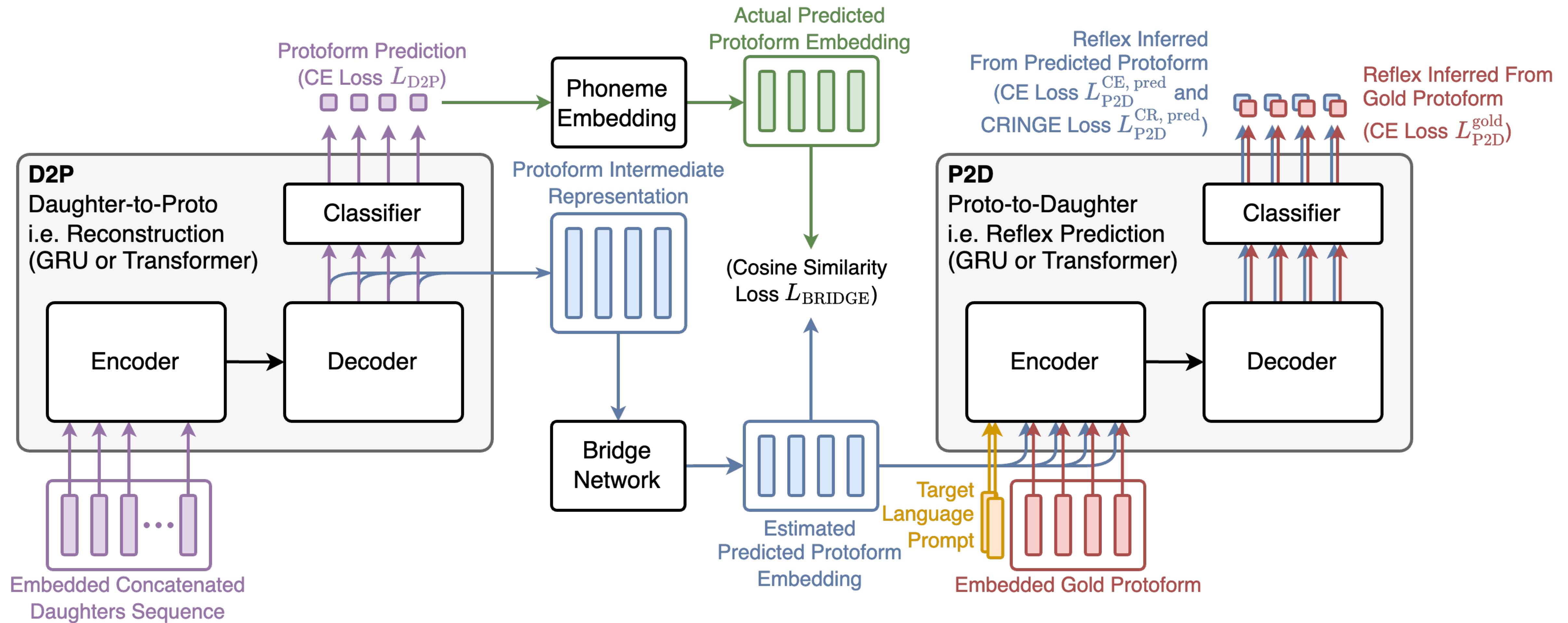
A Hypothetical Example: Tangkhulic Languages

Gloss	'grandchild'	'bone'	'breast'	'laugh'
Labeled?	Yes	Yes	No	No
Kachai	ð e	r e	n e	n i
Huishu	r u k	r u k	n u k	n u k
Ukhrul	r u	r u	n u	n u
Protoform Label	d u	r u	n u	n i
Our Model: DPD	d u	r u	n u	n *

Reflexes $\xrightarrow{\text{Reconstruction}}$ Protoform? $\xrightarrow{\text{Reflex Prediction}}$ Reflexes?

Daughter-to-Proto-to-Daughter (DPD)

The DPD (Dughter-to-Proto-to-Daughter) Architecture



Training Strategies

1. Supervised only (**SUPV**)] **Naïve baseline**
2. Bootstrapping (**BST**)]
3. Π -model (**Π M**)] **Semisupervised techniques**
4. Π -model with Bootstrapping (**Π M-BST**)] **Combination of semisupervised techniques**
5. **DPD**]
6. DPD with Bootstrapping (**DPD-BST**)]
7. DPD merged with Π -model (**DPD- Π M**)] **DPD-based strategies**
8. DPD- Π M with Bootstrapping (**DPD- Π M-BST**)]

Semisupervised Datasets

We randomly remove labels from existing fully-labeled datasets to create semisupervised train sets.

	Chinese (Sinitic)	Romance
5%	181	304
10%	362	607
20%	723	1,214
30%	1,084	1,821
100%	3,615	6,071

Number of labeled training examples (i.e. cognate sets with an associated gold protoform) in the train set for each labeling setting and dataset, as well as the total number of cognate sets for reference (100%).

Results

For all data settings, sub-model architectures, and language families, **the winning model** (by best average accuracy) **is always a DPD model**.

The **DPD- Π M-BST** model that combines **DPD**, **Π -Model**, and **Bootstrapping** is usually the best.

Language Family	Sub-Model Architecture	5% Labeled	10% Labeled	20% Labeled	30% Labeled
Chinese	Transformer	DPD	DPD- Π M-BST	DPD- Π M-BST	DPD- Π M
	GRU	DPD- Π M-BST	DPD- Π M-BST	DPD- Π M-BST	DPD- Π M
Romance	Transformer	DPD- Π M-BST	DPD- Π M-BST	DPD- Π M-BST	DPD- Π M-BST
	GRU	DPD-BST	DPD-BST	DPD- Π M-BST	DPD-BST

Conclusion

DPD marks a step forward toward **building practical computational reconstruction systems** that can assist early-stage proto-language reconstruction projects.

Even in the age of Transformers, **linguistic ideas that predate NLP**, such as the comparative method, can still **provide valuable insights for approaching NLP challenges**.