

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

Abstract

Existing work implementing comparative reconstruction of ancestral languages (proto-languages) has usually required full supervision. However, historical reconstruction models are only of practical value if they can be trained with a limited amount of labeled data. We propose a semisupervised historical reconstruction task in which the model is trained on only a small amount of labeled data (cognate sets with proto-forms) and a large amount of unlabeled data (cognate sets without proto-forms). We propose a neural architecture for comparative reconstruction (DPD-BiReconstructor) incorporating an essential insight from linguists’ comparative method: that reconstructed words should not only be reconstructable from their daughter words, but also deterministically transformable back into their daughter words. We show that this architecture is able to leverage unlabeled cognate sets to outperform strong semisupervised baselines on this novel task¹.

1 Introduction

In the 19th century, European philologists made a discovery that would change the direction of the human sciences: they discovered that languages change in systematic ways and that, by leveraging these systematic patterns, it was possible to reproducibly reconstruct ancestors of families of languages (protolanguages) even when no record of those languages survived. This technique, called the comparative method, provided an unprecedented window into the human past—its cultures, its migrations, and its encounters between populations.

The assumption that historical changes in pronunciation (“sound changes”) are regular, known as ‘the regularity principle’ or ‘the Neogrammarian hypothesis’, is fundamental to the comparative

¹The code is available at <https://github.com/cmu-llab/dpd>.

	‘grandchild’	‘bone’	‘breast’	‘laugh’
Kachai	ðe	re	ne	ni
Huishu	ruk	ruk	nuk	nuk
Ukhrul	ru	ru	nu	nu
Reference	*du	*ru	*nu	*ni
Label	Y	Y	N	N
D2P	*du	*ru	*nu	*nu
DPD	*du	*ru	*nu	*nU

Table 1: Hypothetical illustration of the shortcomings of daughter-to-proto (D2P) models drawn from Tangkhulic Languages. “U” represents a sound other than “u” predicted by a hypothetical daughter-to-proto-to-daughter (DPD) model. “Y” and “N” indicate whether the reference protoform is labeled.

method (Campbell, 2021). As the 19th century Neogrammarians Hermann Osthoff and Karl Brugmann put it:

“Every sound change, in so far as it proceeds mechanically, is completed in accordance with laws admitting of no exceptions; i.e. the direction in which the change takes place is always the same for all members of a language community, apart from the case of dialect division, and all words in which the sound subject to change occurs in the same conditions are affected by the change without exception.” (*Morphologische Untersuchungen auf dem Gebiete der indogermanischen Sprachen* i, Brugmann and Osthoff, 1878, p. xiii, translated and quoted in Szemerényi, 1996)

The comparative method, however, is challenging for humans to apply. This is true largely because it involves dealing with large volumes of data and modeling numerous interactions between competing patterns. One must balance the need for phonetic similarity between reconstructed words and their descendants (reflexes) with the need to be able to deterministically derive reflexes from

reconstructed words with a single set of sound changes. It imposes a heavy cognitive load. For this reason, researchers have long aspired to implement the comparative method computationally.

With some exceptions (He et al., 2023; Akavaram and Bhattacharya, 2023), recent attempts at automatic reconstruction models mostly take the form of neural sequence-to-sequence transduction models similar to those used for machine translation, trained with full supervision on a dataset where every cognate set is paired with a gold reconstruction (Meloni et al., 2021; Chang et al., 2022; Fourrier, 2022; Cui et al., 2022; Kim et al., 2023).

The development of supervised reconstruction systems has given the field insights into how reconstruction of protolanguages (unattested ancestor languages) can be modelled computationally. However, in a realistic scenario, these models only become usable once the hardest part of reconstruction has been done (since they rely on the linguist having already identified enough of the sound changes in the data to reconstruct a substantial part of the lexicon). Computational models of comparative reconstruction are most useful if they can be deployed without training data, or if only a small volume of labeled data is needed to prime the comparative pump.

We introduce a semisupervised protoform (reconstructed parent word) reconstruction task wherein the reconstruction model has access at training time to both a small number of cognate sets (sets of daughter words—reflexes—of a single parent) labeled with a protoform and a large number of unlabeled cognate sets, mirroring the situation of historical linguists early in their reconstruction of a protolanguage. Though similar to semisupervised machine translation, the semisupervised reconstruction formulation entails the absence of target-side monolingual data. Most semisupervised machine translation techniques rely on monolingual data in the target language, such as back-translation (Edunov et al., 2018; Sennrich et al., 2016) and pre-trained target-side language models (Skorokhodov et al., 2018; Gülçehre et al., 2015). In contrast, in this task, models have access only to cognate sets (no monolingual text), meaning the structure of the problem is quite different.

In this paper, we propose to incorporate the comparative method into a semisupervised reconstruction model via end-to-end multi-task learning. Our

proposed model, named DPD-BiReconstructor, learns to improve its reconstructions by performing reflex predictions on an intermediate representation of its predicted reconstructions. Reflex prediction losses are propagated into the reconstruction network, allowing the model to train on cognate sets without protoform labels. A hypothetical example from three Tangkhulic languages is shown in Table 1. In this example, the phonetic information in the ‘grandchild’ and ‘bone’ sets is insufficient to reconstruct ‘laugh’ with a distinct vowel, so daughter-to-*proto* (D2P) models will typically reconstruct the two words identically. Models incorporating reflex prediction, however, are able—in principle—learn to reconstruct words like ‘laugh’ as distinct from words like ‘breast’. Experiments show that DPD is an attractive approach for semisupervised reconstruction, and a combination of DPD with existing semisupervised strategies performs significantly better than baseline strategies in almost all situations. Additionally, analyses show some indication that DPD-based models could help improve supervised reconstruction.

2 Methods

2.1 Model

We propose a multi-task reconstruction strategy that learns to recover reflexes from its own reconstructions, effectively utilizing unlabeled cognate sets. Our model comprises a reconstruction sub-network (D2P for daughter-to-*proto*form) and a reflex-prediction sub-network (P2D for *proto*form-to-daughter) with shared phoneme embeddings. On labeled data, the model learns sound changes from accurate reconstructions to reflexes, in addition to learning reconstruction. In the absence of labeled *proto*forms, the reflex prediction sub-network acts as a weak supervision by informing the reconstruction sub-network on whether correct reflexes can be derived from proposed reconstructions. This workflow directly mirrors the comparative method’s constraint that reconstruction must yield *proto*forms such that reflexes can be derived from them through regular sound changes. We refer to our training strategy and its architectural realization as Daughters-to-*Proto*form-to-Daughters Bidirectional Reconstructor (DPD).

For D2P to learn from P2D, we propagate gradients from P2D into D2P. It is not feasible, however, to simply feed token predictions from D2P into P2D, as token outputs are not differentiable.

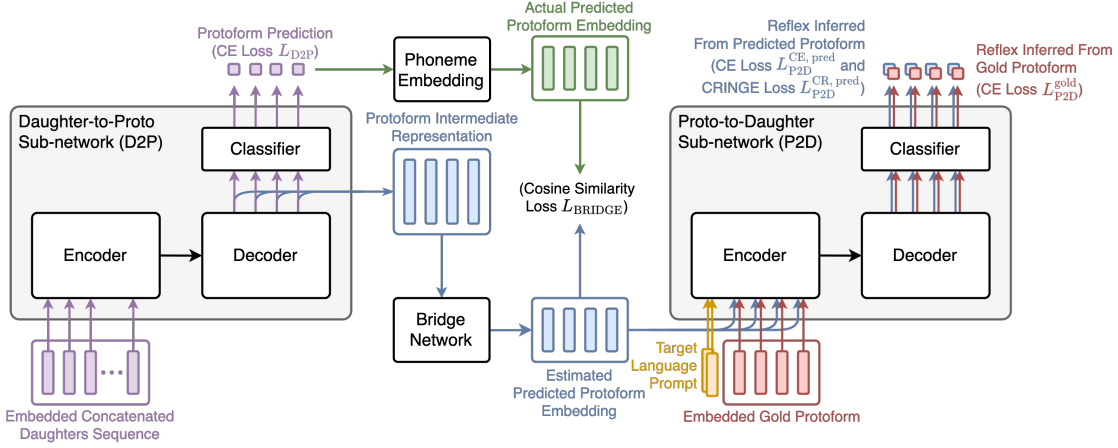


Figure 1: The DPD architecture, with a reconstruction sub-network (D2P), a reflex prediction sub-network (P2D), and a dense bridge network. The bridge network connects the final-layer decoder state from D2P to the encoder input of P2D. For labeled data, supervised cross-entropy (CE) is computed for D2P’s protoform prediction L_{D2P} and P2D’s reflex prediction from the gold protoform L_{P2D}^{gold} . On both labeled and unlabeled data, reflex prediction losses (based on the predicted protoform) consisting of a CE loss $L_{P2D}^{CE, pred}$ and a CRINGE loss $L_{P2D}^{CR, pred}$ (for incorrectly predicted protoforms) are propagated into both sub-networks. An additional cosine similarity loss L_{BRIDGE} is used to train the bridge network. The overall loss is calculated as a weighted sum of all the losses, with the weights being tunable hyperparameters. Reflexes in the same cognate are concatenated together (with separators) into one sequence and embedded with both phoneme and language embedding as in Meloni et al. (2021) and Kim et al. (2023) (and an additional positional embedding for Transformer). Phoneme embedding is shared among D2P and P2D, whereas language embedding is only used in D2P.

Instead, we represent reconstruction outputs as a continuous latent space, inspired by end-to-end spoken language understanding systems in which a semantic understanding sub-network receives a latent representation of text transcriptions from a speech-recognition sub-network (Saxon et al., 2021; Arora et al., 2022). In particular, D2P’s final-layer decoder output is connected to P2D’s encoder via a trainable dense network, referred to as the bridge network (See Figure 1).

There is no easy way to transform gold protoforms used to train P2D so that they match the bridge network’s output latent representation. To encourage a consistent input representation for P2D, we add a cosine similarity loss between the bridge network’s output and the actual phoneme embeddings of D2P’s predictions, effectively training the bridge network to serve as an alternative embedding for reconstructed protoforms.

As a final training objective, we discourage P2D from producing correct reflexes given incorrect protoforms so that D2P can be better informed when producing incorrect reconstructions. This is achieved via a CRINGE loss on the P2D outputs, designed to penalize negative examples (Adolphs et al., 2022). For negative examples, the CRINGE loss contrasts each negative token with

an estimated positive token (by sampling a non-negative token in the model’s top- k predictions) and applies a contrastive learning objective (cross-entropy loss) to lower the probability of the negative token. When training the P2D sub-network, a correctly predicted reflex is considered a negative example if the input was an incorrect protoform.

The DPD strategy does not specify the architecture of the D2P and P2D sub-networks. In this paper, we adapt existing encoder-decoder architectures proposed for neural reconstruction, including GRU sub-networks based on Meloni et al. (2021)² and Transformer sub-networks based on Kim et al. (2023).

2.2 Semisupervised Strategies

Aside from our DPD strategy, a naïve approach is to discard unlabeled data and perform supervised training. We take this, the **supervised-only strategy**, as our first baseline. To compare with established semisupervised machine learning techniques, we implement two more strategies: Bootstrapping and Π -models, representing

²We use Chang et al. (2022)’s PyTorch reimplementation obtain from <https://github.com/cmu-llab/middle-chinese-reconstruction>, with minor modifications to support multi-layer.

the **proxy-labelling** and **consistency regularization** approaches respectively.

Bootstrapping adds the model’s most confident predictions on unlabeled data to the train set as pseudo-labels (Lee, 2013). In our Bootstrapping setup, the model’s most confident (i.e. probable) protoform reconstructions for unlabeled cognate sets, filtered by a minimum confidence threshold and capped at a maximum number of top predictions per epoch, are added as pseudo-labels to the training set at the end of each epoch starting from a set number of warmup epochs (See Appendix C for Bootstrapping hyperparameters).

Π -model optimizes the model’s consistency by creating two stochastically augmented inputs from the same training example, feeding both of augmented inputs into the model, and minimizing the mean square difference between the two outputs (Laine and Aila, 2017). For continuous inputs, stochastic augmentation could be simple noise. Stochastic changes to phonemes, however, would defy protoform reconstruction’s goal of finding regular sound changes. Instead, we implement stochastic cognate set augmentation by randomly permuting the order of reflexes and randomly dropping daughter languages.

Observe that some of the above strategies can be combined: Bootstrapping can always be used on top of every other strategy, while our proposed DPD architecture can be merged with Π -model into a model that performs both reflex prediction and consistency regularization (See Appendix O for detail). In total, we test 8 strategies: supervised only (SUPV), Bootstrapping (BST), Π -model (IIM), Π -model with Bootstrapping (IIM-BST), DPD-BiReconstructor (DPD), DPD with Bootstrapping (DPD-BST), DPD merged with Π -model (DPD-IIM), and DPD-IIM with Bootstrapping (DPD-IIM-BST). Among the 8 strategies, we consider single baseline strategies (SUPV, BST, IIM) to be weak baselines, and the combination of non-DPD semisupervised techniques (IIM-BST) to be the strong baseline. Combined with 2 encoder-decoder architectures, GRU and Transformer, we experiment with 16 strategy-architecture combinations which we identify by an architecture prefix (GRU- or Trans-) followed by the strategy name.

2.3 Experiments

Datasets: We test both Romance and Sinitic languages, represented by the phonetic version³ of Meloni et al. (2021)’s Romance dataset (Rom-phon) for Latin reconstruction and Chang et al. (2022)’s WikiHan dataset for Middle Chinese reconstruction⁴. We simulate the semisupervised reconstruction scenario by hiding a random subset of protoform labels from the fully labeled train set. We refer to the percentage of labels retained after label removal as the **labeling setting**.

Cross-strategy comparisons: The primary interest of our experiments is how reconstruction performance differs between strategies with a fixed percentage of labels. We fix the labeled percentage at 10%, which entails approximately 516 and 870 labeled cognate sets for WikiHan and Rom-phon respectively. Due to randomness in semisupervised dataset generation, we repeat the experiment on four randomly generated semisupervised trained sets for each of WikiHan and Rom-phon, labeled group 1 to group 4. Each of the 16 strategy-architecture combinations is tested 10 times in each group. We then compare the performance of strategies within the same architecture.

Cross-labeling setting comparisons: We test all strategies on datasets with 5%, 10%, 20%, and 30% of labels to study the relationship between the labeling setting and the model’s performance⁵. With these labeling settings, the numbers of labeled cognate sets ranging from 181 to 1,084 for WikiHan and 304 to 1,821 for Rom-phon⁶. Randomly selecting labels for each of the labeling settings—especially on already small datasets—introduces variations on the learnable information contained in the labels and could introduce noise. To mitigate this, we enforce a monotonic subset selection constraint: given the complete train label set L and a semisupervised label set $L_{p_i} \subseteq L$ retaining p_i per cent of the labels, semisupervised sets L_{p_1}, L_{p_2} of increasing percentages ($p_1 \leq p_2$) must satisfy $L_{p_1} \subseteq L_{p_2}$ (See Figure 2). The 10% labeled

³We focus on the phonetic version due to a large number of semisupervised strategies and since DPD is primarily motivated by phonetic reconstruction.

⁴Rom-phon is not licenced for redistribution; WikiHan is licenced under Creative Commons CC0. For more dataset details, see Appendix A.

⁵It is often the case that, in manual reconstruction projects, the majority of sound changes can be discovered from a minority of the available cognate sets. We hope to capture this observation with our choice of labeling settings.

⁶See Table 6 for detail.

dataset we use in cross-labeling setting comparison corresponds to group 1 in cross-strategy comparison (See Appendix B for detail).

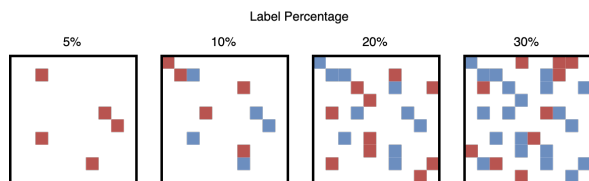


Figure 2: An illustration of the monotonic constraint for creating semisupervised datasets of varied labeling settings. In a hypothetical dataset with 100 cognate sets, represented as a 10×10 grid, shaded cells indicate cognate sets with associated labels (i.e. gold protoform). “■” indicates a protoform label not present in a previous subset. Observe that as the percentage of labels increases, no label is removed.

Hyperparameters: We use Bayesian search to tune hyperparameters for all 16 strategy-architecture combinations on a fixed 10% labeled semisupervised dataset⁷. For each strategy-architecture combination, for 100 iterations, we select the hyperparameter leading to the best validation phoneme edit distance. See Appendix C for details on hyperparameters.

Evaluation metrics: Evaluation metrics we use follow directly from supervised reconstruction literature, including token edit distance (TED), a count of the number of insertion, deletion, or substitution operations between the prediction and the target (Levenshtein, 1966); token error rate (TER), a length-normalized token edit distance (Cui et al., 2022); accuracy (ACC); feature error rate (FER), a measure of phonetic distance by Pan-Phon (Mortensen et al., 2016); and B-Cubed F Score (BCFS), a measure of structural similarity between the prediction and the target (Amigó et al., 2009; List, 2019).

Statistical tests: For each combination of dataset group, labeling setting, and architecture, we test all strategies against each other for differences. We use both the Wilcoxon Rank-Sum test (Wilcoxon, 1992) and Bootstrap test for mean difference (Sivaganesan, 1994) to test for significance in performance difference. We use a $\alpha = 0.01$ significance threshold and consider results to be significant if indicated by both tests.

⁷Due to a large number of strategy-architecture combinations, tuning the model on every labeling settings is costly. We fix the dataset generation seed to 0 for tuning.

3 Results

Cross-strategy comparisons: Table 2 shows the performance of all strategy-architecture combinations on four 10% labeled datasets. For both architectures on WikiHan, DPD-IIM-BST performs the best and significantly better than all baselines on all metrics. Transformer trained with DPD attains similar performance to DPD-IIM-BST, outperforming all baseline strategies. GRU trained with DPD performs similarly to IIM-BST, both of which perform better than weak baselines in a majority of situations. On Rom-phon, Transformer performed the best when trained with DPD-IIM-BST while GRU performed the best when trained with DPD-BST, both of which are significantly better than all baselines across all metrics.

Interestingly, while Kim et al. (2023) finds that a supervised Transformer model outperforms Meloni et al. (2021)’s GRU model on Rom-phon but not WikiHan, we observe the opposite in a 10% labeled semisupervised reconstruction setup, with Transformer outperforming GRU on WikiHan but not Rom-phon under most strategies. This appears to contradict Kim et al. (2023)’s hypothesis that a Transformer reconstruction model requires more data compared to an RNN.

Consistent with our hypothesis that access to different subsets of labels affects learning outcomes, we observe high variations in performance between dataset groups. Figures 6 and 7 visualize the performance of strategy-architecture combination by group, revealing that most strategies tend to do better on some dataset seeds.

Performance for varied labeling settings: Despite only being tuned with 10% of labels, DPD-based strategies generalized to other labeling settings, often outperforming strong baseline strategies on all metrics on at least one of the architectures (See Appendix F for detail). We see a non-linear scaling between performance and the percentage of labels and notice higher performance variations between strategies for lower percentages of labels. At a 5% labeling setting, for example, GRU-DPD-IIM-BST almost doubles the accuracy of GRU-SUPV on WikiHan. At a 30% labeling setting, some strategies attain accuracy close to existing fully supervised reconstruction models, with Trans-DPD-IIM 5.14 percentage points behind Meloni et al. (2021)’s supervised GRU on WikiHan (Chang et al., 2022) and GRU-DPD-BST 7.30 percentage points behind Kim et al. (2023)’s

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	40.50% $\text{\textcircled{8}}$	1.0075 $\text{\textcircled{8}}$	0.2360 $\text{\textcircled{8}}$	0.0970 $\text{\textcircled{8}}$	0.6707 $\text{\textcircled{8}}$
	DPD-BST (ours)	39.06% $\text{\textcircled{8}}$	1.0367 $\text{\textcircled{8}}$	0.2428 $\text{\textcircled{8}}$	0.0997 $\text{\textcircled{8}}$	0.6630 $\text{\textcircled{8}}$
	DPD-IIM (ours)	37.72% $\text{\textcircled{8}}$	1.0791 $\text{\textcircled{1e}}$	0.2528 $\text{\textcircled{1e}}$	0.1022 $\text{\textcircled{18}}$	0.6472 $\text{\textcircled{e}}$
	DPD (ours)	39.50% $\text{\textcircled{8}}$	1.0356 $\text{\textcircled{8}}$	0.2426 $\text{\textcircled{8}}$	0.0993 $\text{\textcircled{8}}$	0.6564 $\text{\textcircled{8}}$
	IIM-BST	34.21%	1.1489	0.2691	0.1106	0.6371
	BST (Lee, 2013)	34.78%	1.1455	0.2683	0.1109	0.6334
	IIM (Laine and Aila, 2017)	34.30%	1.1699	0.2740	0.1122	0.6209
	SUPV	33.25%	1.1891	0.2785	0.1140	0.6138
GRU	DPD-IIM-BST (ours)	39.74% $\text{\textcircled{8}}$	1.0280 $\text{\textcircled{8}}$	0.2408 $\text{\textcircled{8}}$	0.0972 $\text{\textcircled{8}}$	0.6683 $\text{\textcircled{8}}$
	DPD-BST (ours)	35.89% $\text{\textcircled{1e}}$	1.1025 $\text{\textcircled{18}}$	0.2582 $\text{\textcircled{18}}$	0.1039 $\text{\textcircled{18}}$	0.6493 $\text{\textcircled{18}}$
	DPD-IIM (ours)	37.90% $\text{\textcircled{8}}$	1.0697 $\text{\textcircled{18}}$	0.2506 $\text{\textcircled{18}}$	0.1006 $\text{\textcircled{18}}$	0.6517 $\text{\textcircled{18}}$
	DPD (ours)	34.51% $\text{\textcircled{18}}$	1.1538 $\text{\textcircled{18}}$	0.2703 $\text{\textcircled{18}}$	0.1091 $\text{\textcircled{34}}$	0.6278 $\text{\textcircled{18}}$
	IIM-BST	34.99% $\text{\textcircled{2}}$	1.1479 $\text{\textcircled{3}}$	0.2689 $\text{\textcircled{3}}$	0.1077 $\text{\textcircled{3}}$	0.6354 $\text{\textcircled{2}}$
	BST (Lee, 2013)	28.18%	1.3092	0.3067	0.1208	0.5939
	IIM (Laine and Aila, 2017)	32.59%	1.2047	0.2822	0.1137	0.6166
	SUPV	28.16%	1.3257	0.3105	0.1234	0.5835
Transformer	DPD-IIM-BST (ours)	34.63% $\text{\textcircled{8}}$	1.3115 $\text{\textcircled{8}}$	0.1463 $\text{\textcircled{8}}$	0.0588 $\text{\textcircled{8}}$	0.7850 $\text{\textcircled{8}}$
	DPD-BST (ours)	33.51% $\text{\textcircled{8}}$	1.3605 $\text{\textcircled{8}}$	0.1517 $\text{\textcircled{8}}$	0.0599 $\text{\textcircled{8}}$	0.7763 $\text{\textcircled{8}}$
	DPD-IIM (ours)	29.24%	1.5888	0.1772	0.0732	0.7423
	DPD (ours)	31.94% $\text{\textcircled{34}}$	1.5111	0.1685	0.0678 $\text{\textcircled{2}}$	0.7529
	IIM-BST	32.10% $\text{\textcircled{34}}$	1.4005 $\text{\textcircled{34}}$	0.1562 $\text{\textcircled{34}}$	0.0636 $\text{\textcircled{34}}$	0.7716 $\text{\textcircled{34}}$
	BST (Lee, 2013)	29.95%	1.5066	0.1680	0.0704	0.7555
	IIM (Laine and Aila, 2017)	26.97%	1.7134	0.1911	0.0796	0.7239
	SUPV	26.99%	1.7331	0.1933	0.0794	0.7218
GRU	DPD-IIM-BST (ours)	36.78% $\text{\textcircled{e}}$	1.2380 $\text{\textcircled{8}}$	0.1381 $\text{\textcircled{8}}$	0.0483 $\text{\textcircled{8}}$	0.7980 $\text{\textcircled{8}}$
	DPD-BST (ours)	37.60% $\text{\textcircled{8}}$	1.2149 $\text{\textcircled{8}}$	0.1355 $\text{\textcircled{8}}$	0.0457 $\text{\textcircled{8}}$	0.8014 $\text{\textcircled{8}}$
	DPD-IIM (ours)	31.51%	1.4892	0.1661	0.0628	0.7586
	DPD (ours)	31.12%	1.4837	0.1655	0.0608	0.7591
	IIM-BST	35.50%	1.2970	0.1447	0.0531	0.7909 $\text{\textcircled{1}}$
	BST (Lee, 2013)	35.87%	1.2893	0.1438	0.0509	0.7908
	IIM (Laine and Aila, 2017)	29.40%	1.5440	0.1722	0.0643	0.7517
	SUPV	30.69%	1.5018	0.1675	0.0612	0.7558

Table 2: Performance of all strategies on 10% labeled WikiHan (top) and Rom-phon (bottom) for each architecture, averaged across all runs in four groups (10 runs per strategy-architecture combination per group). Bold: best-performing strategy for the corresponding architecture and dataset; $\text{\textcircled{1}}$: significantly better than all weak baselines (SUPV, BST, and IIM) on group 1 with $p < 0.01$; $\text{\textcircled{e}}$: significantly better than the IIM-BST strong baseline and all weak baselines on group 1 with $p < 0.01$; $\text{\textcircled{2}}$, $\text{\textcircled{3}}$, $\text{\textcircled{4}}$, $\text{\textcircled{2}}$, $\text{\textcircled{3}}$, $\text{\textcircled{4}}$: likewise for groups 2, 3, and 4.

supervised Transformer on Rom-phon (See Table 15 and Appendix J for detail).

4 Analysis

4.1 DPD Training

We investigate the interaction between D2P and P2D during training. Figure 3 shows the train and validation accuracy trajectories of reflexes reconstructed from the intermediate representation of protoform reconstructions. We perform a CRINGE loss ablation by repeating the same run but with CRINGE loss disabled (same hyperpa-

rameter and model parameter initialization). In both cases, reflexes are more accurately predicted from correct protoform reconstructions, corroborating our motivation that better reflex reconstruction should promote better protoform reconstruction. CRINGE loss appears to lead to slightly lower reflex accuracy given incorrect reconstruction at earlier epochs, albeit at the expense of slightly lower reflex accuracy given correct reconstruction. In practice, we find the CRINGE loss weight to be a relatively insignificant hyperparameter.

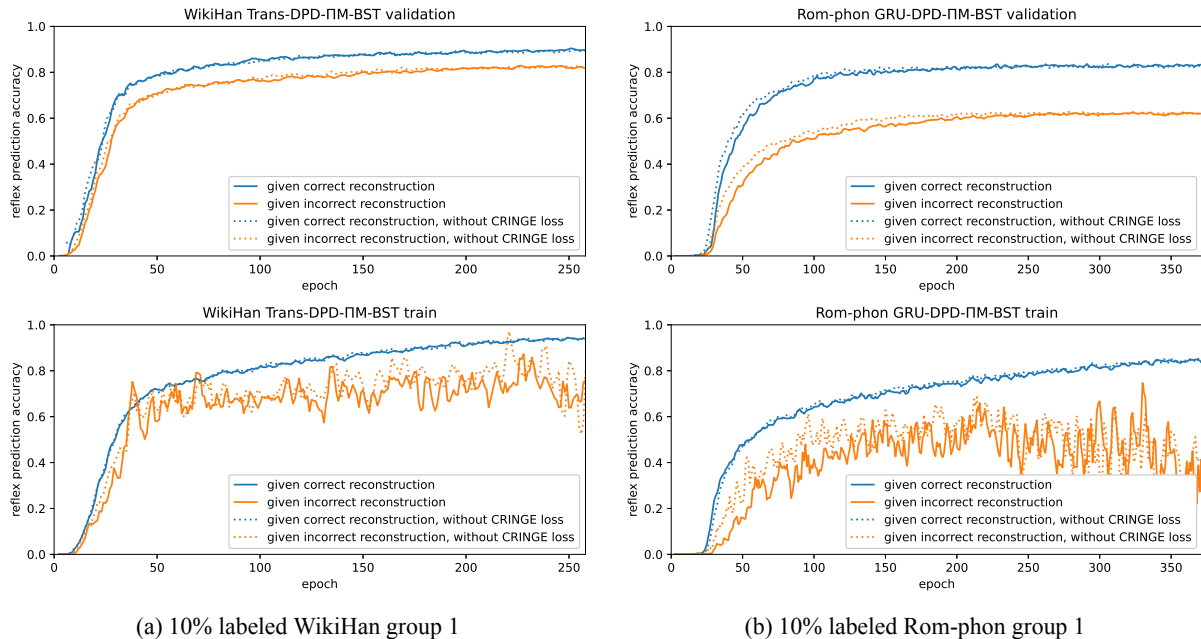


Figure 3: Validation (top) and train (bottom) reflex reconstruction accuracy given correct versus incorrect protoform prediction during training for a randomly selected run within the most accurate strategy. A rolling average of window size 3 is used for smoothing.

4.2 Reflex Prediction Performance

While the DPD architecture is designed for reflex prediction-assisted reconstruction, we observe good reflex performance in some situations. Table 3 shows the strategy-architecture combinations with the most accurate P2D sub-network when evaluated on gold protoforms⁸. Compared to reconstruction, we obtain semisupervised reflex prediction performance that are much closer to supervised performance, even with small percentages of labels—consistent with the assumption that sound changes in the proto-to-daughters direction are easier to model.

It is worth noting, however, that reflex prediction performance based on gold protoforms depends on the weight of the corresponding loss. In fact, some of the best DPD models perform poorly

⁸It is arguably more interesting to evaluate reflex prediction performance based on the model’s latent protoform representation. Unfortunately, obtaining latent representations of correct protoforms is not always possible.

when the P2D sub-network is evaluated on gold protoforms. We conclude that P2D’s ability to assist D2P during training is not contingent on P2D’s performance on gold protoforms as discrete input.

4.3 Learned Phonetic Representations

Inspired by Meloni et al. (2021), we probe the model’s learned embeddings for a hierarchical organization of phonemes using sklearn’s Ward variance minimization algorithm (Ward, 1963). Figure 4 shows the results for two selected daughter languages on the most accurate model in each dataset (GRU-DPD-BST for Rom-photon and Trans-DPD-IIM-BST for WikiHan) and the best model from their non-DPD counterpart (GRU-BST and Trans-IIM-BST respectively).

For French, phoneme embeddings trained with DPD-BST reveals a clear division between vowels and consonants similar to Meloni et al. (2021)’s supervised reconstruction model. Except for [ø], nasal vowels are grouped together. Specific

		5%	10%	20%	30%	100%
WikiHan	Top performer	GRU-DPD-IIM-BST	GRU-DPD-IIM-BST	Trans-DPD-IIM-BST	GRU-DPD-BST	-
	ACC	45.37%	56.74%	61.83%	62.02%	66.43%
Rom-photon	Top performer	Trans-DPD-IIM-BST	Trans-DPD-IIM-BST	Trans-DPD-BST	Trans-DPD-BST	-
	ACC	53.51%	57.95%	60.08%	61.13%	63.85%

Table 3: Strategy-architecture combinations with the highest reflex prediction accuracy in group 1 for each labeling setting, along with reference supervised (100% labeled) reflex prediction accuracy.

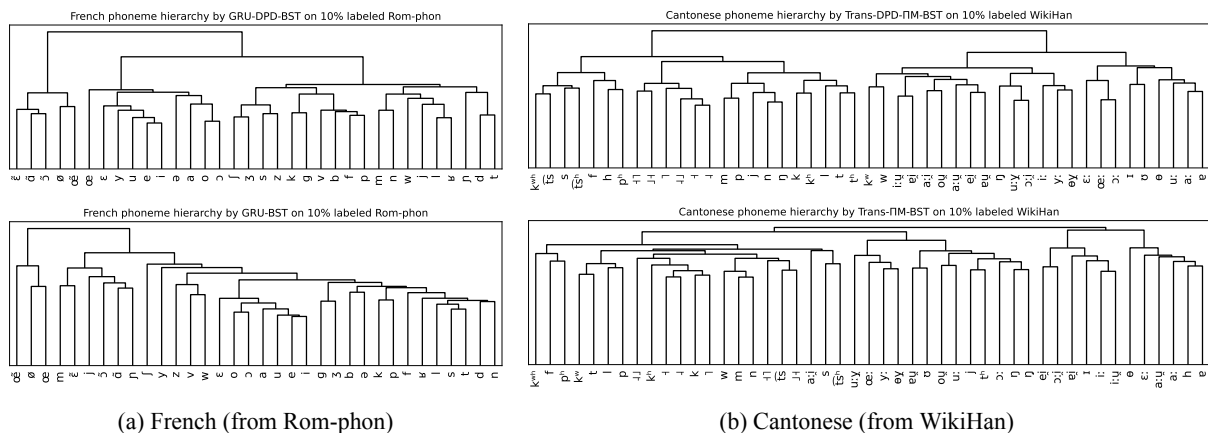


Figure 4: Hierarchical clustering revealing learned phoneme organization for two selected daughter languages obtained from the best run (within group 1 of 10% labeling setting) in the best DPD-based strategy-architecture combination (top) and the best run from their non-DPD counterpart (bottom).

phoneme pairs with minimal difference in features are also placed together, such as the alveolar fricatives [s] and [z], post-alveolar fricatives [ʃ] and [ʒ], and velar plosives [k] and [g], and all of which only differ in voicing. In the embeddings trained with BST, some but not all vowels are clustered together, and sister groups (i.e. immediate relative in the tree) are less interpretable, such as [m] and [ɛ̃] as well as [b] and [ə]. For Cantonese, we see a similar pattern where vowels and consonants have a clearer division when trained with DPD-IIM-BST. Additionally, tones are organized into the same cluster by DPD-IIM-BST but not by IIM-BST. In Appendix N, we extend our probe to set of all phonemes present in daughter languages⁹ and find that the above observations generalize to phoneme organization beyond a single language, which Meloni et al. (2021) does not consider as part of their analysis.

We conclude from phonetic probing that DPD-based strategies are better at capturing linguistically meaningful representations of phonemes. It is possible that DPD-based strategies need good phonetic representations to perform well on multiple phonology-intensive tasks, which could in turn better inform protoform reconstruction.

4.4 Ablation on Unlabeled Data

To study whether the performance gains of semisupervised strategies are because of their effective use of unlabeled cognate sets, we perform ablation experiments at a 10% labeling setting but with all

⁹Phonemes present only in the protolanguage are not included because non-DPD do not update their embeddings.

unlabeled training data removed¹⁰, effectively creating a small supervised training set¹¹.

We find that, in the absence of unlabeled data, IIM, DPD, and DPD-IIM can sometimes perform significantly better than SUPV, but almost always perform significantly worse than when unlabeled data is used or when unlabeled data is used in conjunction with Bootstrapping (see Tables 13 and 14). This seems to suggest that IIM, DPD, and DPD-IIM learn effectively from both labeled and unlabeled data. It is possible that, on labeled data, the P2D sub-network in DPD can still inform the D2P sub-subnetwork, and the stochastic data augmentation in our implementation of II-model can augment labeled training examples.

4.5 Applicability in Supervised Reconstruction

Seeing some indication that the semisupervised reconstruction strategies are applicable for supervised reconstruction on a small subset of the training set (see Section 4.4), we test whether their advantages generalize to supervised reconstruction on the full training set¹².

Table 15 compares the supervised reconstruction performance of IIM, DPD, and DPD-IIM with existing supervised reconstruction methods, including Meloni et al. (2021)’s GRU model, Kim et al. (2023)’s Transformer model, and Lu et al. (2024)’s state-of-the-art reranked reconstruction systems. We find that, with the right architecture,

¹⁰We exclude strategies with Bootstrapping because it has no effect when there is no unlabeled data.

¹¹See Appendix I for experimental details.

¹²See Appendix J for experimental details.

IIM, DPD, and DPD-IIM can often outperform SUPV. On average, Trans-DPD-IIM performs the best on WikiHan for all metrics, and GRU-DPD performs the best on Rom-phon for all metrics except FER. On WikiHan, Trans-DPD-IIM performs significantly better than Lu et al. (2024) on FER only. On Rom-phon, GRU-DPD and GRU-DPD-IIM both perform significantly better than Lu et al. (2024) on ACC only. We conclude that, despite being motivated by semisupervised reconstruction, IIM and DPD could be useful for supervised reconstruction. We leave it to future work to understand the role of data augmentation and the DPD architecture in supervised settings.

5 Related Work

Computational Historical Linguistics: Protoform reconstruction and reflex prediction are two central tasks in computational historical linguistics. Protoform reconstruction predicts the protoform given reflexes in a cognate set, while reflex prediction models the changes from the protoform to its reflexes. Computational reconstruction and reflex prediction methods vary and include rule-based systems (Heeringa and Joseph, 2007; Marr and Mortensen, 2020, 2023), probabilistic models operating on phylogenetic trees (Bouchard-Côté et al., 2007b,a, 2009, 2013), automated alignment systems (Bodt and List, 2022; List et al., 2022), and more recently, neural networks (Fourrier, 2022).

Supervised Reconstruction and reflex prediction: Ciobanu and Dinu (2018) and Ciobanu et al. (2020) formulate protoform reconstruction as a sequence labelling task and use conditional random field to reconstruct protoform phonemes at each position in the daughter sequences. Meloni et al. (2021) reformulates the reconstruction task as a sequence-to-sequence task by concatenating reflexes into a single input sequence separated by language tags and uses GRU to reconstruct Latin on a new Romance dataset. A group of subsequent researchers refined this task with additional datasets (Chang et al., 2022) and improved neural methods (Kim et al., 2023; Akavarapu and Bhatlacharya, 2023; Lu et al., 2024).

Reflex prediction can be viewed as sequence-to-sequence transduction in the reverse direction. Cathcart and Rama (2020) propose an LSTM encoder-decoder model to infer Indo-Aryan reflexes from Old Indo-Aryan, aided by semantic

embedding. Arora et al. (2023) replicates Cathcart and Rama (2020)’s experiments on South Asia languages with both GRU and Transformer models. As reflex prediction maps from protolanguage to multiple daughter languages, a prompting token is often attached to the input to specify the target daughter language.

Non-Supervised Computational Historical Linguistics: Past work in which comparative reconstruction was performed without full supervision included Bouchard-Côté et al. (2009) and He et al. (2023). To the best of our knowledge, no research focuses specifically on semisupervised neural reconstruction.

Semisupervised Learning: Effective semisupervised learning should utilize unlabeled training data. One approach is proxy-labelling, by which synthesized labels are added to unlabeled training examples via heuristics (Ouali et al., 2020). Another approach is consistency regularisation, largely based on the smoothness assumption that is applicable regardless of whether labels are present: similar input data in high-density regions should have similar labels, whereas input data separated by low-density regions should not (Tsai and Lin, 2019; Ouali et al., 2020; Luo et al., 2018).

Our task is related to semisupervised machine translation (Edunov et al., 2018; Sennrich et al., 2016; Skorokhodov et al., 2018; Gülçehre et al., 2015; Cheng et al., 2016). However, it differs crucially in that the model has no access to monolingual text, only labeled and unlabeled cognate sets.

6 Conclusion

We introduce the task of semisupervised reconstruction, marking a step forward toward a practical computational reconstruction system that can assist early-stage protolanguage reconstruction projects. We design the DPD-BiReconstructor architecture to implement historical linguists’ comparative method, yielding performance that surpasses existing sequence-to-sequence reconstruction models and established semisupervised learning techniques, especially when protoform labels are scarce.

Limitations

Due to a large number of possible strategies, we have limited our focus of semisupervised reconstruction experiments with DPD to 10% labeled

WikiHan and Rom-phon datasets. It is left to future work to expand the research on semisupervised reconstruction to other datasets.

Though DPD has a clear motivation and demonstrates superior empirical performance, interpretations of what sound changes DPD learns from the unlabeled cognate sets that enable its better reconstruction performance are less clear. Our theory is that neural networks within the system are sufficiently expressive to learn better reconstruction in a bidirectional manner, but we have not yet obtained evidence that the model’s reasoning matches that of a linguist beyond having a better representation of phonemes and taking the step to derive the reflexes from the reconstruction. Nevertheless, we demonstrate that inferring reflexes from reconstructions is not just a powerful methodology for historical linguistics, but also for computational historical linguistics.

Although we observe less accurate reflexes being predicted from incorrect protoforms compared to correct protoforms, reflexes derived from incorrect protoforms are still highly accurate. Future work could explore ways to mitigate this issue and improve the reflex prediction sub-network’s ability to discriminate between correct and incorrect protoforms.

It is also notable that our implementation of the Π -Model included only reflex permutation and daughter deletion as noising strategies. It is possible that other possible strategies may have strengthened this baseline.

The protoform reconstruction task is far from solved—with humans’ success at the reconstruction of ancient languages, a truly intelligent reconstruction system should in the future be able to perform reconstruction without the help of labels.

Ethics

Historical reconstruction involves very limited risks to humans. The risks that do exist are both individual and political. On the one hand, we cannot guarantee that all of the data used in this study were collected in an ethical fashion. However, we did not do any data collection and relied upon existing resources which, to the best of our knowledge, were collected under standard scholarly and scientific protocols. On the other hand, the results of historical reconstruction can be politically fraught. For example, historical reconstructions can be used to show, in some cases, that linguis-

tic boundaries between people groups do not align with cultural and political loyalties. This can be disruptive and may even be associated with violence. Because our work does not concentrate directly on phylogeny—the main source of political complications in comparative reconstruction—we believe that the risks are minimal.

Acknowledgements

This work is supported by Carnegie Mellon University’s SURF grant. We would like to thank Chenxuan Cui, Calvin Chang, and Graham Neubig for helpful ideas and discussions. We are grateful to our anonymous reviewers for many comments that helped improve the writing and informed additional experiments and analyses.

References

- Leonard Adolphs, Tianyu Gao, Jing Xu, Kurt Shuster, Sainbayar Sukhbaatar, and Jason Weston. 2022. [The CRINGE Loss: Learning what language not to model](#).
- V. S. D. S. Mahesh Akavarapu and Arnab Bhattacharya. 2023. [Cognate Transformer for Automated Phonological Reconstruction and Cognate Reflex Prediction](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6852–6862.
- Enrique Amigó, Julio Gonzalo, Javier Artiles, and Felisa Verdejo. 2009. A comparison of extrinsic clustering evaluation metrics based on formal constraints. *Information retrieval*, 12:461–486.
- Aryaman Arora, Adam Farris, Samopriya Basu, and Suresh Kolichala. 2023. [Jambu: A historical linguistic database for South Asian languages](#). In *Proceedings of the 20th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 68–77, Toronto, Canada. Association for Computational Linguistics.
- Siddhant Arora, Siddharth Dalmia, Brian Yan, Florian Metze, Alan W. Black, and Shinji Watanabe. 2022. [Token-level Sequence Labeling for Spoken Language Understanding using Compositional End-to-End Models](#).
- Lukas Biewald. 2020. Experiment tracking with weights and biases.
- Timotheus A. Bodt and Johann-Mattis List. 2022. [Reflex prediction: A case study of Western Kho-Bwa](#). *Diachronica*, 39(1):1–38.
- Alexandre Bouchard-Côté, Thomas L. Griffiths, and Dan Klein. 2009. Improved Reconstruction of Proto-language Word Forms. In *Proceedings of Human*

- Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 65–73, Boulder, Colorado. Association for Computational Linguistics.
- Alexandre Bouchard-Côté, David Hall, Thomas L. Griffiths, and Dan Klein. 2013. [Automated reconstruction of ancient languages using probabilistic models of sound change](#). *Proceedings of the National Academy of Sciences*, 110(11):4224–4229.
- Alexandre Bouchard-Côté, Percy Liang, Thomas Griffiths, and Dan Klein. 2007a. A Probabilistic Approach to Diachronic Phonology. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 887–896, Prague, Czech Republic. Association for Computational Linguistics.
- Alexandre Bouchard-Côté, Percy S Liang, Dan Klein, and Thomas Griffiths. 2007b. A Probabilistic Approach to Language Change. In *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc.
- Karl Brugmann and Hermann Osthoff. 1878. *Morphologische Untersuchungen auf dem Gebiete der indogermanischen Sprachen*, volume 1. Hirzel.
- L. Campbell. 2021. *Historical Linguistics: An Introduction*. Edinburgh University Press.
- Chundra Cathcart and Taraka Rama. 2020. [Disentangling dialects: A neural approach to Indo-Aryan historical phonology and subgrouping](#). In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 620–630, Online. Association for Computational Linguistics.
- Kalvin Chang, Chenxuan Cui, Youngmin Kim, and David R. Mortensen. 2022. WikiHan: A New Comparative Dataset for Chinese Languages. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3563–3569, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Yong Cheng, Wei Xu, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. [Semi-Supervised Learning for Neural Machine Translation](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1965–1974, Berlin, Germany. Association for Computational Linguistics.
- Alina Maria Ciobanu and Liviu P. Dinu. 2018. Ab Initio: Automatic Latin Proto-word Reconstruction. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1604–1614, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Alina Maria Ciobanu, Liviu P. Dinu, and Laurentiu Zoicas. 2020. Automatic Reconstruction of Missing Romanian Cognates and Unattested Latin Words. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3226–3231, Marseille, France. European Language Resources Association.
- Chenxuan Cui, Ying Chen, Qinxin Wang, and David R. Mortensen. 2022. Neural Proto-Language Reconstruction. Technical report, Carnegie Mellon University, Pittsburgh, PA.
- Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. [Understanding Back-Translation at Scale](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 489–500, Brussels, Belgium. Association for Computational Linguistics.
- Clémentine Fourrier. 2022. *Neural Approaches to Historical Word Reconstruction*. Ph.D. thesis, Université PSL (Paris Sciences & Lettres).
- Clémentine Fourrier and Benoît Sagot. 2022. Probing multilingual cognate prediction models. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3786–3801.
- Çağlar Gülçehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loïc Barrault, Hui-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2015. On Using Monolingual Corpora in Neural Machine Translation. *ArXiv*.
- Andre He, Nicholas Tomlin, and Dan Klein. 2023. [Neural Unsupervised Reconstruction of Protolanguage Word Forms](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1636–1649, Toronto, Canada. Association for Computational Linguistics.
- Wilbert Heeringa and Brian Joseph. 2007. The Relative Divergence of Dutch Dialect Pronunciations from their Common Source: An Exploratory Study. In *Proceedings of Ninth Meeting of the ACL Special Interest Group in Computational Morphology and Phonology*, pages 31–39, Prague, Czech Republic. Association for Computational Linguistics.
- Young Min Kim, Calvin Chang, Chenxuan Cui, and David R. Mortensen. 2023. Transformed Proto-form Reconstruction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 24–38, Toronto, Canada. Association for Computational Linguistics.
- Samuli Laine and Timo Aila. 2017. [Temporal Ensembling for Semi-Supervised Learning](#).
- Dong-Hyun Lee. 2013. Pseudo-Label : The Simple and Efficient Semi-Supervised Learning Method for Deep Neural Networks.

- Vladimir I. Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet Physics Doklady*, volume 10, pages 707–710. Soviet Union.
- Johann-Mattis List. 2019. [Beyond edit distances: Comparing linguistic reconstruction systems](#). *Theoretical Linguistics*, 45(3-4):247–258.
- Johann-Mattis List and Robert Forkel. 2021. [LingPy. A Python Library for Historical Linguistics](#). Zenodo.
- Johann-Mattis List, Robert Forkel, and Nathan Hill. 2022. [A New Framework for Fast Automated Phonological Reconstruction Using Trimmed Alignments and Sound Correspondence Patterns](#). In *Proceedings of the 3rd Workshop on Computational Approaches to Historical Language Change*, pages 89–96, Dublin, Ireland. Association for Computational Linguistics.
- Liang Lu, Jingzhi Wang, and David R. Mortensen. 2024. [Improved Neural Protoform Reconstruction via Reflex Prediction](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 8683–8707, Torino, Italia. ELRA and ICCL.
- Yucen Luo, Jun Zhu, Mengxi Li, Yong Ren, and Bo Zhang. 2018. [Smooth Neighbors on Teacher Graphs for Semi-Supervised Learning](#). 2018 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8896–8905.
- Clayton Marr and David Mortensen. 2023. [Large-scale computerized forward reconstruction yields new perspectives in French diachronic phonology](#). *Diachronica*, 40(2):238–285.
- Clayton Marr and David R. Mortensen. 2020. [Computerized Forward Reconstruction for Analysis in Diachronic Phonology, and Latin to French Reflex Prediction](#). In *Proceedings of LT4HALA 2020 - 1st Workshop on Language Technologies for Historical and Ancient Languages*, pages 28–36, Marseille, France. European Language Resources Association (ELRA).
- Carlo Meloni, Shauli Ravfogel, and Yoav Goldberg. 2021. [Ab Antiquo: Neural Proto-language Reconstruction](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4460–4473, Online. Association for Computational Linguistics.
- David R. Mortensen, Patrick Littell, Akash Bharadwaj, Kartik Goyal, Chris Dyer, and Lori Levin. 2016. [PanPhon: A Resource for Mapping IPA Segments to Articulatory Feature Vectors](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3475–3484, Osaka, Japan. The COLING 2016 Organizing Committee.
- Yassine Ouali, Céline Hudelot, and Myriam Tami. 2020. [An Overview of Deep Semi-Supervised Learning](#).
- Michael Saxon, Samridhi Choudhary, Joseph P. McKenna, and Athanasios Mouchtaris. 2021. [End-to-End Spoken Language Understanding for Generalized Voice Assistants](#). In *Interspeech 2021*, pages 4738–4742.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Improving Neural Machine Translation Models with Monolingual Data](#).
- Siva Sivaganesan. 1994. [An Introduction to the Bootstrap \(Bradley Efron and Robert J. Tibshirani\)](#). *SIAM Review*, 36(4):677–678.
- Ivan Skorokhodov, Anton Rykachevskiy, Dmitry Emelyanenko, Sergey Slotin, and Anton Ponkratov. 2018. [Semi-Supervised Neural Machine Translation with Language Models](#). In *Proceedings of the AMTA 2018 Workshop on Technologies for MT of Low Resource Languages (LoResMT 2018)*, pages 37–44, Boston, MA. Association for Machine Translation in the Americas.
- Oswald J. L. Szemerényi. 1996. *Introduction to Indo-European Linguistics*. Oxford University Press UK.
- Kuen-Han Tsai and Hsuan-Tien Lin. 2019. [Learning from Label Proportions with Consistency Regularization](#). *ArXiv*.
- Joe H. Ward, Jr. 1963. [Hierarchical Grouping to Optimize an Objective Function](#). *Journal of the American Statistical Association*, 58(301):236–244.
- Frank Wilcoxon. 1992. [Individual Comparisons by Ranking Methods](#). *Breakthroughs in Statistics*, pages 196–202.

A Dataset Details

Both WikiHan and Rom-phon are split by 70%, 10%, and 20% into train, validation, and test sets. We remove labels from the supervised train set to create semisupervised train sets. The validation and test sets are left unmodified. The splits for WikiHan (Chang et al., 2022) match the original work. The splits and preprocessing for Meloni et al. (2021) and matches Kim et al. (2023). Table 4 shows the number of cognate sets in each split. WikiHan includes 8 daughter languages: Cantonese, Gan, Hakka, Jin, Mandarin, Hokkien, Wu, and Xiang. Rom-phon includes 5 daughter languages: French, Italian, Spanish, Romanian, and Portuguese.

B Dataset Groups

Semisupervised reconstruction datasets are generated pseudo-randomly based on a seed, which is

	WikiHan	Rom-phon
Train	3,615	6,071
Validation	517	878
Test	1,033	1,754
Total	5,165	8,703

Table 4: Number of cognate sets in the train, validation, and test split of both datasets.

itself chosen randomly. That is, a dataset seed deterministically generates a semisupervised dataset. To generate a semisupervised dataset, we initialize PyTorch’s pseudo-random number generator with the dataset seed, assign a uniformly distributed number between 0 and 1 to each training example with `torch.rand` in the same order as they appear in the dataset, and keep the protoform label on cognate sets whose assigned number is above a threshold such that the desired percentage of labels remain. The dataset seed is independent of seeds used to initialize model parameters in the experiments. Table 5 details the dataset seed used to select the subset of training labels to include for each experiment setup, and Table 6 details the number of labels at each labeling setting. We use group 1 for comparisons between labeling settings, with the same dataset seed ensuring label sets are monotonic subsets with respect to labeling settings with increasing percentages of labels. For the 10% labeled setting, four distinct dataset seeds simulate variations in the training data.

C Hyperparameters

We tune hyperparameters using Bayesian search on WandB (Biewald, 2020) with 100 runs for each strategy-architecture combination. We validate the model every 3 epochs and use early stopping if no improvement is made after 24 epochs. The dataset used for tuning semisupervised models is a 10% labeled train set generated using 0 as the dataset seed, making it different from the semisupervised datasets used in the experiments. Hyperparameters for semisupervised models can be found in Tables 16 and 17. Additional hyperparameter tunings for analysis purposes (Sections I and J) follow the same procedure but remove unlabeled training data or use a different labeling setting.

D Training

Parameter counts for the models can be found in Table 22. Hyperparameter tuning and experiments

are performed on a mix of NVIDIA GeForce GTX 1080 Ti, NVIDIA GeForce RTX 2080 Ti, NVIDIA RTX 6000 Ada Generation, NVIDIA RTX A6000, Quadro RTX 8000, and Tesla V100-SXM2-32GB GPUs for a total of 411 GPU days.

E Package Usage

Our model is implemented in PyTorch (See the code for details at <https://github.com/cmu-llab/dpd>). Sequence alignment is done using `lingpy` (List and Forkel, 2021) with default parameters. Hierarchical clustering is done using `AgglomerativeClustering` from `sklearn` with `Ward` linkage and distance threshold set to 0. Bootstrap tests are done using `scipy` with random state set to 0, a 99% confidence interval, and 9,999 resamples (default). Wilcoxon Rank-Sum tests are done using `scipy` at default parameter. Hierarchical clustering is visualized using `scipy`. Plots are created using `Matplotlib`.

F Additional Performance Data

Tables 8 and 9 show the performance of each strategy on group 1 for the 5%, 20%, and 30% labeling settings, along with indicators of statistical significance.

G Transductive Evaluation

Evaluation for semisupervised learning can be categorized as transductive and inductive (Ouali et al., 2020). Inductive evaluation tests the model’s performance on unseen data in a test set, whereas transductive evaluation tests the model’s ability to predict labels for the unlabeled data in the train set. In the context of semisupervised protoform reconstruction, transductive evaluation corresponds to the early stage of a real-world reconstruction project where cognate sets in higher abundance than known protoforms. In this section, we report the transductive performance of all the strategies.

We find that transductive performance differences between strategies are largely similar to inductive performance on the test set. Significance as to whether our DPD strategies perform better than baseline is also similar to that of inductive evaluation on the test set. For detailed transductive performance and their statistical significance, see Tables 10 (10% labeled WikiHan and Rom-phon), 11 (5%, 20%, and 30% labeled WikiHan), and 12 (5%, 20%, and 30% labeled Rom-phon).

		5%	10%	20%	30%
WikiHan	Group 1	2706283079	2706283079	2706283079	2706283079
	Group 2		2188396888		
	Group 3		2718489156		
	Group 4		1416758132		
Rom-phon	Group 1	1893608599	1893608599	1893608599	1893608599
	Group 2		1517201602		
	Group 3		2341117665		
	Group 4		3045950670		

Table 5: Dataset seeds used to create semisupervised train sets for each group and labeling setting. Using the same seed on Group 1 guarantees the monotonic increasing subset selection constraint for comparison between different labeling settings.

	WikiHan	Rom-phon
5%	181	304
10%	362	607
20%	723	1,214
30%	1,084	1,821
100%	3,615	6,071

Table 6: 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%).

We observed no clear pattern as to which strategies perform better on transductive evaluation compared to inductive evaluation. On 10% labeled WikiHan, average transductive accuracies are all 0.0-1.14% higher than their test accuracies. On 10% labeled Rom-phon, transductive accuracies for all strategies are 1.20-1.90% worse than their test accuracies. Given that differences between test and transductive performance are relatively consistent across strategies, including for supervised strategies in which the unlabeled portion of the train set effectively acts as another test set, our hypothesis is that the fixed train-test split played a role in the evaluation.

H Aligned Error Analysis

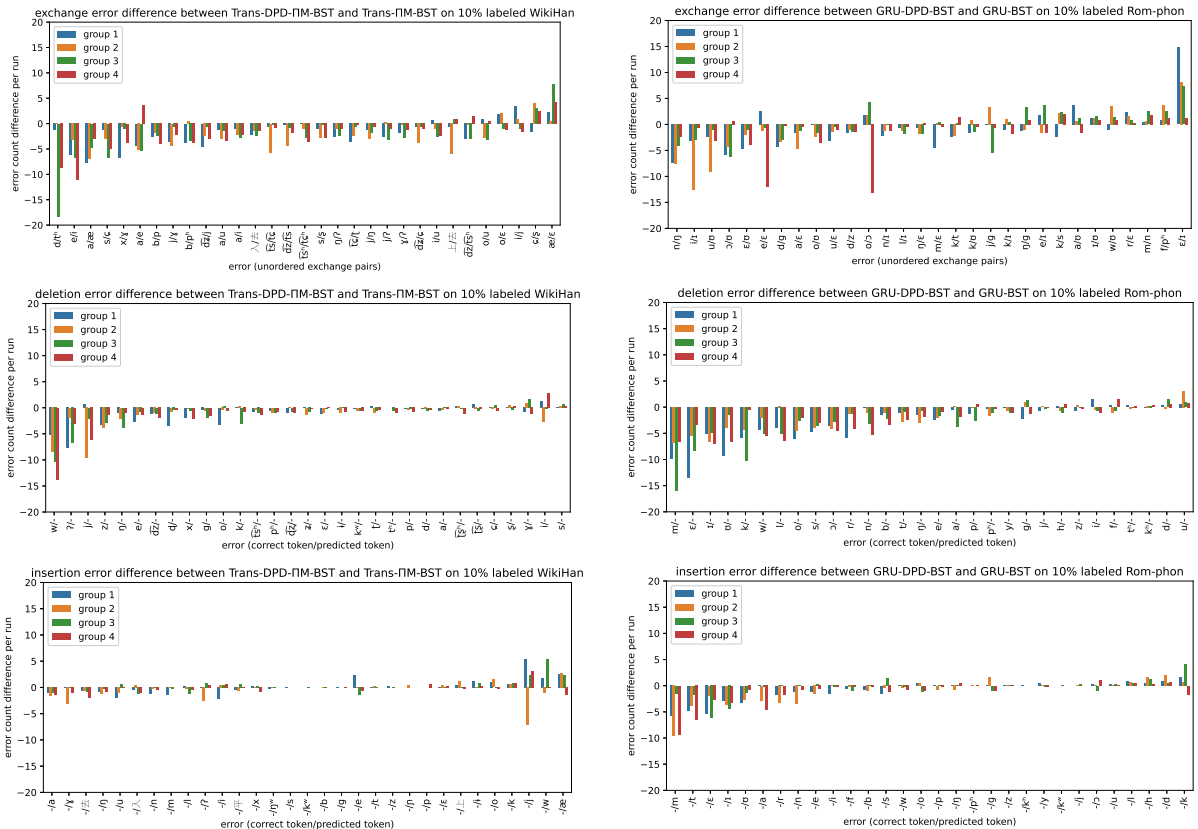
As an additional analysis, we align the protoform predictions and their targets on the test set using *lingpy* (List and Forkel, 2021) and identify errors made by reconstruction models. Consistent with Meloni et al. (2021)’s (supervised reconstruction) error analysis on Rom-phon, tense-lax errors occur most frequently for both the best DPD strategy and its non-DPD counterpart, with [i]/[ɪ], [e]/[ɛ], [o]/[ɔ], and [u]/[ʊ] being the top-four exchange er-

rors¹³. The top errors in WikiHan include inserting or deleting [j] and [w], vowel height exchange errors between pairs such as [e]/[i], [o]/[u], and [æ]/[a], tone errors, along with [a]/[o] errors. Exchange errors make up 68% and 63% of all errors for WikiHan and Rom-phon respectively.

We compare the average number of exchange errors between the best DPD-based strategy and their non-DPD counterpart. Figure 5 shows the most prominent exchange and non-exchange (i.e. insertion or deletion) error differences between the strategies. On average, GRU-DPD-BST makes 130 fewer mistakes than GRU-BST among 1731 test entries on Rom-phon, with a majority (72%) being insertion or deletion error reductions. TransDPD-PIM-BST makes 144 fewer mistakes than Trans-PIM-BST on WikiHan among 1033 test entries, with a majority (74%) being exchange error reductions. Some error differences between DPD and its non-DPD counterpart appear highly dependent on the dataset, such as GRU-DPD-BST making on average about 13 less [o]/[ɔ] errors on group 4 but more such errors on all other groups. Differences between groups could indicate that the distribution of additional unlabeled cognate sets in the training data plays an important role in the error patterns of DPD-BST strategies.

On Rom-phon, DPD is better at deciding whether to insert [m] in the reconstruction. This is the exact type of problem DPD is designed to handle: in situations where it is not apparent from the reflexes whether a phoneme absent in the reflexes should be added, it is often the case that one of the decisions will lead to lost information in the reconstruction. In DPD, the reflex prediction sub-network should be able to detect such information

¹³The [i]/[ɪ] exchange error, for example, refers to [i] being predicted in place of [ɪ] or vice versa.



(a) Error count differences between Trans-DPD-IIM-BST and Trans-IIM-BST on WikiHan.

(b) Error count differences between GRU-DPD-BST and GRU-BST on Rom-phon.

Figure 5: Most notable error count differences per run between the best-performing strategy-architecture combination and their non-DPD counterpart, averaged across 10 runs for each group on 10% labeled datasets. A negative count difference indicates fewer mistakes made by a DPD-based strategy. Exchange error pairs (top) are in no specific order. insertion or deletion errors (middle and bottom) are ordered and the x-labels indicate the target phoneme followed by the predicted phoneme, separated by a slash.

loss, which would otherwise lead to reflexes not being inferable from the reconstruction.

We compare errors made by the best GRU-DPD-BST against GRU-BST (within group 1 at 10% labeling setting). Among test examples where only one of GRU-DPD-BST and GRU-BST is correct, we see 12¹⁴ instances where GRU-DPD-BST produces the correct reconstruction but GRU-BST makes a mistake on [m] insertion or deletion. The opposite is true in only 2 instances¹⁵. Table 7 shows three test examples where GRU-DPD-BST performs better at deciding whether to add an [m] in the reconstruction as well as one example where it fails.

We hypothesize that since DPD makes better use of unlabeled data, an error is less likely to occur if there is an abundance of unlabeled examples involving information about the underlying

sound change pattern governing the context of a possible error. However, without explicit rules being learned by the models, it is difficult to assess what additional information is available to DPD in the unlabeled portion of the dataset and how the model learns from it. In an attempt to test this hypothesis, we estimate the abundance of learning resources in the unlabeled dataset by counting the number of examples involving the same sound correspondences (using lingpy Meloni et al. (2021)’s multi-sequence alignment) as a heuristic. We find no indication that count differences in exchange, insertion, or deletion errors are correlated with a higher abundance of training examples with the same sound correspondence in the unlabeled portion of the train set. Previous work has observed that statistical baselines can perform better than neural models at identifying the sound correct correspondences Fourier and Sagot (2022), implying that sound correspondence is not necessarily a

¹⁴excluding two false positives by lingpy

¹⁵excluding one false positive by lingpy

	konsentire	populare	immaterialem
French	k ð - s ã - - t i κ - -	p o p y l e κ - -	i m m a t e κ j e l - -
Italian	k o n s e - n t i r e -	p o p o l a r e -	i m m a t e r i a l e -
Romanian	k o n s i - m t - s i -	p o p u l a r - -	i - m a t e r j a l - -
Spanish	k o n s e - n t i r - -	p o p u l a r - -	i n m a t e r j a l - -
Portuguese	k u ŋ s e i ŋ t i r - -	p u p u l a r - -	i - m a t e r i a l - -
Latin	k o n s e - n t i r e -	p o p o l a r e m	i m m a t e r i a l e m
GRU-DPD-BST prediction	k o n s e - n t i r e -	p o p o l a r e m	i m m a t e r i a l e m
GRU-BST prediction	k o n s e - m t i r e m	p o p o l a r e -	i - m a t e r i a l e m

	traditor
French	t κ e - - - t - κ - -
Italian	t r a - d i t o r e -
Romanian	t r ə - d ə t o r -
Spanish	t r a i ð - - o r - -
Portuguese	t r a i d - - o r - -
Latin	t r a - d i t o r e m
GRU-DPD-BST prediction	t r a i d - - o r e -
GRU-BST prediction	t r a - d i t o r e m

Table 7: Instances where GRU-DPD-BST produces the correct protoform but GRU-BST makes a [m] insertion or deletion error (top) and an instance where the opposite is true (bottom). The success examples are *immaterialem* ‘immaterial’, *consentire* ‘to agree’, and *populare* ‘popular’. The failure example is *traditor* ‘traitor’. Words are aligned manually, and ‘-’ indicates an empty position in the alignment. The positions where an [m] insertion or deletion error occur are shaded.

proxy for sound change rules. It is likely that our DPD models learnt to use unlabeled cognate sets via means other than just sound correspondences.

I Details on Unlabeled Data Ablation Experiments

We perform ablations on unlabeled training at a 10% labeling setting using the group 1 dataset seed, effectively keeping the same subset of labeled data as semisupervised group 1 experiments. This is equivalent to supervised reconstruction but with only 362 and 607 training examples for WikiHan and Rom-phon respectively.

We perform additional hyperparameter tuning for non-trivial semisupervised strategies¹⁶ to account for a large difference in dataset size. The hyperparameter we use are reported in Tables 18 and 19. We then perform 10 runs (random seed) for each strategy other than SUPV.

Table 13 compares the performance of various strategies after ablation, and Table 14 compares the performance with base strategies including IIM, DPD, and DPD-IIM under the configuration of whether to exclude unlabeled data, include unla-

beled data, or include and pseudo-label unlabeled data (+BST, include unlabelled).

In many situations, the non-trivial semisupervised strategies still perform significantly better than SUPV despite the lack of unlabeled data. On WikiHan, DPD-IIM performs the best on all metrics except Trans-DPD-IIM on BCFS. On Rom-Phon, IIM performs the best on all metrics. For strategies involving DPD but not IIM, performance is always higher when unlabeled cognate sets are used. Interestingly, strategies involving IIM sometimes perform worse when using unlabeled cognate sets on Rom-phon. In almost all situations, configurations that include and pseudo-label unlabeled cognate sets perform significantly better than when unlabeled data is excluded.

J Details on 100% Labeled Supervised Reconstruction Experiments

Similar to the setup for unlabeled data ablation at 10% labeling setting (see Section I), we tune additional hyperparameters for IIM, DPD, and DPD-IIM to account for a difference in data. The hyperparameters we obtain are reported in Tables 20 and 21.

The SUPV strategy is equivalent to existing reconstruction models: GRU-SUPV is equivalent to

¹⁶For SUPV, including unlabeled data has no effect, so we reuse hyperparameters and runs from semisupervised reconstruction experiments.

Meloni et al. (2021) and Trans-SUPV is equivalent Kim et al. (2023). Existing hyperparameters and checkpoints exist for SUPV: for WikiHan, we obtain 10 checkpoints from Kim et al. (2023) and 10 additional checkpoints from Lu et al. (2024); for Rom-phon, we obtain 20 checkpoints from Lu et al. (2024) (for hyperparameters, refer to Kim et al. (2023) and Lu et al. (2024)). For all other strategies, we perform 10 runs (random seed).

We compare the results of IIM, DPD, and DPD-IIM (both GRU and Transformer) against existing supervised reconstruction systems in Table 15. We use Meloni et al. (2021) and Kim et al. (2023) as SUPV baselines, and Lu et al. (2024)’s state-of-the-art supervised reconstruction systems (both GRU-BS + GRU Reranker and GRU-BS + Trans. Reranker) as strong baselines. We obtain evaluation results for 20 runs per setup from Lu et al. (2024).

K Sample Outputs at Different Labeling Settings

We compare the protoform predictions of the best-performing strategy-architecture combination (best at a 10% labeling setting) trained on different labeling settings (with group 1 dataset seed)¹⁷. We present sample predictions stratified by the following categories:

- The predictions are correct at all %
- The predictions are correct only above a certain % threshold
- The predictions are correct only below a certain % threshold
- The predictions are incorrect at all %
- All other patterns

Table 23 shows the distribution of test examples in these categories. Tables 25 and 24 show sample protoform predictions in each of these categories (proportionate sampling). For both datasets, correct only above a certain % threshold is the most common category when the predictions differ between labeling setting.

L Sample Outputs from Different Strategies

We show sample predictions from the best performing runs in group 1. Tables 28 and 27 show predictions of the best-performing strategy-architecture

¹⁷At a 100% labeling setting, BST has no effect, so Trans-DPD-IIM-BST is equivalent to Trans-DPD-IIM and GRU-DPD-BST is equivalent to GRU-DPD.

combination and its non-DPD counterpart for each dataset, proportionately sampled by the confusion matrices in Table 26. Tables 30 and 29 compares predictions across all 16 strategy-architecture pairs for 4 randomly selected test examples among those with the most diverse (approximately in the upper quartile¹⁸) protoform predictions.

M Performance Data Visualizations

Figures 6 and 7 show the distribution of semisupervised reconstruction performance by group on 10% labeled datasets. Figures 8 and 9 visualize the group 1 performance of each strategy-architecture combination on datasets with different percentages of label.

N Additional Phonetic Probing Results

A complete set of hierarchical clusterings of learned phoneme embeddings can be found in Figures 10 (Rom-phon) and 11 (WikiHan). Similar to French and Cantonese, phoneme embeddings learned using a DPD-based strategy also appear more interpretable for other languages. For instance, Trans-DPD-IIM-BST creates a big cluster of tones within all phonemes present in daughter languages in WikiHan, and GRU-DPD-BST creates a cluster encompassing most palatalized consonants in Romanian.

O DPD-IIM Implementation

In a hybrid model combining the DPD architecture with Π -model, we first augment the input cognate set to obtain two augmented inputs A and B . Input A is used to train DPD as described in Figure 1. Input B is fed through D2P (and not P2D). A mean square difference loss then is used to minimize the difference between the D2P classifier logits given A and B as inputs.

P Responsible Use of AI

GitHub Copilot has been used as a coding assistant for our model implementation and data analysis. All code generated by Github Copilot is checked manually.

¹⁸The precise cutoff is top 25.65% for WikiHan and top 24.91% for Rom-phon, accounting for ties and integer divisions.

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	21.85% $\textcircled{1}$	1.5150 $\textcircled{1}$	0.3549 $\textcircled{1}$	0.1575 $\textcircled{1}$	0.5485 $\textcircled{1}$
	DPD-BST (ours)	23.41% $\textcircled{1}$	1.4501 $\textcircled{1}$	0.3397 $\textcircled{1}$	0.1441 $\textcircled{1}$	0.5676 $\textcircled{1}$
	DPD-IIM (ours)	23.90% $\textcircled{1}$	1.4488 $\textcircled{1}$	0.3394 $\textcircled{1}$	0.1445 $\textcircled{1}$	0.5622 $\textcircled{1}$
	DPD (ours)	25.46% $\textcircled{1}$	1.3682 $\textcircled{1}$	0.3205 $\textcircled{1}$	0.1374 $\textcircled{1}$	0.5747 $\textcircled{1}$
	IIM-BST	16.04%	1.7841	0.4179	0.1805	0.4872
	BST (Lee, 2013)	16.47%	1.7365	0.4068	0.1736	0.4897
	IIM (Laine and Aila, 2017)	18.18%	1.6594	0.3887	0.1672	0.5076
	SUPV	16.12%	1.7178	0.4024	0.1703	0.4927
GRU	DPD-IIM-BST (ours)	27.70% $\textcircled{1}$	1.2834 $\textcircled{1}$	0.3006 $\textcircled{1}$	0.1207 $\textcircled{1}$	0.6144 $\textcircled{1}$
	DPD-BST (ours)	22.85% $\textcircled{1}$	1.4037 $\textcircled{1}$	0.3288 $\textcircled{1}$	0.1338 $\textcircled{1}$	0.5865 $\textcircled{1}$
	DPD-IIM (ours)	25.33% $\textcircled{1}$	1.3600 $\textcircled{1}$	0.3186 $\textcircled{1}$	0.1286 $\textcircled{1}$	0.5901 $\textcircled{1}$
	DPD (ours)	21.39%	1.4842	0.3477	0.1419	0.5570
	IIM-BST	21.52%	1.4863	0.3481	0.1462	0.5666
	BST (Lee, 2013)	16.22%	1.6901	0.3959	0.1656	0.5152
	IIM (Laine and Aila, 2017)	20.57%	1.5301	0.3584	0.1464	0.5503
	SUPV	15.99%	1.6812	0.3938	0.1616	0.5140

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	47.02% $\textcircled{1}$	0.8648 $\textcircled{1}$	0.2026 $\textcircled{1}$	0.0816	0.7049 $\textcircled{1}$
	DPD-BST (ours)	45.90%	0.8773	0.2055	0.0834	0.7027
	DPD-IIM (ours)	44.72%	0.9246	0.2166	0.0881	0.6884
	DPD (ours)	45.03%	0.9087	0.2129	0.0896	0.6916
	IIM-BST	45.48%	0.8905	0.2086	0.0848	0.6979
	BST (Lee, 2013)	45.30%	0.8887	0.2082	0.0840	0.6988
	IIM (Laine and Aila, 2017)	42.26%	0.9758	0.2286	0.0949	0.6739
	SUPV	42.62%	0.9641	0.2258	0.0945	0.6745
GRU	DPD-IIM-BST (ours)	46.17% $\textcircled{1}$	0.9069 $\textcircled{1}$	0.2124 $\textcircled{1}$	0.0851 $\textcircled{1}$	0.6941 $\textcircled{1}$
	DPD-BST (ours)	45.23% $\textcircled{1}$	0.9045 $\textcircled{1}$	0.2119 $\textcircled{1}$	0.0861 $\textcircled{1}$	0.6940 $\textcircled{1}$
	DPD-IIM (ours)	44.12% $\textcircled{1}$	0.9554 $\textcircled{1}$	0.2238 $\textcircled{1}$	0.0909 $\textcircled{1}$	0.6794 $\textcircled{1}$
	DPD (ours)	43.11% $\textcircled{1}$	0.9605 $\textcircled{1}$	0.2250 $\textcircled{1}$	0.0906 $\textcircled{1}$	0.6764 $\textcircled{1}$
	IIM-BST	44.32% $\textcircled{1}$	0.9475 $\textcircled{1}$	0.2220 $\textcircled{1}$	0.0900 $\textcircled{1}$	0.6839 $\textcircled{1}$
	BST (Lee, 2013)	40.47%	1.0143	0.2376	0.0953	0.6614
	IIM (Laine and Aila, 2017)	41.37%	1.0056	0.2356	0.0957	0.6661
	SUPV	39.24%	1.0448	0.2447	0.0992	0.6530

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	49.93%	0.8229	0.1928	0.0789	0.7163
	DPD-BST (ours)	49.99%	0.8170	0.1914	0.0782	0.7177
	DPD-IIM (ours)	50.44%	0.8090	0.1895	0.0778	0.7212
	DPD (ours)	48.66%	0.8457	0.1981	0.0811	0.7078
	IIM-BST	49.64%	0.8268	0.1937	0.0790	0.7149
	BST (Lee, 2013)	49.50%	0.8270	0.1937	0.0782	0.7134
	IIM (Laine and Aila, 2017)	47.56%	0.8700	0.2038	0.0834	0.7025
	SUPV	46.92%	0.8842	0.2071	0.0856	0.6958
GRU	DPD-IIM-BST (ours)	49.43% $\textcircled{1}$	0.8389 $\textcircled{1}$	0.1965 $\textcircled{1}$	0.0789 $\textcircled{1}$	0.7126 $\textcircled{1}$
	DPD-BST (ours)	48.79% $\textcircled{1}$	0.8342 $\textcircled{1}$	0.1954 $\textcircled{1}$	0.0797 $\textcircled{1}$	0.7116 $\textcircled{1}$
	DPD-IIM (ours)	49.84% $\textcircled{1}$	0.8370 $\textcircled{1}$	0.1961 $\textcircled{1}$	0.0791 $\textcircled{1}$	0.7136 $\textcircled{1}$
	DPD (ours)	47.67% $\textcircled{1}$	0.8696	0.2037	0.0824	0.7020 $\textcircled{1}$
	IIM-BST	47.73%	0.8788	0.2059	0.0825	0.7005
	BST (Lee, 2013)	45.59%	0.9066	0.2124	0.0853	0.6897
	IIM (Laine and Aila, 2017)	46.28%	0.9045	0.2119	0.0857	0.6924
	SUPV	44.82%	0.9199	0.2155	0.0866	0.6867

Table 8: Performance of all strategies on 5% (top), 20% (middle), and 30% (bottom) labeled WikiHan for each architecture, averaged across 10 runs in group 1. Bold: best-performing strategy for the corresponding architecture; $\textcircled{1}$: significantly better than all weak baselines (SUPV, BST, and IIM) with $p < 0.01$; $\textcircled{1}$: significantly better than the IIM-BST strong baseline and all weak baselines with $p < 0.01$.

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	27.38% $\textcircled{1}$	1.5508 $\textcircled{1}$	0.1730 $\textcircled{1}$	0.0717 $\textcircled{1}$	0.7523 $\textcircled{1}$
	DPD-BST (ours)	26.04% $\textcircled{1}$	1.6169 $\textcircled{1}$	0.1803 $\textcircled{1}$	0.0751 $\textcircled{1}$	0.7403 $\textcircled{1}$
	DPD-IIM (ours)	16.33%	2.3046	0.2570	0.1149	0.6474
	DPD (ours)	23.52% $\textcircled{1}$	1.8284 $\textcircled{1}$	0.2039 $\textcircled{1}$	0.0854 $\textcircled{1}$	0.7110 $\textcircled{1}$
	IIM-BST	18.19%	2.1148	0.2359	0.1127	0.6724
	BST (Lee, 2013)	16.65%	2.2502	0.2510	0.1222	0.6560
	IIM (Laine and Aila, 2017)	10.66%	2.8262	0.3152	0.1468	0.5806
	SUPV	14.54%	2.4150	0.2694	0.1188	0.6314
GRU	DPD-IIM-BST (ours)	30.68%	1.3867 $\textcircled{1}$	0.1547 $\textcircled{1}$	0.0562 $\textcircled{1}$	0.7788 $\textcircled{1}$
	DPD-BST (ours)	30.94%	1.3731 $\textcircled{1}$	0.1531 $\textcircled{1}$	0.0518 $\textcircled{1}$	0.7803 $\textcircled{1}$
	DPD-IIM (ours)	25.48%	1.6870	0.1882	0.0740	0.7328
	DPD (ours)	24.13%	1.7210	0.1919	0.0733	0.7269
	IIM-BST	28.96%	1.4869	0.1658	0.0654	0.7673
	BST (Lee, 2013)	30.29%	1.4572	0.1625	0.0605	0.7687
	IIM (Laine and Aila, 2017)	22.81%	1.8119	0.2021	0.0809	0.7141
	SUPV	24.11%	1.7354	0.1936	0.0747	0.7255

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	40.60% $\textcircled{1}$	1.1887 $\textcircled{1}$	0.1326 $\textcircled{1}$	0.0531 $\textcircled{1}$	0.8004 $\textcircled{1}$
	DPD-BST (ours)	39.49% $\textcircled{1}$	1.2052 $\textcircled{1}$	0.1344 $\textcircled{1}$	0.0515 $\textcircled{1}$	0.7968
	DPD-IIM (ours)	38.84%	1.2421	0.1385	0.0556	0.7921
	DPD (ours)	37.79%	1.3035	0.1454	0.0562	0.7824
	IIM-BST	40.22% $\textcircled{1}$	1.2023 $\textcircled{1}$	0.1341 $\textcircled{1}$	0.0536	0.7982 $\textcircled{1}$
	BST (Lee, 2013)	38.18%	1.2467	0.1390	0.0553	0.7919
	IIM (Laine and Aila, 2017)	37.98%	1.2797	0.1427	0.0585	0.7865
	SUPV	35.14%	1.3910	0.1551	0.0609	0.7699
GRU	DPD-IIM-BST (ours)	42.57% $\textcircled{1}$	1.1351 $\textcircled{1}$	0.1266 $\textcircled{1}$	0.0468	0.8080 $\textcircled{1}$
	DPD-BST (ours)	42.13% $\textcircled{1}$	1.1079 $\textcircled{1}$	0.1236 $\textcircled{1}$	0.0426 $\textcircled{1}$	0.8134 $\textcircled{1}$
	DPD-IIM (ours)	38.51%	1.2476	0.1392	0.0519	0.7922
	DPD (ours)	37.47%	1.2702	0.1417	0.0510	0.7890
	IIM-BST	41.01%	1.1721	0.1307	0.0490	0.8045
	BST (Lee, 2013)	40.33%	1.1779	0.1314	0.0474	0.8031
	IIM (Laine and Aila, 2017)	36.71%	1.2994	0.1449	0.0529	0.7860
	SUPV	37.25%	1.2947	0.1444	0.0524	0.7834

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	45.05% $\textcircled{1}$	1.0812 $\textcircled{1}$	0.1206 $\textcircled{1}$	0.0475	0.8162 $\textcircled{1}$
	DPD-BST (ours)	43.28%	1.1281	0.1258	0.0482	0.8089
	DPD-IIM (ours)	42.65%	1.1182 $\textcircled{1}$	0.1247 $\textcircled{1}$	0.0501	0.8132
	DPD (ours)	43.07%	1.1588	0.1293	0.0492	0.8034
	IIM-BST	44.21% $\textcircled{1}$	1.1196	0.1249	0.0499	0.8102
	BST (Lee, 2013)	42.76%	1.1515	0.1284	0.0485	0.8045
	IIM (Laine and Aila, 2017)	41.61%	1.1447	0.1277	0.0519	0.8095
	SUPV	40.61%	1.2381	0.1381	0.0526	0.7912
GRU	DPD-IIM-BST (ours)	45.19% $\textcircled{1}$	1.0966	0.1223	0.0457	0.8127
	DPD-BST (ours)	45.74% $\textcircled{1}$	1.0417 $\textcircled{1}$	0.1162 $\textcircled{1}$	0.0400 $\textcircled{1}$	0.8214 $\textcircled{1}$
	DPD-IIM (ours)	43.31%	1.1358	0.1267	0.0466	0.8089
	DPD (ours)	41.57%	1.1681	0.1303	0.0473	0.8037
	IIM-BST	44.70%	1.1134	0.1242	0.0472	0.8117
	BST (Lee, 2013)	44.05%	1.1081	0.1236	0.0446	0.8125
	IIM (Laine and Aila, 2017)	41.11%	1.1838	0.1320	0.0482	0.8023
	SUPV	40.94%	1.1985	0.1337	0.0479	0.7983

Table 9: Performance of all strategies on 5% (top), 20% (middle), and 30% (bottom) labeled Rom-phon for each architecture, averaged across 10 runs in group 1. Bold: best-performing strategy for the corresponding architecture; $\textcircled{1}$: significantly better than all weak baselines (SUPV, BST, and IIM) with $p < 0.01$; $\textcircled{!}$: significantly better than the IIM-BST strong baseline and all weak baselines with $p < 0.01$.

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD- Π M-BST (ours)	41.15% \mathbb{B}	0.9955 \mathbb{B}	0.2326 \mathbb{B}	0.0920 \mathbb{B}	0.6703 \mathbb{B}
	DPD-BST (ours)	39.83% \mathbb{B}	1.0212 \mathbb{B}	0.2386 \mathbb{B}	0.0944 \mathbb{B}	0.6636 \mathbb{B}
	DPD- Π M (ours)	38.28% \mathbb{B}	1.0612 \mathbb{B}	0.2480 \mathbb{B}	0.0968 \mathbb{B}	0.6478 \mathbb{B}
	DPD (ours)	39.61% \mathbb{B}	1.0295 \mathbb{B}	0.2406 \mathbb{B}	0.0946 \mathbb{B}	0.6547 \mathbb{B}
	Π M-BST	34.76%	1.1447	0.2675	0.1078	0.6336
	BST (Lee, 2013)	34.93%	1.1490	0.2685	0.1086	0.6281
	Π M (Laine and Aila, 2017)	34.43%	1.1721	0.2739	0.1087	0.6151
	SUPV	33.57%	1.1920	0.2785	0.1110	0.6078
GRU	DPD- Π M-BST (ours)	40.50% \mathbb{B}	0.9977 \mathbb{B}	0.2331 \mathbb{B}	0.0894 \mathbb{B}	0.6716 \mathbb{B}
	DPD-BST (ours)	36.64% \mathbb{B}	1.0794 \mathbb{B}	0.2522 \mathbb{B}	0.0978 \mathbb{B}	0.6512 \mathbb{B}
	DPD- Π M (ours)	38.67% \mathbb{B}	1.0460 \mathbb{B}	0.2444 \mathbb{B}	0.0933 \mathbb{B}	0.6533 \mathbb{B}
	DPD (ours)	35.12% $\textcircled{1}$	1.1339 $\textcircled{1}$	0.2650 $\textcircled{1}$	0.1030	0.6282 $\textcircled{1}$
	Π M-BST	35.28% $\textcircled{1}$	1.1266 $\textcircled{3}$	0.2632 $\textcircled{3}$	0.1016 $\textcircled{1}$	0.6360 $\textcircled{3,2}$
	BST (Lee, 2013)	29.24%	1.2929	0.3021	0.1174	0.5931
	Π M (Laine and Aila, 2017)	33.69%	1.1774	0.2751	0.1063	0.6173
	SUPV	28.98%	1.3066	0.3053	0.1185	0.5827
<hr/>						
Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD- Π M-BST (ours)	32.94% \mathbb{B}	1.3682 \mathbb{B}	0.1533 \mathbb{B}	0.0597 \mathbb{B}	0.7717 \mathbb{B}
	DPD-BST (ours)	31.80% \mathbb{B}	1.4636 \mathbb{B}	0.1640 \mathbb{B}	0.0634 \mathbb{B}	0.7553 \mathbb{B}
	DPD- Π M (ours)	27.58%	1.6277	0.1824	0.0732 $\textcircled{1,2}$	0.7328
	DPD (ours)	30.15% \mathbb{B}	1.5413 \mathbb{B}	0.1727 \mathbb{B}	0.0678 \mathbb{B}	0.7449 \mathbb{B}
	Π M-BST	30.36% \mathbb{B}	1.5174 \mathbb{B}	0.1700 \mathbb{B}	0.0681 \mathbb{B}	0.7487 \mathbb{B}
	BST (Lee, 2013)	28.70%	1.6216	0.1817	0.0754	0.7332
	Π M (Laine and Aila, 2017)	25.18%	1.7561	0.1968	0.0804	0.7144
	SUPV	25.37%	1.7756	0.1990	0.0802	0.7122
GRU	DPD- Π M-BST (ours)	35.34% \mathbb{B}	1.2634 \mathbb{B}	0.1416 \mathbb{B}	0.0480 \mathbb{B}	0.7899 \mathbb{B}
	DPD-BST (ours)	35.83% \mathbb{B}	1.2402 \mathbb{B}	0.1390 \mathbb{B}	0.0458 \mathbb{B}	0.7929 \mathbb{B}
	DPD- Π M (ours)	30.19%	1.5103	0.1692	0.0619	0.7511
	DPD (ours)	29.83%	1.5071	0.1689	0.0604	0.7512
	Π M-BST	34.31%	1.3134 $\textcircled{1}$	0.1472 $\textcircled{1}$	0.0523	0.7846 $\textcircled{1,4}$
	BST (Lee, 2013)	34.22%	1.3163	0.1475	0.0514	0.7819
	Π M (Laine and Aila, 2017)	28.16%	1.5648	0.1754	0.0636	0.7447
	SUPV	29.30%	1.5292	0.1714	0.0610	0.7473

Table 10: Transductive evaluation of all strategies on 10% labeled WikiHan (top) and Rom-phon (bottom) for each architecture, averaged across all runs in four groups (10 runs per strategy-architecture combination per group). Bold: best-performing strategy for the corresponding architecture and dataset; $\textcircled{1}$: significantly better than all weak baselines (SUPV, BST, and Π M) on group 1 with $p < 0.01$; \mathbb{B} : significantly better than the Π M-BST strong baseline and all weak baselines on group 1 with $p < 0.01$; $\textcircled{2}$, $\textcircled{3}$, $\textcircled{4}$, $\textcircled{2}$, $\textcircled{3}$, $\textcircled{4}$: likewise for groups 2, 3, and 4.

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	21.40% ①	1.5246 ①	0.3564 ①	0.1525 ①	0.5403 ①
	DPD-BST (ours)	23.74% ①	1.4558 ①	0.3403 ①	0.1409 ①	0.5608 ①
	DPD-IIM (ours)	23.98% ①	1.4538 ①	0.3398 ①	0.1407 ①	0.5564 ①
	DPD (ours)	25.35% ①	1.3911 ①	0.3252 ①	0.1348 ①	0.5636 ①
	IIM-BST	15.97%	1.7986	0.4204	0.1786	0.4774
	BST (Lee, 2013)	16.43%	1.7464	0.4082	0.1681	0.4805
	IIM (Laine and Aila, 2017)	17.80%	1.6913	0.3953	0.1665	0.4935
	SUPV	16.42%	1.7418	0.4071	0.1681	0.4815
GRU	DPD-IIM-BST (ours)	28.38% ①	1.2499 ①	0.2922 ①	0.1126 ①	0.6175 ①
	DPD-BST (ours)	24.32% ①	1.3716 ①	0.3206 ①	0.1258 ①	0.5876 ①
	DPD-IIM (ours)	25.73% ①	1.3493 ①	0.3154 ①	0.1238 ①	0.5887 ①
	DPD (ours)	22.27%	1.4420	0.3371	0.1329 ①	0.5599
	IIM-BST	22.06%	1.4639	0.3422	0.1395	0.5641 ①
	BST (Lee, 2013)	16.50%	1.6554	0.3870	0.1596	0.5169
	IIM (Laine and Aila, 2017)	20.89%	1.5094	0.3528	0.1407	0.5486
	SUPV	16.72%	1.6520	0.3861	0.1572	0.5120

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	48.17% ①	0.8513 ①	0.1986 ①	0.0749 ①	0.7039 ①
	DPD-BST (ours)	46.61%	0.8666 ①	0.2022 ①	0.0768	0.7004 ①
	DPD-IIM (ours)	44.99%	0.8983	0.2096	0.0798	0.6912
	DPD (ours)	44.88%	0.9125	0.2129	0.0838	0.6852
	IIM-BST	45.93%	0.8904	0.2077	0.0788	0.6939
	BST (Lee, 2013)	45.43%	0.8995	0.2099	0.0791	0.6904
	IIM (Laine and Aila, 2017)	41.61%	0.9829	0.2293	0.0898	0.6671
	SUPV	42.25%	0.9698	0.2263	0.0888	0.6683
GRU	DPD-IIM-BST (ours)	47.92% ①	0.8527 ①	0.1990 ①	0.0755 ①	0.7031 ①
	DPD-BST (ours)	46.07% ①	0.8783 ①	0.2049 ①	0.0791 ①	0.6966 ①
	DPD-IIM (ours)	44.82% ①	0.9194 ①	0.2145 ①	0.0827 ①	0.6839 ①
	DPD (ours)	44.01% ①	0.9309 ①	0.2172 ①	0.0855 ①	0.6801 ①
	IIM-BST	44.73% ①	0.9232 ①	0.2154 ①	0.0831 ①	0.6836 ①
	BST (Lee, 2013)	41.60%	0.9888	0.2307	0.0908	0.6646
	IIM (Laine and Aila, 2017)	42.34%	0.9793	0.2285	0.0889	0.6666
	SUPV	40.12%	1.0283	0.2399	0.0938	0.6518

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	50.21%	0.8336	0.1941	0.0738	0.7087
	DPD-BST (ours)	50.42% ①	0.8217 ①	0.1913 ①	0.0716 ①	0.7121 ①
	DPD-IIM (ours)	50.60% ①	0.8105 ①	0.1887 ①	0.0716 ①	0.7166 ①
	DPD (ours)	48.51%	0.8571	0.1995	0.0760	0.7013
	IIM-BST	49.81%	0.8389	0.1953	0.0744	0.7071
	BST (Lee, 2013)	49.40%	0.8476	0.1973	0.0751	0.7046
	IIM (Laine and Aila, 2017)	46.97%	0.8905	0.2073	0.0804	0.6923
	SUPV	46.27%	0.9019	0.2100	0.0805	0.6881
GRU	DPD-IIM-BST (ours)	49.55% ①	0.8296 ①	0.1931 ①	0.0737 ①	0.7119 ①
	DPD-BST (ours)	49.01% ①	0.8413 ①	0.1959 ①	0.0746 ①	0.7066 ①
	DPD-IIM (ours)	49.55% ①	0.8310 ①	0.1935 ①	0.0732 ①	0.7107 ①
	DPD (ours)	47.95%	0.8635	0.2010	0.0770	0.7009
	IIM-BST	48.31% ①	0.8653 ①	0.2015 ①	0.0761 ①	0.6996 ①
	BST (Lee, 2013)	46.35%	0.9051	0.2107	0.0812	0.6879
	IIM (Laine and Aila, 2017)	46.00%	0.9041	0.2105	0.0803	0.6893
	SUPV	44.53%	0.9415	0.2192	0.0847	0.6787

Table 11: Transductive evaluation of all strategies on 5% (top), 20% (middle), and 30% (bottom) labeled WikiHan for each architecture, averaged across 10 runs in group 1. Bold: best-performing strategy for the corresponding architecture; **①**: significantly better than all weak baselines (SUPV, BST, and IIM) with $p < 0.01$; **②**: significantly better than the IIM-BST strong baseline and all weak baselines with $p < 0.01$.

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	25.73% $\textcircled{1}$	1.6015 $\textcircled{1}$	0.1797 $\textcircled{1}$	0.0735 $\textcircled{1}$	0.7399 $\textcircled{1}$
	DPD-BST (ours)	24.62% $\textcircled{1}$	1.7219 $\textcircled{1}$	0.1932 $\textcircled{1}$	0.0797 $\textcircled{1}$	0.7199 $\textcircled{1}$
	DPD-IIM (ours)	15.46%	2.3351	0.2620	0.1148 $\textcircled{1}$	0.6387
	DPD (ours)	22.38% $\textcircled{1}$	1.8658 $\textcircled{1}$	0.2093 $\textcircled{1}$	0.0861 $\textcircled{1}$	0.7002 $\textcircled{1}$
	IIM-BST	17.29%	2.2002	0.2468	0.1160	0.6544
	BST (Lee, 2013)	15.85%	2.3260	0.2610	0.1248	0.6394
	IIM (Laine and Aila, 2017)	9.67%	2.8616	0.3210	0.1461	0.5706
	SUPV	14.06%	2.4557	0.2755	0.1203	0.6219
GRU	DPD-IIM-BST (ours)	29.62% $\textcircled{1}$	1.3946 $\textcircled{1}$	0.1565 $\textcircled{1}$	0.0553 $\textcircled{1}$	0.7724 $\textcircled{1}$
	DPD-BST (ours)	29.67% $\textcircled{1}$	1.3977 $\textcircled{1}$	0.1568 $\textcircled{1}$	0.0519 $\textcircled{1}$	0.7712 $\textcircled{1}$
	DPD-IIM (ours)	24.86%	1.7148	0.1924	0.0731	0.7230
	DPD (ours)	24.02%	1.7394	0.1951	0.0734	0.7185
	IIM-BST	28.42%	1.5011	0.1684	0.0648	0.7601
	BST (Lee, 2013)	28.49%	1.5016	0.1685	0.0624	0.7568
	IIM (Laine and Aila, 2017)	22.36%	1.8316	0.2055	0.0807	0.7054
	SUPV	23.49%	1.7649	0.1980	0.0751	0.7146

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	38.96% $\textcircled{1}$	1.2502 $\textcircled{1}$	0.1404 $\textcircled{1}$	0.0539 $\textcircled{1}$	0.7860 $\textcircled{1}$
	DPD-BST (ours)	37.75% $\textcircled{1}$	1.3183	0.1480	0.0551 $\textcircled{1}$	0.7740
	DPD-IIM (ours)	36.67%	1.2899 $\textcircled{1}$	0.1448 $\textcircled{1}$	0.0552 $\textcircled{1}$	0.7804 $\textcircled{1}$
	DPD (ours)	36.11%	1.3521	0.1518	0.0569	0.7706
	IIM-BST	38.04% $\textcircled{1}$	1.2899 $\textcircled{1}$	0.1448 $\textcircled{1}$	0.0557 $\textcircled{1}$	0.7796 $\textcircled{1}$
	BST (Lee, 2013)	35.99%	1.3854	0.1555	0.0609	0.7647
	IIM (Laine and Aila, 2017)	35.57%	1.3284	0.1491	0.0580	0.7748
	SUPV	33.61%	1.4521	0.1630	0.0621	0.7569
GRU	DPD-IIM-BST (ours)	41.19% $\textcircled{1}$	1.1635 $\textcircled{1}$	0.1306 $\textcircled{1}$	0.0462 $\textcircled{1}$	0.7991 $\textcircled{1}$
	DPD-BST (ours)	41.00% $\textcircled{1}$	1.1345 $\textcircled{1}$	0.1274 $\textcircled{1}$	0.0425 $\textcircled{1}$	0.8041 $\textcircled{1}$
	DPD-IIM (ours)	36.77%	1.2877	0.1446	0.0513	0.7811
	DPD (ours)	35.94%	1.3036	0.1464	0.0511	0.7788
	IIM-BST	40.07%	1.1982	0.1345	0.0482	0.7960
	BST (Lee, 2013)	39.21%	1.2014	0.1349	0.0472	0.7945
	IIM (Laine and Aila, 2017)	35.17%	1.3368	0.1501	0.0530	0.7751
	SUPV	35.60%	1.3277	0.1491	0.0524	0.7736

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	DPD-IIM-BST (ours)	41.96% $\textcircled{1}$	1.1563 $\textcircled{1}$	0.1302 $\textcircled{1}$	0.0485 $\textcircled{1}$	0.7998 $\textcircled{1}$
	DPD-BST (ours)	40.53% $\textcircled{1}$	1.2425	0.1399	0.0513	0.7853
	DPD-IIM (ours)	39.29%	1.1843 $\textcircled{1}$	0.1333 $\textcircled{1}$	0.0503 $\textcircled{1}$	0.7974 $\textcircled{1}$
	DPD (ours)	39.57%	1.2368	0.1392	0.0509	0.7867
	IIM-BST	41.02% $\textcircled{1}$	1.1950	0.1345	0.0510	0.7929
	BST (Lee, 2013)	39.94%	1.2865	0.1448	0.0542	0.7773
	IIM (Laine and Aila, 2017)	37.92%	1.2134	0.1366	0.0523	0.7933
	SUPV	37.44%	1.3133	0.1478	0.0545	0.7755
GRU	DPD-IIM-BST (ours)	43.29% $\textcircled{1}$	1.1312	0.1273	0.0447	0.8023
	DPD-BST (ours)	43.92% $\textcircled{1}$	1.0702 $\textcircled{1}$	0.1205 $\textcircled{1}$	0.0402 $\textcircled{1}$	0.8119 $\textcircled{1}$
	DPD-IIM (ours)	40.62%	1.1738	0.1321	0.0461	0.7974
	DPD (ours)	38.84%	1.2149	0.1368	0.0470	0.7903
	IIM-BST	42.93% $\textcircled{1}$	1.1497	0.1294	0.0464	0.8006
	BST (Lee, 2013)	41.81%	1.1402	0.1284	0.0439	0.8016
	IIM (Laine and Aila, 2017)	39.08%	1.2188	0.1372	0.0475	0.7911
	SUPV	38.77%	1.2268	0.1381	0.0471	0.7878

Table 12: Transductive evaluation of all strategies on 5% (top), 20% (middle), and 30% (bottom) labeled Rom-phon for each architecture, averaged across 10 runs in group 1. Bold: best-performing strategy for the corresponding architecture; $\textcircled{1}$: significantly better than all weak baselines (SUPV, BST, and IIM) with $p < 0.01$; $\textcircled{2}$: significantly better than the IIM-BST strong baseline and all weak baselines with $p < 0.01$.

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	SUPV	33.25%	1.1891	0.2785	0.1140	0.6138
	IIM (Laine and Aila, 2017)	32.21%	1.2141	0.2844	0.1140	0.6077
	DPD (ours)	33.57%	1.1621	0.2722	0.1108	0.6246 \bullet
	DPD-IIM (ours)	34.93% \bullet	1.1307 \bullet	0.2649 \bullet	0.1097	0.6344 \bullet
GRU	SUPV	28.16%	1.3257	0.3105	0.1234	0.5835
	IIM (Laine and Aila, 2017)	29.46% $\textcircled{1}$	1.2475 $\textcircled{1}$	0.2922 $\textcircled{1}$	0.1157 $\textcircled{1}$	0.6124 $\textcircled{1}$
	DPD (ours)	26.57%	1.3424	0.3144	0.1259	0.5830
	DPD-IIM (ours)	30.27% $\textcircled{1}$	1.2393 $\textcircled{1}$	0.2903 $\textcircled{1}$	0.1156 $\textcircled{1}$	0.6067 $\textcircled{1}$

Architecture	Strategy	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	SUPV	26.99%	1.7331	0.1933	0.0794	0.7218
	IIM (Laine and Aila, 2017)	30.01% $\textcircled{1}$	1.5261 $\textcircled{1}$	0.1702 $\textcircled{1}$	0.0699 $\textcircled{1}$	0.7536 $\textcircled{1}$
	DPD (ours)	28.22%	1.6713	0.1864	0.0745	0.7308
	DPD-IIM (ours)	27.42%	1.5860 $\textcircled{1}$	0.1769 $\textcircled{1}$	0.0741	0.7465 $\textcircled{1}$
GRU	SUPV	30.69%	1.5018	0.1675	0.0612	0.7558
	IIM (Laine and Aila, 2017)	32.38% $\textcircled{1}$	1.4232 $\textcircled{1}$	0.1587 $\textcircled{1}$	0.0591	0.7718 $\textcircled{1}$
	DPD (ours)	30.68%	1.5010	0.1674	0.0629	0.7588
	DPD-IIM (ours)	32.35% $\textcircled{1}$	1.4294 $\textcircled{1}$	0.1594 $\textcircled{1}$	0.0596	0.7702 $\textcircled{1}$

Table 13: Performance when unlabeled cognate sets are excluded for 10% labeled group 1 WikiHan (top) and 10% labeled group 1 Rom-phon (bottom), averaged across 10 runs. Bold: best-performing strategy for the corresponding architecture; $\textcircled{1}$: significantly better than SUPV ($p < 0.01$); \bullet : significantly better than both SUPV and IIM ($p < 0.01$).

Architecture	Base Strategy	Configuration	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	SUPV	exclude/include unlabeled	32.41%	1.1922	0.2793	0.1145	0.6122
		+BST, include unlabeled	34.87% $\text{\textcircled{1}}$	1.1277 $\text{\textcircled{1}}$	0.2641 $\text{\textcircled{1}}$	0.1088 $\text{\textcircled{1}}$	0.6388 $\text{\textcircled{1}}$
	IIM	exclude unlabeled	32.21%	1.2141	0.2844	0.1140	0.6077
		include unlabeled	34.50% $\text{\textcircled{1}}$	1.1547 $\text{\textcircled{1}}$	0.2705 $\text{\textcircled{1}}$	0.1102	0.6256 $\text{\textcircled{1}}$
		+BST, include unlabeled	34.41%	1.1348 $\text{\textcircled{1}}$	0.2658 $\text{\textcircled{1}}$	0.1067 $\text{\textcircled{1}}$	0.6418 $\text{\textcircled{1}}$
	DPD	exclude unlabeled	33.57%	1.1621	0.2722	0.1108	0.6246
		include unlabeled	39.56% $\text{\textcircled{1}}$	1.0153 $\text{\textcircled{1}}$	0.2378 $\text{\textcircled{1}}$	0.0972 $\text{\textcircled{1}}$	0.6628 $\text{\textcircled{1}}$
		+BST, include unlabeled	39.61% $\text{\textcircled{1}}$	1.0051 $\text{\textcircled{1}}$	0.2354 $\text{\textcircled{1}}$	0.0948 $\text{\textcircled{1}}$	0.6722 $\text{\textcircled{1}}$
	DPD-IIM	exclude unlabeled	34.93%	1.1307	0.2649	0.1097	0.6344
		include unlabeled	37.16% $\text{\textcircled{1}}$	1.0819 $\text{\textcircled{1}}$	0.2534 $\text{\textcircled{1}}$	0.1027 $\text{\textcircled{1}}$	0.6469 $\text{\textcircled{1}}$
		+BST, include unlabeled	39.93% $\text{\textcircled{1}}$	0.9997 $\text{\textcircled{1}}$	0.2342 $\text{\textcircled{1}}$	0.0959 $\text{\textcircled{1}}$	0.6747 $\text{\textcircled{1}}$
	GRU	SUPV	exclude/include unlabeled	27.42%	1.3288	0.3112	0.1238
+BST, include unlabeled			29.44% $\text{\textcircled{1}}$	1.2600 $\text{\textcircled{1}}$	0.2951 $\text{\textcircled{1}}$	0.1167 $\text{\textcircled{1}}$	0.6071 $\text{\textcircled{1}}$
IIM		exclude unlabeled	29.46%	1.2475	0.2922	0.1157	0.6124
		include unlabeled	31.32%	1.2195	0.2856	0.1129	0.6145
		+BST, include unlabeled	35.46% $\text{\textcircled{1}}$	1.1168 $\text{\textcircled{1}}$	0.2616 $\text{\textcircled{1}}$	0.1026 $\text{\textcircled{1}}$	0.6452 $\text{\textcircled{1}}$
DPD		exclude unlabeled	26.57%	1.3424	0.3144	0.1259	0.5830
		include unlabeled	33.62% $\text{\textcircled{1}}$	1.1666 $\text{\textcircled{1}}$	0.2733 $\text{\textcircled{1}}$	0.1102 $\text{\textcircled{1}}$	0.6265 $\text{\textcircled{1}}$
		+BST, include unlabeled	35.59% $\text{\textcircled{1}}$	1.1035 $\text{\textcircled{1}}$	0.2585 $\text{\textcircled{1}}$	0.1031 $\text{\textcircled{1}}$	0.6499 $\text{\textcircled{1}}$
DPD-IIM		exclude unlabeled	30.27%	1.2393	0.2903	0.1156	0.6067
		include unlabeled	36.19% $\text{\textcircled{1}}$	1.0992 $\text{\textcircled{1}}$	0.2575 $\text{\textcircled{1}}$	0.1007 $\text{\textcircled{1}}$	0.6440 $\text{\textcircled{1}}$
		+BST, include unlabeled	39.92% $\text{\textcircled{1}}$	1.0093 $\text{\textcircled{1}}$	0.2364 $\text{\textcircled{1}}$	0.0952 $\text{\textcircled{1}}$	0.6735 $\text{\textcircled{1}}$

Architecture	Base Strategy	Configuration	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
Transformer	SUPV	exclude/include unlabeled	27.66%	1.6753	0.1869	0.0758	0.7303
		+BST, include unlabeled	30.56% $\text{\textcircled{1}}$	1.4712 $\text{\textcircled{1}}$	0.1641 $\text{\textcircled{1}}$	0.0681 $\text{\textcircled{1}}$	0.7610 $\text{\textcircled{1}}$
	IIM	exclude unlabeled	30.01%	1.5261	0.1702	0.0699	0.7536
		include unlabeled	27.69%	1.6483	0.1838	0.0751	0.7332
		+BST, include unlabeled	31.90% $\text{\textcircled{1}}$	1.3935 $\text{\textcircled{1}}$	0.1554 $\text{\textcircled{1}}$	0.0636 $\text{\textcircled{1}}$	0.7740 $\text{\textcircled{1}}$
	DPD	exclude unlabeled	28.22%	1.6713	0.1864	0.0745	0.7308
		include unlabeled	31.81% $\text{\textcircled{1}}$	1.5031 $\text{\textcircled{1}}$	0.1677 $\text{\textcircled{1}}$	0.0668 $\text{\textcircled{1}}$	0.7543 $\text{\textcircled{1}}$
		+BST, include unlabeled	33.96% $\text{\textcircled{1}}$	1.3332 $\text{\textcircled{1}}$	0.1487 $\text{\textcircled{1}}$	0.0591 $\text{\textcircled{1}}$	0.7812 $\text{\textcircled{1}}$
	DPD-IIM	exclude unlabeled	27.42%	1.5860	0.1769	0.0741	0.7465
		include unlabeled	30.09% $\text{\textcircled{1}}$	1.5505	0.1729	0.0708 $\text{\textcircled{1}}$	0.7482
		+BST, include unlabeled	34.33% $\text{\textcircled{1}}$	1.3121 $\text{\textcircled{1}}$	0.1463 $\text{\textcircled{1}}$	0.0592 $\text{\textcircled{1}}$	0.7856 $\text{\textcircled{1}}$
	GRU	SUPV	exclude/include unlabeled	30.01%	1.5156	0.1690	0.0607
+BST, include unlabeled			35.21% $\text{\textcircled{1}}$	1.3284 $\text{\textcircled{1}}$	0.1482 $\text{\textcircled{1}}$	0.0530 $\text{\textcircled{1}}$	0.7857 $\text{\textcircled{1}}$
IIM		exclude unlabeled	32.38%	1.4232	0.1587	0.0591	0.7718
		include unlabeled	29.42%	1.5511	0.1730	0.0643	0.7529
		+BST, include unlabeled	35.38% $\text{\textcircled{1}}$	1.3072 $\text{\textcircled{1}}$	0.1458 $\text{\textcircled{1}}$	0.0537 $\text{\textcircled{1}}$	0.7896 $\text{\textcircled{1}}$
DPD		exclude unlabeled	30.68%	1.5010	0.1674	0.0629	0.7588
		include unlabeled	30.74%	1.4816	0.1652	0.0591 $\text{\textcircled{1}}$	0.7589
		+BST, include unlabeled	36.45% $\text{\textcircled{1}}$	1.2469 $\text{\textcircled{1}}$	0.1391 $\text{\textcircled{1}}$	0.0473 $\text{\textcircled{1}}$	0.7977 $\text{\textcircled{1}}$
DPD-IIM		exclude unlabeled	32.35%	1.4294	0.1594	0.0596	0.7702
		include unlabeled	31.81%	1.4930	0.1665	0.0628	0.7583
		+BST, include unlabeled	36.42% $\text{\textcircled{1}}$	1.2529 $\text{\textcircled{1}}$	0.1397 $\text{\textcircled{1}}$	0.0493 $\text{\textcircled{1}}$	0.7957 $\text{\textcircled{1}}$

Table 14: Performance comparison between whether to exclude unlabeled data, include unlabeled data, or include and pseudo-label unlabeled data (+BST, include unlabelled), evaluated on 10% labeled group 1 WikiHan (top) and 10% labeled group 1 Rom-phon (bottom), averaged across 10 runs. Bold: best performance for the corresponding base strategy; $\text{\textcircled{1}}$: significantly better than when unlabeled data are not used ($p < 0.01$). For SUPV without BST, including and excluding unlabeled data are equivalent.

Dataset	Reconstruction System	ACC% \uparrow	TED \downarrow	TER \downarrow	FER \downarrow	BCFS \uparrow
WikiHan	GRU-SUPV (Meloni et al., 2021)	55.58%	0.7360	0.1724	0.0686	0.7426
	Trans-SUPV (Kim et al., 2023)	54.62%	0.7453	0.1746	0.0696	0.7393
	GRU-BS + GRU Reranker (Lu et al., 2024)	57.14% \circ	0.7045 \circ	0.1650 \circ	0.0661 \circ	0.7515 \circ
	GRU-BS + Trans. Reranker (Lu et al., 2024)	57.26% \circ	0.7029 \circ	0.1646 \circ	0.0658 \circ	0.7520 \circ
	GRU-IIM (Laine and Aila, 2017)	57.37% \circ	0.7109 \circ	0.1665 \circ	0.0646 \circ	0.7505 \circ
	GRU-DPD (ours)	55.22%	0.7405	0.1734	0.0680	0.7410
	GRU-DPD-IIM (ours)	56.63% \circ	0.7206	0.1688	0.0645 \circ	0.7469
	Trans-IIM (Laine and Aila, 2017)	56.89% \circ	0.7144 \circ	0.1673 \circ	0.0661 \circ	0.7487 \circ
	Trans-DPD (ours)	56.36%	0.7183 \circ	0.1683 \circ	0.0670	0.7479 \circ
Trans-DPD-IIM (ours)	57.63% \circ	0.6967 \circ	0.1632 \circ	0.0631 \bullet	0.7544 \circ	
Rom-phon	GRU-SUPV (Meloni et al., 2021)	51.92%	0.9775	0.1244	0.0390	0.8275
	Trans-SUPV (Kim et al., 2023)	53.04%	0.9050	0.1148	0.0377	0.8417
	GRU-BS + GRU Reranker (Lu et al., 2024)	53.95% \circ	0.8775 \circ	0.0979 \circ	0.0336 \circ	0.8460 \circ
	GRU-BS + Trans. Reranker (Lu et al., 2024)	53.85% \circ	0.8765 \circ	0.0978 \circ	0.0333 \circ	0.8461 \circ
	GRU-IIM (Laine and Aila, 2017)	54.43% \circ	0.9226	0.1029 \circ	0.0388	0.8390
	GRU-DPD (ours)	56.20% \bullet	0.8658 \circ	0.0966 \circ	0.0350 \circ	0.8462 \circ
	GRU-DPD-IIM (ours)	55.68% \bullet	0.8810 \circ	0.0983 \circ	0.0371	0.8453 \circ
	Trans-IIM (Laine and Aila, 2017)	54.95% \circ	0.8864 \circ	0.0989 \circ	0.0379	0.8450
	Trans-DPD (ours)	53.44%	0.9318	0.1039 \circ	0.0391	0.8367
Trans-DPD-IIM (ours)	53.25%	0.8921	0.0995 \circ	0.0376	0.8459 \circ	

Table 15: Supervised reconstruction (100% labeling setting) performance of baseline methods and semisupervised strategies, averaged across 20 or 10 runs (see Section J for detail). Bold: best-performing reconstruction system for the corresponding dataset; \circ : significantly better than both GRU-SUPV (Meloni et al., 2021) and Trans-SUPV (Kim et al., 2023) ($p < 0.01$); \bullet : significantly better than all of GRU-SUPV (Meloni et al., 2021), Trans-SUPV (Kim et al., 2023), GRU-BS + GRU Reranker (Lu et al., 2024), and GRU-BS + Trans. Reranker (Lu et al., 2024) ($p < 0.01$).

Performance distribution by group for 10% labeled WikiHan

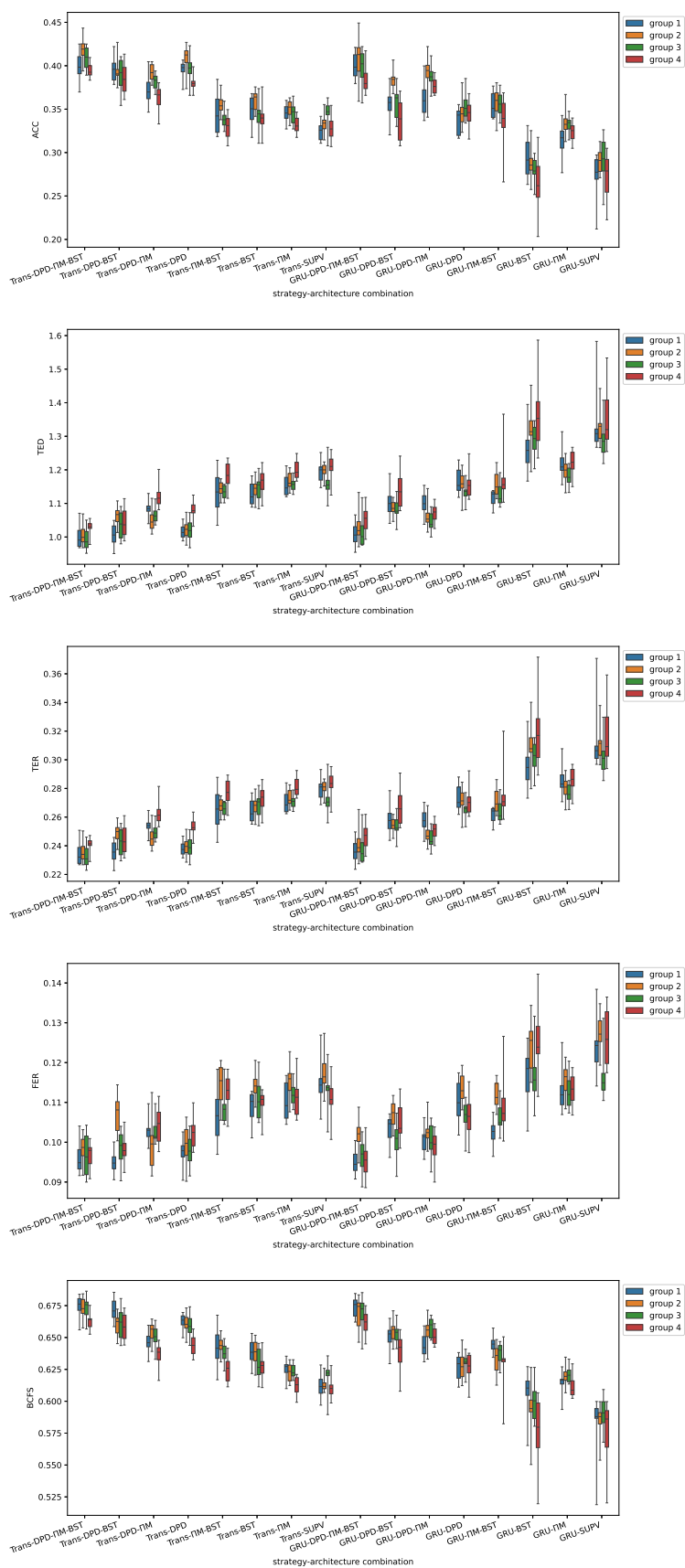


Figure 6: Box plots showing performance distribution for each metric and for each group on 10% labeled WikiHan.

Performance distribution by group for 10% labeled Rom-phon

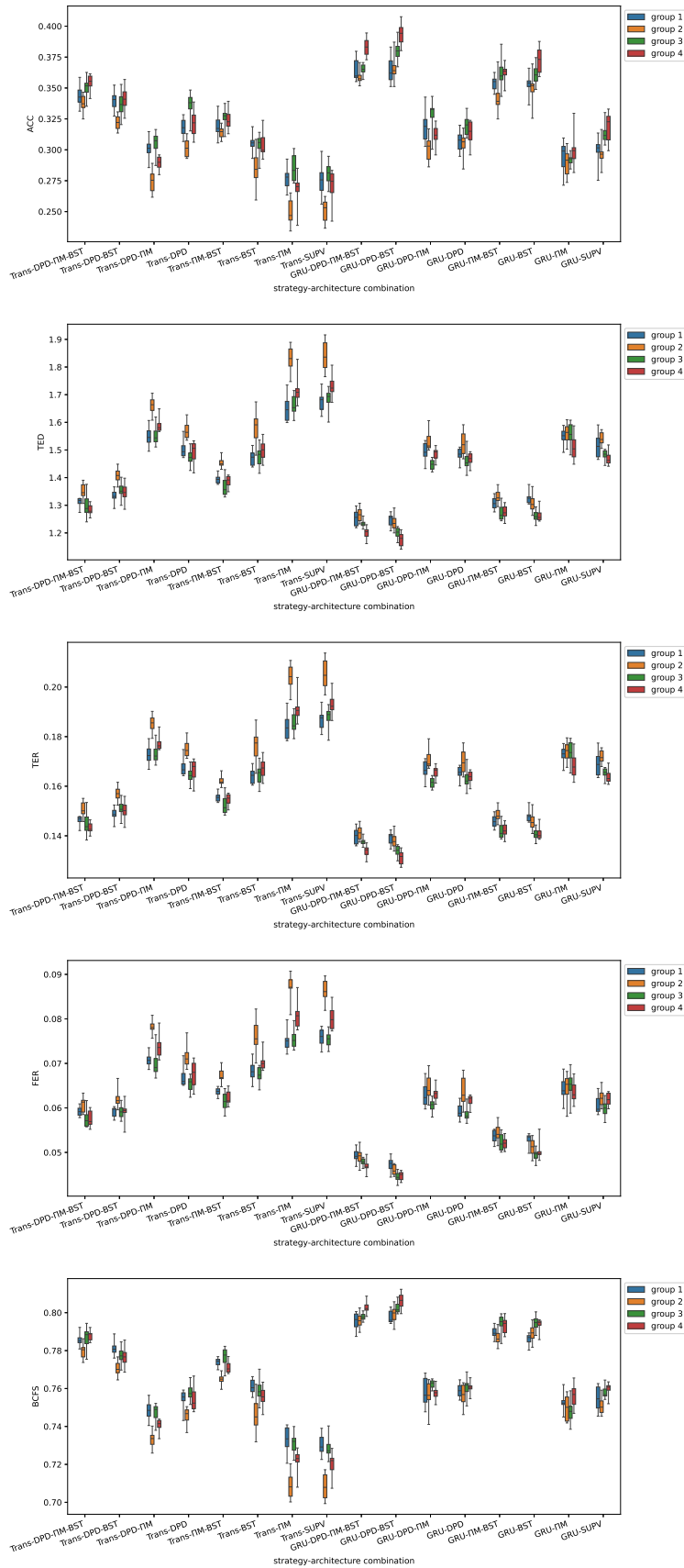


Figure 7: Box plots showing performance distribution for each metric and for each group on 10% labeled Rom-phon.

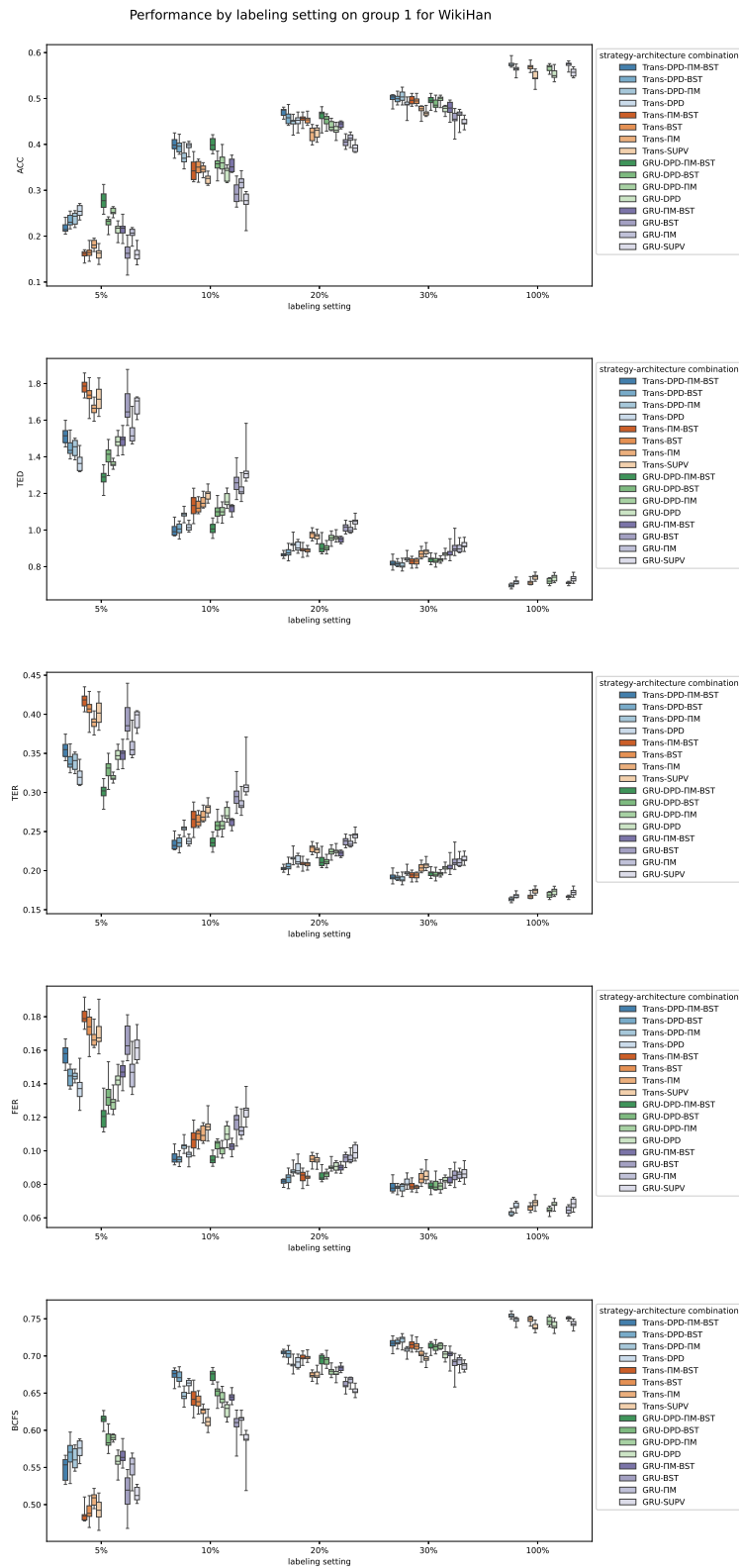


Figure 8: Box plots showing performance distribution for each metric given varied percentages of labels on WikiHan group 1.

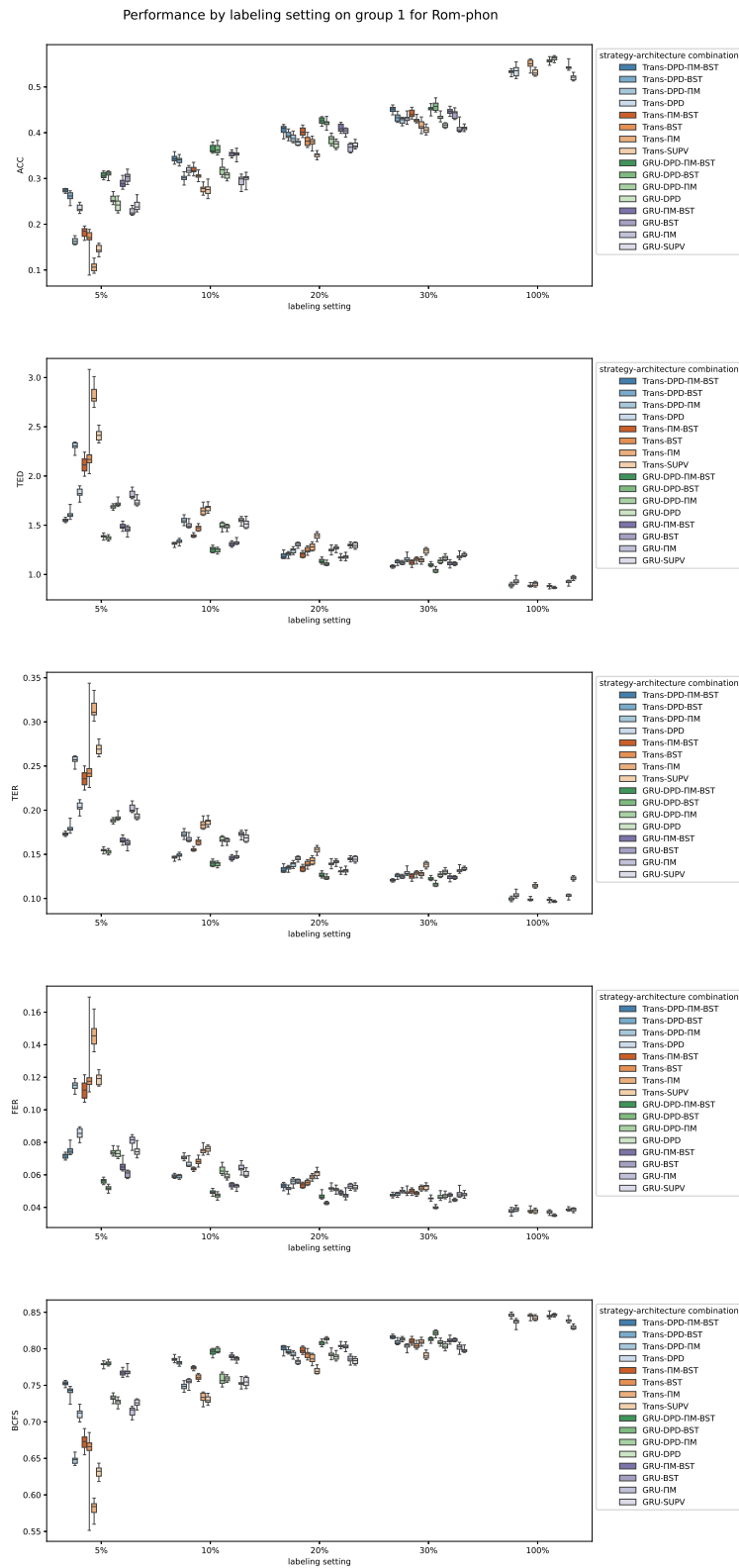
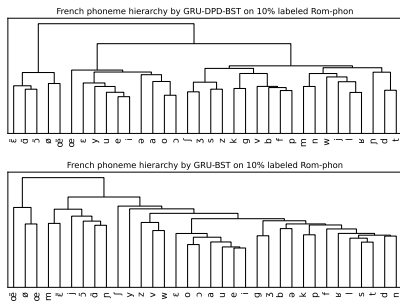
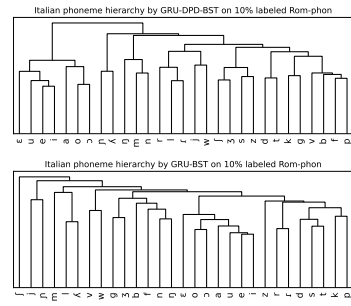


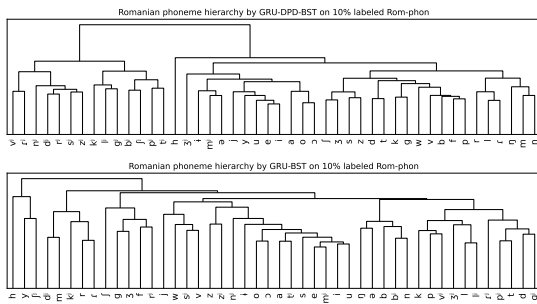
Figure 9: Box plots showing performance distribution for each metric given varied percentages of labels on Rom-phon group 1.



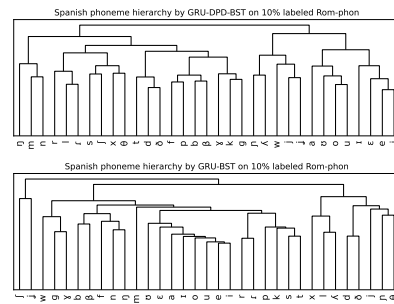
(a) French



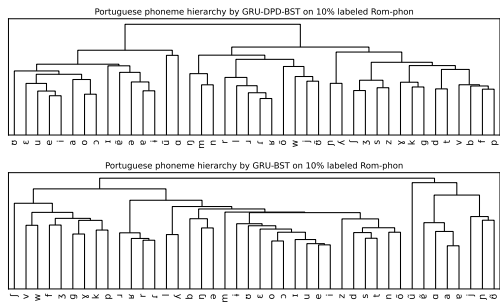
(b) Italian



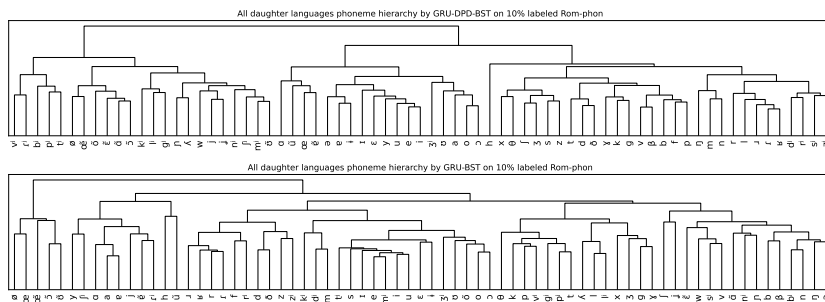
(c) Romanian



(d) Spanish



(e) Portuguese



(f) All daughters

Figure 10: Hierarchical clustering revealing phoneme organization learned by the best run in the best DPD-based strategy-architecture combination (within group 1 and on 10% labeled Rom-phon) (top) and the best run from their non-DPD counterpart (bottom). Note that a comparison for Latin is not possible since non-DPD strategies do not learn embeddings for Latin phonemes.

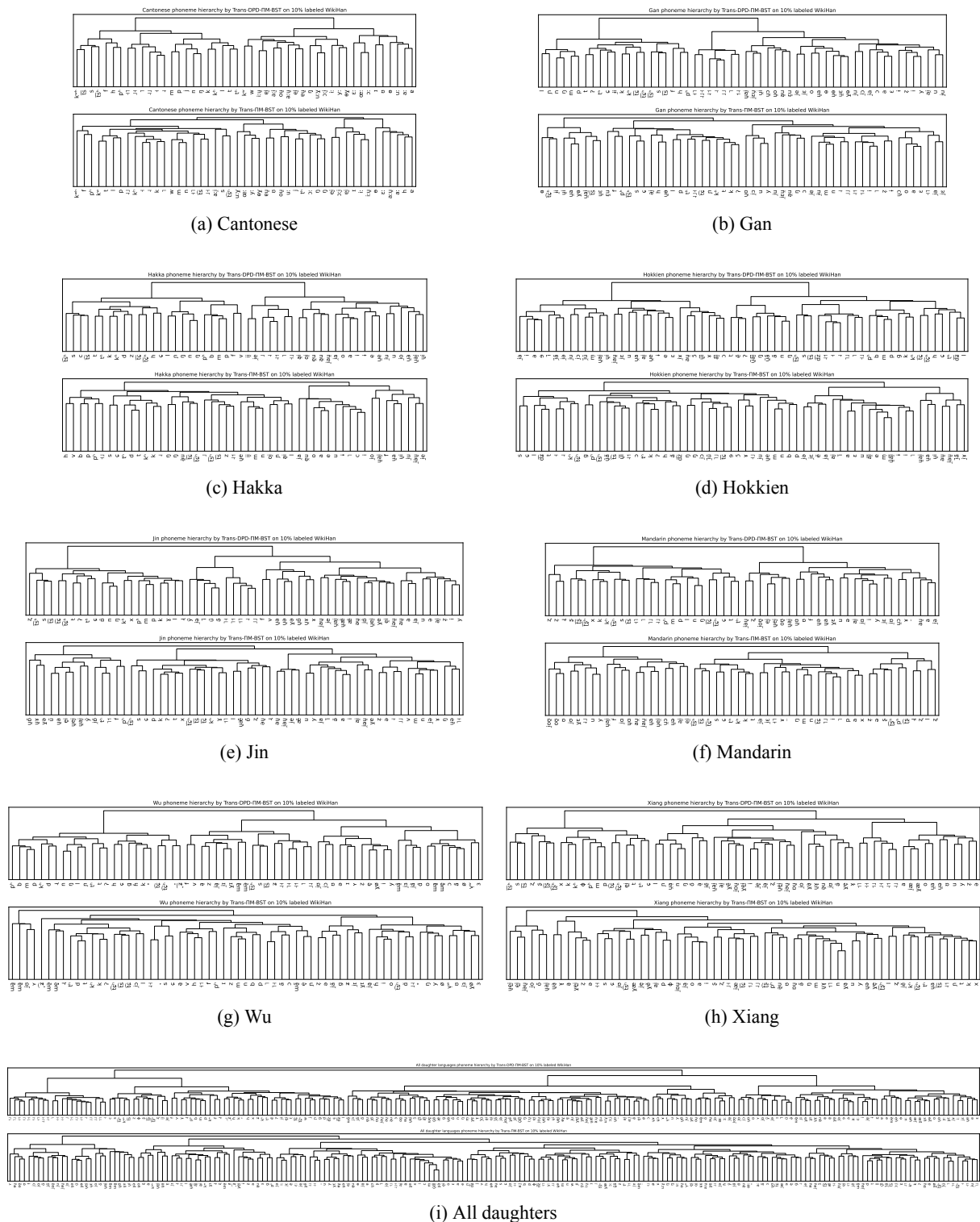


Figure 11: Hierarchical clustering revealing phoneme organization learned by the best run in the best DPD-based strategy-architecture combination (within group 1 and on 10% labeled WikiHan) (top) and the best run from their non-DPD counterpart (bottom). Note that a comparison for Middle Chinese is not possible since non-DPD strategies do not learn embeddings for Middle Chinese phonemes.

	GRU-SUPV	GRU-LIM	GRU-BST	GRU-LIM-BST	GRU-DPD	GRU-DPD-LIM	GRU-DPD-BST	GRU-DPD-LIM-BST	Trans-SUPV	Trans-LIM	Trans-BST	Trans-LIM-BST	Trans-DPD	Trans-DPD-LIM	Trans-DPD-BST	Trans-DPD-LIM-BST
batch size	128	64	128	64	128	128	64	64	128	128	256	256	64	256	64	64
max epochs	221	247	238	253	345	268	384	253	205	288	374	341	261	256	390	259
warmup epochs (learning rate)	29	12	22	16	37	10	2	27	6	6	17	26	18	7	31	22
β_1 (Adam)	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000
β_2 (Adam)	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990
ϵ (Adam)	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08
learning rate	0.0009315	0.0007376	0.0008285	0.0005646	0.0007678	0.0009974	0.0008380	0.0008118	0.0007602	0.0005026	0.0009074	0.0008765	0.0007032	0.0009288	0.0006642	0.0008475
weight decay (Adam)	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.000e-07	0.000e-07	0.000e-07	0.0000	0.0000	0.0000
DPD protoform reconstruction weight					1.100	0.7723	0.6766	0.5396					0.6270	0.5239	1.338	0.6952
DPD bridge network weight					0.3119	0.06223	0.6173	0.3858					0.5003	0.5839	0.1296	0.2182
DPD reflex prediction from gold protoform weight					0.04734	0.5011	0.09115	0.6411					0.6265	0.4270	0.4811	0.2464
DPD reflex prediction from reconstruction weight					1.498	1.249	0.9654	1.468					1.196	1.475	0.8612	1.322
DPD CRINCE loss weight α					0.7015	0.3820	0.5720	0.3732					0.3769	0.2332	0.6782	0.3056
DPD CRINCE loss top k					3	4	1	3					1	4	3	3
DPD shared embedding size					128	256	128	256					128	128	256	256
D2P encoder layers count	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
D2P dropout	0.3117	0.2189	0.2336	0.3137	0.3228	0.2997	0.1999	0.2471	0.2539	0.1951	0.1949	0.1877	0.1594	0.2649	0.1574	0.1692
D2P inference decode max length	15	15	15	15	15	15	15	15	15	15	15	15	15	15	15	15
D2P feedforward dimension	512	512	512	512	512	512	512	512	512	512	512	512	512	512	512	512
D2P embedding size	384	128	128	128	64	64	64	64	256	384	256	128	512	512	512	512
D2P model size	128	64	128	128	64	64	64	64	256	384	256	128	512	512	512	512
D2P encoder layers count																
D2P number of heads																
D2P max input length																
D2P encoder layers count																
P2D dropout					2	2	2	2					2	2	2	2
P2D inference decode max length					0.1728	0.3136	0.3145	0.1546					0.3374	0.2532	0.1752	0.3339
P2D feedforward dimension					15	15	15	15					15	15	15	15
P2D embedding size					512	512	512	512					512	512	512	512
P2D model size					128	64	128	128					128	128	128	128
P2D encoder layers count																
P2D number of heads																
P2D max input length																
Bootstrapping starting epoch																
Bootstrapping log probability threshold																
Bootstrapping max new pseudo-labels per epoch																
Hi-model consistency ramp-up epochs	1	38	75	7	87	11	29	24	49	16	4	79	86	25	86	27
Hi-model max consistency scaling	382.0	-0.006900	-0.008731	222.1	181.7	223.3	-0.007531	-0.005348	181.7	-0.002607	-0.009312	-0.003975	-0.007042	207.5	-0.003975	88.70

Table 16: Hyperparameters for semisupervised reconstruction experiments on WikiHAn.

	GRU-SUPV	GRU-LIM	GRU-BST	GRU-LIM-BST	GRU-DPDP	GRU-DPDP-LIM	GRU-DPDP-BST	GRU-DPDP-LIM-BST	Trans-SUPV	Trans-LIM	Trans-BST	Trans-LIM-BST	Trans-DPDP	Trans-DPDP-LIM	Trans-DPDP-BST	Trans-DPDP-LIM-BST
batch size	256	256	64	64	256	256	128	128	256	128	256	128	256	256	128	128
max epochs	316	283	289	383	224	334	371	373	219	266	382	344	305	206	257	216
warmup epochs (learning rate)	39	15	20	29	34	29	26	13	33	4	36	37	37	29	10	32
β_1 (Adam)	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000
β_2 (Adam)	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990
ϵ (Adam)	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08
learning rate	0.0006089	0.0005475	0.0007060	0.0009183	0.0006284	0.0005706	0.0008698	0.0006257	0.0007210	0.0005128	0.0009196	0.0006014	0.0005052	0.0006181	0.0005023	0.0008986
weight decay (Adam)	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
DPPD protoform reconstruction weight					1.358	0.5163	0.6345	0.8579					0.0000	1.000e-07		1.000e-07
DPPD bridge network weight					0.4293	0.3316	0.5736	0.5268					0.6164	1.033	0.5132	0.5674
DPPD reflex prediction from gold protoform weight					0.6922	0.6284	0.1439	0.4346					0.7455	0.4612	0.6321	0.5607
DPPD reflex prediction from reconstruction weight					1.317	0.9957	1.476	0.9865					0.2108	0.5989	0.3991	0.3401
DPPD CRINCE loss weight α					0.02017	0.002749	0.5966	0.5339					1.201	0.8703	1.294	1.421
DPPD CRINCE loss top k					5	1	5	5					0.03967	0.3295	0.6362	0.02023
DPPD shared embedding size					256	384	384	384					5	384	4	1
D2P encoder layers count	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
D2P dropout	0.3342	0.2129	0.2680	0.2261	0.3300	0.1774	0.1970	0.2402	0.2838	0.3484	0.1963	0.2587	0.2667	0.3453	0.2627	0.2551
D2P inference decode max length	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30
D2P feedforward dimension	512	512	512	512	512	512	512	512	512	512	512	512	512	512	512	512
D2P embedding size	384	128	256	128	64	64	128	128	256	128	384	384	512	512	512	512
D2P model size	128	128	64	64	64	64	128	128	2	2	2	2	2	2	2	2
D2P encoder layers count									2	2	2	2	2	2	2	2
D2P number of heads									8	8	8	8	8	8	8	8
D2P max input length									128	128	128	128	128	128	128	128
P2D encoder layers count					2	2	2	2					2	2	2	2
P2D dropout					0.2214	0.2429	0.3359	0.3346					0.3470	0.3168	0.2199	0.2197
P2D inference decode max length					30	30	30	30					30	30	30	30
P2D feedforward dimension					512	512	512	512					512	512	512	512
P2D embedding size					128	64	64	64					512	512	512	512
P2D model size													2	2	2	2
P2D encoder layers count													2	2	2	2
P2D number of heads													8	8	8	8
P2D max input length													128	128	128	128
Booststrapping starting epoch		37		20												
Booststrapping log probability threshold		-0.005144		-0.009153												
Booststrapping max new pseudo-labels per epoch		51		95					28	-0.003401		16	-0.005465			
H-model consistency ramp-up epochs	18			14		29		86			37		37		48	
H-model max consistency scaling	168.8			168.4		198.2		211.6		2	393.6		248.6		23	
															301.3	
															73	
															46	
															18	
															192.2	

Table 17: Hyperparameters for semisupervised reconstruction experiments on Rom-phon.

	GRU-STUPV	GRU-LIM	GRU-DPD	GRU-DPD-LIM	Trans-STUPV	Trans-LIM	Trans-DPD	Trans-DPD-LIM
batch size	128	128	64	256	128	64	64	64
max epochs	221	206	323	363	205	383	217	363
warmup epochs (learning rate)	29	17	5	11	26	26	35	37
β_1 (Adam)	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000
β_2 (Adam)	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990
ϵ (Adam)	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08
learning rate	0.0009315	0.0008794	0.0005477	0.0009955	0.0007602	0.0009179	0.0008318	0.0009033
weight decay (Adam)	0.000	0.000	0.000	0.000	1.000e-07	1.000e-07	1.000e-07	1.000e-07
DPD protoform reconstruction weight							0.6045	0.5390
DPD bridge network weight			0.01908	0.6016			0.2167	0.3800
DPD reflex prediction from gold protoform weight			0.4224	0.1871			0.1942	0.7438
DPD reflex prediction from reconstruction weight			1.140	0.7006			1.326	1.089
DPD CRINCE loss weight α			0.4218	0.03425			0.1233	0.6634
DPD shared embedding size	2	2	384	256	2	2	384	128
D2P dropout	0.3117	0.3494	0.2859	0.2805	0.2539	0.3113	0.2902	0.2546
D2P inference decode max length	15	15	15	15	15	15	15	15
D2P feedforward dimension	512	512	512	512	512	512	512	512
D2P embedding size	384	128	384	128	256	384	512	512
D2P model size	128	64	64	128	2	2	2	2
D2P encoder layers count					2	2	2	2
D2P number of heads					8	8	8	8
D2P max input length					128	128	128	128
P2D dropout			2	2			2	2
P2D inference decode max length			0.1934	0.2122			0.2663	0.2011
P2D feedforward dimension			15	15			15	15
P2D embedding size			512	512			512	512
P2D model size			64	128			2	2
P2D encoder layers count							2	2
P2D number of heads							8	8
P2D max input length							128	128
Booststrapping starting epoch								
Booststrapping log probability threshold								
Booststrapping max new pseudo-labels per epoch								
I-model consistency ramp-up epochs	11	352.9		24		28		4
I-model max consistency scaling				257.3		202.9		197.1

Table 18: Additional hyperparameters for 10% labeled WikitHan when unlabeled cognate sets are excluded.

Table 19: Additional hyperparameters for 10% labeled Rom-phon when unlabeled cognate sets are excluded.

	GRL- <small>STUPV</small>	GRL- <small>IM</small>	GRL- <small>DPD</small>	GRL- <small>DPD-IM</small>	Trans- <small>STUPV</small>	Trans- <small>IM</small>	Trans- <small>DPD</small>	Trans- <small>DPD-IM</small>
batch size	256	256	128	128	256	64	64	256
max epochs	316	210	357	373	219	355	221	310
warmup epochs (learning rate)	39	17	29	14	33	31	22	33
β_1 (Adam)	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000
β_2 (Adam)	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990
ϵ (Adam)	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08
learning rate	0.00060689	0.0009945	0.0006927	0.0005954	0.0007210	0.0007154	0.0007520	0.0009266
weight decay (Adam)	0.000	0.000	0.000	0.000	1.000e-07	0.000	0.000	0.000
DPD protoform reconstruction weight			0.9617	0.8749			0.6145	0.7346
DPD bridge network weight			0.07922	0.5155			0.2638	0.4678
DPD reflex prediction from gold protoform weight			0.2101	0.6491			0.5297	0.1037
DPD reflex prediction from reconstruction weight			1.248	0.6504			0.9362	1.036
DPD CRINCE loss weight α			0.6470	0.2729			0.6543	0.1272
DPD shared embedding size			1	2			3	4
D2P encoder layers count			384	384			384	384
D2P dropout	2	2	2	2	2	2	2	2
D2P inference decode max length	0.3342	0.2882	0.3309	0.2913	0.2838	0.3305	0.2617	0.3448
D2P feedforward dimension	30	30	30	30	30	30	30	30
D2P embedding size	512	512	512	512	512	512	512	512
D2P model size	384	128	384	64	256	384	384	512
D2P encoder layers count								
D2P number of heads						2	2	2
D2P max input length						8	8	8
D2P encoder layers count						128	128	128
P2D dropout			2	2			2	2
P2D inference decode max length			0.3384	0.2994			0.3333	0.2712
P2D feedforward dimension			30	30			30	30
P2D embedding size			512	512			512	512
P2D model size			64	128			2	2
P2D encoder layers count							2	2
P2D number of heads							8	8
P2D max input length							128	128
Booststrapping starting epoch								
Booststrapping log probability threshold								
Booststrapping max new pseudo-labels per epoch		21				24		6
II-model consistency ramp-up epochs		266.8				370.4		334.1
II-model max consistency scaling					145.5			

	GRU-IJM	GRU-DPD	GRU-DPD-IJM	Trans-IJM	Trans-DPD	Trans-DPD-IJM
batch size	236	64	256	128	128	64
max epochs	283	269	275	382	232	315
warmup epochs (learning rate)	4	7	4	22	40	22
β_1 (Adam)	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000
β_2 (Adam)	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990
ϵ (Adam)	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08
learning rate	0.0006593	0.0005142	0.0008696	0.0007581	0.0007199	0.0007997
weight decay (Adam)	0.000	0.000	0.000	0.000	0.000	1.000e-07
DPD prototom reconstruction weight	0.8200	0.9314		0.5602	1.427	
DPD bridge network weight	0.3285	0.5967		0.4405	0.5442	
DPD reflex prediction from gold prototom weight	0.5770	0.3550		0.5421	0.4134	
DPD reflex prediction from reconstruction weight	0.8717	0.6495		1.009	0.8214	
DPD CRINGE loss weight α	0.3174	0.2638		0.001917	0.2851	
DPD CRINGE loss top k	5	3		4	4	
DPD shared embedding size	128	128		128	128	
D2P encoder layers count	2	2		2	2	
D2P dropout	0.2404	0.2346	0.1761	0.1685	0.2935	0.3055
D2P inference decode max length	15	15	15	15	15	15
D2P feedforward dimension	512	512	512	512	512	512
D2P embedding size	128			256		
D2P model size	128	64	128			
D2P encoder layers count				2	2	2
D2P number of heads				8	8	8
D2P max input length				128	128	128
P2D encoder layers count	2	2		2	2	
P2D dropout	0.1692	0.3488		0.3083	0.1804	
P2D inference decode max length	15	15	15	15	15	15
P2D feedforward dimension	512	512	512	512	512	512
P2D embedding size		128	64			
P2D model size				2	2	2
P2D encoder layers count				8	8	8
P2D number of heads				128	128	128
P2D max input length						
Bootstrapping starting epoch						
Bootstrapping log probability threshold						
Bootstrapping max new pseudo-labels per epoch	4			14		14
l1-model consistency ramp-up epochs	257.4		22	200.2		142.0
l1-model max consistency scaling			156.4			

Table 20: Additional hyperparameters for 100% labeled WikiHan.

	GRU-LIM	GRU-DPD	GRU-DPD-LIM	Trans-LIM	Trans-DPD	Trans-DPD-LIM
batch size	256	128	64	64	64	128
max epochs	379	309	307	295	218	365
warmup epochs (learning rate)	39	37	7	33	13	3
β_1 (Adam)	0.9000	0.9000	0.9000	0.9000	0.9000	0.9000
β_2 (Adam)	0.9990	0.9990	0.9990	0.9990	0.9990	0.9990
ϵ (Adam)	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08	1.000e-08
learning rate	0.0008681	0.0008735	0.0005797	0.0006855	0.0006586	0.0006777
weight decay (Adam)	0.000	0.000	0.000	1.000e-07	0.000	1.000e-07
DPD protoform reconstruction weight		0.5814	1.184		0.5771	1.473
DPD bridge network weight		0.7379	0.5496		0.5159	0.1906
DPD reflex prediction from gold protoform weight		0.3473	0.2647		0.3375	0.5105
DPD reflex prediction from reconstruction weight		1.147	0.5560		0.5103	1.405
DPD CRINGE loss weight α		0.1322	0.5651		0.4757	0.09542
DPD CRINGE loss top k		4	4		5	1
DPD shared embedding size		256	384		384	256
D2P dropout	2	2	2	2	2	2
D2P inference decode max length	0.2888	0.2923	0.3016	0.1566	0.2144	0.2418
D2P feedforward dimension	30	30	30	30	30	30
D2P embedding size	512	512	512	512	512	512
D2P model size	128	128	128	128	128	128
D2P encoder layers count				2	2	2
D2P number of heads				8	8	8
D2P max input length				128	128	128
P2D dropout	2	2	2	2	2	2
P2D inference decode max length	0.1589	0.2117		0.1996	0.2921	
P2D feedforward dimension	30	30	30	30	30	30
P2D embedding size	512	512	512	512	512	512
P2D model size	64	64	64	2	2	2
P2D encoder layers count				8	8	8
P2D number of heads				128	128	128
P2D max input length						
Bootstrapping starting epoch						
Bootstrapping log probability threshold	9	15	30	29		
Bootstrapping max new pseudo-labels per epoch	83.25	92.97	67.85	64.89		
L1-model consistency ramp-up epochs						
L1-model max consistency scaling						

Table 21: Additional hyperparameters for 100% labeled Rom-phon.

	Semisupervised		10%, exclude unlabeled		100%	
	WikiHan	Rom-phon	WikiHan	Rom-phon	WikiHan	Rom-phon
GRU-SUPV	1,724,288	1,605,184	1,724,288	1,605,184		
GRU-IIM	1,231,232	1,605,184	1,231,232	1,605,184	1,724,288	1,605,184
GRU-DPD	2,851,584	3,036,416	3,130,112	2,929,024	2,851,584	3,101,952
GRU-DPD-IIM	2,754,560	2,929,024	3,748,864	3,450,752	2,917,120	3,516,288
Trans-SUPV	3,938,917	3,846,935	3,938,917	3,846,935		
Trans-IIM	3,938,917	3,846,935	3,938,917	3,846,935	3,938,917	3,846,935
Trans-DPD	6,373,066	11,488,686	11,653,706	11,488,686	2,665,290	11,488,686
Trans-DPD-IIM	2,665,290	11,488,686	2,665,290	11,488,686	2,665,290	6,262,958
GRU-BST	1,724,288	1,112,128				
GRU-IIM-BST	1,724,288	1,112,128				
GRU-DPD-BST	2,851,584	3,516,288				
GRU-DPD-IIM-BST	3,218,944	3,516,288				
Trans-BST	3,938,917	3,846,935				
Trans-IIM-BST	3,938,917	3,846,935				
Trans-DPD-BST	6,373,066	11,488,686				
Trans-DPD-IIM-BST	6,373,066	6,262,958				

Table 22: Trainable parameter count for each architecture-strategy combination in each experiment setup.

Dataset	Category	Count	Percentage
WikiHan	Correct at all %	140	13.55%
	Correct only above % threshold	359	34.75%
	Correct only below % threshold	37	3.58%
	Incorrect at all %	278	26.91%
	Other pattern	219	21.20%
Rom-phon	Correct at all %	359	20.47%
	Correct only above % threshold	451	25.71%
	Correct only below % threshold	65	3.71%
	Incorrect at all %	586	33.41%
	Other pattern	293	16.70%

Table 23: Distribution of reconstruction correctness patterns across labeling setting for Trans-DPD-IIM-BST on group 1 WikiHan and GRU-DPD-BST on group 1 Rom-phon.

Category	\hat{y} @ 5%	\hat{y} @ 10%	\hat{y} @ 20%	\hat{y} @ 30%	\hat{y} @ 100%	y (Reference)	
Correct at all %	bjwot 入 dzoj 平 jen 平 k^hwa 平 luk 入 mjun 平 teuw 平 xjuw 去	bjwot 入 dzoj 平 jen 平 k^hwa 平 luk 入 mjun 平 teuw 平 xjuw 去	bjwot 入 dzoj 平 jen 平 k^hwa 平 luk 入 mjun 平 teuw 平 xjuw 去	bjwot 入 dzoj 平 jen 平 k^hwa 平 luk 入 mjun 平 teuw 平 xjuw 去	bjwot 入 dzoj 平 jen 平 k^hwa 平 luk 入 mjun 平 teuw 平 xjuw 去	bjwot 入 dzoj 平 jen 平 k^hwa 平 luk 入 mjun 平 teuw 平 xjuw 去	bjwot 入 dzoj 平 jen 平 k^hwa 平 luk 入 mjun 平 teuw 平 xjuw 去
Correct only above % threshold	bjuw去 bjwin平 laŋ去 maw上 jew上 bjwin平 sjuw去 dzaj去 sen去 nuŋ ^w 平 kwak入 mat入 yəj平 djuŋ ^w 平 fʂ ^h jak入 gij平 kwan去 k ^h iŋ平 mam平 teaŋ平 juŋ ^w 平	bju 去 bjwon 平 ljaŋ 去 mæw 上 new 上 pjun 平 su 去 tsoj 去 zjen 去 joŋ^w 平 kwak入 mat入 yoj平 djuŋ ^w 平 djwak入 gij平 kwan去 k ^h oŋ平 kom平 tiŋ平 joŋ ^w 平	bju 去 bjwon 平 ljaŋ 去 mæw 上 new 上 pjun 平 su 去 tsoj 去 zjen 去 joŋ^w 平 kæk ^w 入 mwat入 yej平 fʂ ^h uŋ ^w 平 teak入 gij平 kwen去 k ^h oŋ平 yam平 teŋ平 joŋ ^w 平	bju 去 bjwon 平 ljaŋ 去 mæw 上 new 上 pjun 平 su 去 tsoj 去 zjen 去 joŋ^w 平 kæk ^w 入 mwat入 yej平 fʂ ^h uŋ ^w 平 te ^h ak入 gij平 kwen去 k ^h æŋ平 yam平 dʒæŋ平 joŋ ^w 平	bju 去 bjwon 平 ljaŋ 去 mæw 上 new 上 pjun 平 su 去 tsoj 去 zjen 去 joŋ^w 平 kæk ^w 入 mwat入 yej平 fʂ ^h uŋ ^w 平 te ^h ak入 gi平 k ^h wen去 k ^h ɛŋ平 yom平 fʂɛŋ平 ʔjoŋ ^w 平	bju 去 bjwon 平 ljaŋ 去 mæw 上 new 上 pjun 平 su 去 tsoj 去 zjen 去 joŋ^w 平 kæk ^w 入 mwat入 yej平 fʂ ^h uŋ ^w 平 te ^h ak入 gi平 k ^h wen去 k ^h ɛŋ平 yom平 fʂɛŋ平 ʔjoŋ ^w 平	bju 去 bjwon 平 ljaŋ 去 mæw 上 new 上 pjun 平 su 去 tsoj 去 zjen 去 joŋ^w 平 kæk ^w 入 mwat入 yej平 fʂ ^h uŋ ^w 平 te ^h ak入 gi平 k ^h wen去 k ^h ɛŋ平 yom平 fʂɛŋ平 ʔjoŋ ^w 平
Correct only below % threshold	jej 去 kij 去	jej 去 kij 去	?ej去 kij去	?ej去 kij去	?ej去 kij去	jej 去 kij 去	
Incorrect at all %	daw上 fʂ ^h jen平 gjo平 kwaŋ去 k ^h waŋ平 pap入 kep入 tjwan上 t ^h e去 ywi去 jen去 næn去 tei去 sjwa上 te平 ʔaŋ去	daw上 dzem平 gjo平 kwaŋ去 k ^h waŋ平 pæ上 kjeŋ入 tsan上 di去 ywej去 njen去 nep入 tee去 swa上 ts ^h i平 ʔaŋ去	daw上 dzem平 gjo平 ywaŋ去 k ^h waŋ平 pæ上 kjeŋ入 tsen上 ts ^h i去 ywej去 njen去 dzep入 di去 swa上 ts ^h i平 ʔaŋ去	daw上 dzem平 gjo平 gwaŋ去 k ^h waŋ平 pæ上 tsjep入 tsan上 ts ^h i去 ywej去 njen去 nep入 di平 swa上 te ^h i平 ʔaŋ去	daw上 dzem平 gjo平 kwaŋ去 k ^h jaŋ平 pæ上 kjeŋ入 tsan上 ts ^h je去 ywej去 njen去 nep入 dij去 swa上 ts ^h i平 ʔaŋ去	daw去 dzen平 gju平 k ^h waŋ上 k ^h jaŋ平 pæ去 tsep入 tswan上 ts ^h ij去 nen去 cep入 dji去 gjo上 fʂ ^h je平 ʔaŋ上	daw去 dzen平 gju平 k ^h waŋ上 k ^h jaŋ平 pæ去 tsep入 tswan上 ts ^h ij去 nen去 cep入 dji去 gjo上 fʂ ^h je平 ʔaŋ上
Other pattern	dan 平 mjew 平 ywæ 平 dzjem 平 pej去 ŋæ 平 yeŋ平 luk ^w 入 kew上 pjeŋ上 kwam平 min上 yjawæn去	don平 mew平 yu平 zjem平 pej去 ŋeŋ平 yeŋ 平 jok^w 入 k ^h æw上 pjeŋ上 k ^h om平 min上 yjen去	dan 平 mjew 平 ywæ 平 dzjem 平 pej去 ŋæ 平 yeŋ 平 juk ^w 入 k ^h aw上 pjaŋ 上 k ^h am平 min上 yjen去	dan 平 mjew 平 ywæ 平 dzjem 平 pej去 ŋæ 平 yeŋ 平 lok入 k ^h æw上 pjeŋ上 k ^h am平 min 平 ywen 去	dan 平 mjew 平 ywæ 平 dzjem 平 pej去 ŋæ 平 yəŋ平 jok^w 入 k ^h aw上 pjaŋ 上 kam平 min上 zjen去	dan 平 mjew 平 ywæ 平 dzjem 平 pej去 ŋæ 平 yeŋ 平 jok^w 入 k ^h aw上 pjaŋ 上 k ^h am平 min 平 ywen 去	dan 平 mjew 平 ywæ 平 dzjem 平 pej去 ŋæ 平 yeŋ 平 jok^w 入 k ^h aw上 pjaŋ 上 k ^h am平 min 平 ywen 去

Table 24: Sample protoform predictions by Trans-DPD-ΠIM-BST for different labeling settings on group 1 WikiHan. Bold: correct protoform; \hat{y} @ {5%, 10%, 20%, 30%, 100%}: protoform prediction when trained with a {5%, 10%, 20%, 30%, 100%} labeling setting; y: gold protoform.

Category	\hat{y} @ 5%	\hat{y} @ 10%	\hat{y} @ 20%	\hat{y} @ 30%	\hat{y} @ 100%	y (Reference)
Correct at all %	balteøm distillationem kakare mørsøm medikamentøm puritatem sanare skapølam taktilem wariabilem wenøstøm wermem	balteøm distillationem kakare mørsøm medikamentøm puritatem sanare skapølam taktilem wariabilem wenøstøm wermem	balteøm distillationem kakare mørsøm medikamentøm puritatem sanare skapølam taktilem wariabilem wenøstøm wermem	balteøm distillationem kakare mørsøm medikamentøm puritatem sanare skapølam taktilem wariabilem wenøstøm wermem	balteøm distillationem kakare mørsøm medikamentøm puritatem sanare skapølam taktilem wariabilem wenøstøm wermem	balteøm distillationem kakare mørsøm medikamentøm puritatem sanare skapølam taktilem wariabilem wenøstøm wermem
Correct only above % threshold	armipotentem enøntiare infantiam sørdøm indikativøm ekwabilitatem diskretionem plenarøm prærogatiwam preskribere dilatorøm feritam ekstrarre ekwiwokare konsensøm iŋkørabilem	armipotentem enøntiare infantiam sørdøm indikativøm ekwabilitatem diskretionem plenarøm prærogatiwam preskribere dilatorøm feritam ekstrarre ekwiwokare konsensøm iŋkørabilem	armipotentem enøntiare infantiam sørdøm indikativøm arkwabilitatem diskretionem plenarøm prærogatiwam praeskribere dilatorøm feritam ekstrarre ekwiwokare konsensøm iŋkørabilem	armipotentem enøntiare infantiam sørdøm indikativøm aikwabilitatem diskretionem plenarøm prærogatiwam praeskribere dilatorøm feritam ekstrarre ekwiwokare konsensøm iŋkørabilem	armipotentem enøntiare infantiam sørdøm indikativøm aikwabilitatem diskretionem plenarøm prærogatiwam praeskribere dilatorøm feritam ekstrarre ekwiwokare konsensøm iŋkørabilem	armipotentem enøntiare infantiam sørdøm indikativøm arkwabilitatem diskretionem plenarøm prærogatiwam praeskribere dilatorøm feritam ekstrarre ekwiwokare konsensøm iŋkørabilem
Correct only below % threshold	illustrare kernere	illustrare kernere	illustrare kernere	illøstrare kernere	illøstrare kernere	illustrare kernere
Incorrect at all %	abissøm adesionem allegere frotsetøm ipnotikøm kyrkøndurre køwilem maritale nukleare parire proklamationem pelegrinøm ekswerkøm skøltørem søllatøm stanøm settennem termitem tristetiam gøwernaiem	abissøm adesionem allegere frotsetøm ipnotikøm kyrkøndurre købilem maritale nøkleare parire proklamationem pelegrinøm ekswerkøm skøltørem søllakøm stanøm settennem termitem tristetiam gøbernam	abissøm adesionem allegere frotsetøm hymnotikøm kyrkøndurre køwilem maritale nøkleare parire proklamationem pelerrinøm ekswerkøm skøltørem søllakøm stagnatem staneøm tristetiam gøwernatem	abissøm adesionem allegere frotsetøm hybnøtikøm kyrkøndurre køwilem maritale nukleare parire proklamationem pelegrinøm skøertøm skøltørem søllakøm staneøm staneøm tristetiam gøwernatem	abissøm adesionem alligere frotketøm hymnotikøm kyrkøndurem købilem maritale nukleare parire proklamationem pelegrinøm ekskodere skøltørem søllakøm stannøm sektennem termitem trystektiam gøbernatikøm	abysøm adhesionem allegare frotketøm hypnotikøm kyrkømdukere købile maritale nuklearem parere proklamationem peregrinøm skørtøm skølpørem sølakøm stagnøm septennem tarmitem tristetiam gøbernakølem
Other pattern	tepørem wadøm immaterialem bustøm iŋkwilnøm eram køstitutionem mølkere petatem empirikøm	tepørem wadøm immaterialem bustøm iŋkwilnøm eram køstitutionem mølkere petasøm empirikøm	tepørem wadøm imaterialem bustøm iŋkwilnøm eram konstitutionem mølkere petasem empirikøm	tepørem wadøm imaterialem bustøm iŋkwilnøm aram konstitutionem mølkere petatem empirikøm	tepørem wadøm immaterialem bustøm iŋkwilnøm heram konstitutionem mølkere pestas empyriskøm	tepørem wadøm immaterialem bustøm iŋkwilnøm aram køstitutionem mølkere petasøm empirikøm

Table 25: Sample protoform predictions by GRU-DPD-BST for different labeling settings on group 1 Rom-phon. Bold: correct protoform; \hat{y} @ {5%, 10%, 20%, 30%, 100%}: protoform prediction when trained with a {5%, 10%, 20%, 30%, 100%} labeling setting; y : gold protoform.

		Trans-IIM-BST	
		correct	incorrect
Trans-DPD-IIM-BST	correct	334 (32.33%)	112 (10.84%)
	incorrect	69 (6.68%)	518 (50.15%)

		GRU-BST	
		correct	incorrect
GRU-DPD-BST	correct	546 (31.13%)	126 (7.18%)
	incorrect	96 (5.47%)	986 (56.21%)

Table 26: Confusion matrix of protoform prediction correctness for Trans-DPD-IIM-BST vs. Trans-IIM-BST on WikiHan (top) and GRU-DPD-BST vs. GRU-BST on Rom-phon (bottom)

Romanian	French	Italian	Spanish	Portuguese	Latin	Latin _{non-DPP}	Latin _{DPP}	Romanian _{DPP}	French _{DPP}	Italian _{DPP}	Spanish _{DPP}	Portuguese _{DPP}
-	adikti	addetto	añiċto	adito	addictum	addictum	addictum	-	adikti	addetto	añiċto	adito
-	apădde	appendere	apender	ardor	appendere	appendere	ardorem	-	apădde	ardore	ardore	ardore
ardare	ardore	ardore	ardor	ardor	ardorem	ardorem	ardorem	ardore	ardore	ardore	ardore	ardore
-	arinate	arinate	arinate	arinate	arinate	arinate	arinate	-	arinate	arinate	arinate	arinate
-	bestiar	bestiaro	bestiaro	bestiaro	bestiarum	bestiarum	bestiarum	bestiaro	bestiar	bestiaro	bestiaro	bestiaro
bovin	bovine	bovino	bovino	bovino	bovium	bovium	bovium	bovin	bovine	bovino	bovino	bovino
-	debl	deblie	deblie	deblie	deblie	deblie	deblie	deblie	debl	deblie	deblie	deblie
-	delikats	delikatsa	delinkwenta	delinkwenta	delinkwenta	delinkwenta	delinkwenta	-	delikats	delikatsa	delinkwenta	delinkwenta
-	fluid	dilidgere	fluido	fluid	dilidgere	dilidgere	dilidgere	fluid	fluid	fluid	fluid	fluid
-	flyid	flyid	flyid	flyid	flyid	flyid	flyid	flyid	flyid	flyid	flyid	flyid
-	ŷleapa	ŷleapă	ŷleapă	ŷleapă	ŷleapă	ŷleapă	ŷleapă	ŷleapă	ŷleapă	ŷleapă	ŷleapă	ŷleapă
-	kwadrilater	kwadrilatero	kwadrilatero	kwadrilatero	kwadrilaterum	kwadrilaterum	kwadrilaterum	-	kwadrilater	kwadrilatero	kwadrilatero	kwadrilatero
-	-	-	-	-	-	-	-	-	-	-	-	-
-	konsubstansjal	konsustansjale	konsustansjal	konsustansjal	konsubstansjal	konsubstansjal	konsubstansjal	konsubstansjal	konsubstansjal	konsubstansjal	konsubstansjal	konsubstansjal
kunmat	kupa	kopato	kupato	kupato	kupatum	kupatum	kupatum	kunmat	kupa	kopato	kopato	kopato
labli	labli	labile	labile	labile	labile	labile	labile	labile	labli	labile	labile	labile
-	mitigare	mitigare	mitiyar	mitiyar	mitigare	mitigare	mitigare	mitigare	mitigare	mitigare	mitiyar	mitiyar
-	nodus	mendafle	nodoso	murdasi	mendekem	mendekem	mendekem	nodus	nodus	nodus	nodoso	nodoso
nota	nota	nota	nota	nota	notam	notam	notam	nota	nota	nota	nota	nota
orajsje	orajzi	orajsione	orajon	orajsje	orajonem	orajonem	orajonem	orajsje	orajzi	orajsione	orajon	orajon
prefektura	prefektur	prefektura	prefektura	prefektura	prefekturam	prefekturam	prefekturam	prefektura	prefektur	prefektura	prefektura	prefektura
-	presele	preledere	preleber	preleber	preledere	preledere	preledere	presele	presele	preledere	preleber	preleber
profes	profeso	profeso	profeso	profeso	profesum	profesum	profesum	profes	profeso	profeso	profeso	profeso
purifika	purifike	purifikare	purifikar	purifika	purifikare	purifikare	purifikare	purifika	purifike	purifikare	purifikar	purifikar
pekuliv	pekyl	pekulio	pekuljo	pekuliv	pekuliam	pekuliam	pekuliam	pekuliv	pekyl	pekulio	pekuljo	pekuljo
-	-	-	-	-	-	-	-	-	-	-	-	-
-	sanasji	sanajsione	sanajon	sanajsje	sanajonem	sanajonem	sanajonem	sanasji	sanasji	sanajsione	sanajon	sanajon
sartfina	sartfina	sartfina	sartfina	sartfina	sartfina	sartfina	sartfina	sartfina	sartfina	sartfina	sartfina	sartfina
-	-	-	-	-	-	-	-	-	-	-	-	-
-	servivud	servitute	servitute	servitute	servitutum	servitutum	servitutum	servivud	servivud	servitute	servitute	servitute
-	triseteta	triseteta	triseteta	triseteta	triseteta	triseteta	triseteta	triseteta	triseteta	triseteta	triseteta	triseteta
tub	tubo	tubo	tubo	tubo	tubum	tubum	tubum	tub	tubo	tubo	tubo	tubo
venjal	venjal	venjale	venjal	venjal	venalem	venalem	venalem	venjal	venjal	venjale	venjal	venjal
-	verdigino	verdiginoso	verdiginoso	verdiginoso	verdiginosum	verdiginosum	verdiginosum	verdigino	verdigino	verdiginoso	verdiginoso	verdiginoso
vaduva	vaduva	vedova	vedova	vedova	veduam	veduam	veduam	veduwa	veduwa	veduwa	veduwa	veduwa
ivortiv	ivortiv	avonio	avonio	avonio	avonium	avonium	avonium	evroek	evroek	avonia	avonia	avonia
-	ekspyrgje	spurgare	ekspuryar	ekspuryar	ekspurgare	ekspurgare	ekspurgare	ekspyrgje	ekspyrgje	spurgare	ekspuryar	ekspuryar
-	ekskaksasji	ekskaksasione	ekskaksasione	ekskaksasione	ekskaksationem	ekskaksationem	ekskaksationem	ekskaksasji	ekskaksasji	ekskaksasione	ekskaksasione	ekskaksasione
episkopat	episkopa	episkopato	episkopato	episkopato	episkopatum	episkopatum	episkopatum	episkopat	episkopa	episkopato	episkopato	episkopato
gladium	gliev	gladio	gladio	gladio	gladium	gladium	gladium	gladium	gliev	gladio	gladio	gladio
grasjos	grasjo	grasioso	grasioso	grasjos	grasiosum	grasiosum	grasiosum	grasjos	grasjo	grasioso	grasioso	grasioso
-	-	geometria	geometria	geometria	geometrium	geometrium	geometrium	-	-	geometria	geometria	geometria
-	-	malbare	malbare	malbare	malbare	malbare	malbare	-	-	malbare	malbare	malbare
-	-	indagare	indagar	indagare	indagare	indagare	indagare	-	-	indagare	indagar	indagare

Table 28: Sample reconstruction predictions by GRU-DPP-BST (denoted DPP) and GRU-BST (denoted non-DPP) and sample reflex prediction predictions by GRU-DPP-BST (with latent reconstruction as input) for Rom-phon (proportionate sampling according to the confusion matrix in Table 26). **Language**: prediction for **Language**; bold: correct prediction; ‘-’: the dataset does not contain this reflex.

	Cantonese	Gan	Hakka	Hokkien	Jin	Mandarin	Wu	Xiang	Middle Chinese
Reference	ts^ha:mɿ	-	-	ts^hamɿ	-	ts^hanɿ	-	-	dzam ^平
GRU-SUPV									dzæm ^平
GRU-IIM									dzjem ^平
GRU-BST									ts ^h æm ^平
GRU-IIM-BST									dzom ^平
GRU-DPD	ts^ha:mɿ	-	-	ts^hamɿ	-	ts ^h anɿ	-	-	dzæm ^平
GRU-DPD-IIM	ts^ha:mɿ	-	-	tsamɿ	-	ts^hanɿ	-	-	ʃs ^h æm ^平
GRU-DPD-BST	ts^ha:mɿ	-	-	ts^hamɿ	-	ts^hanɿ	-	-	qæm ^平
GRU-DPD-IIM-BST	ts^ha:mɿ	-	-	ts^hamɿ	-	ts^hanɿ	-	-	dzam ^平
Trans-SUPV									ʃs ^h æm ^平
Trans-IIM									dzæm ^平
Trans-BST									ʃs ^h æm ^平
Trans-IIM-BST									dzam ^平
Trans-DPD	ts^ha:mɿ	-	-	ts^hamɿ	-	ts^hanɿ	-	-	ʃs ^h æm ^平
Trans-DPD-IIM	ts^ha:mɿ	-	-	ts^hamɿ	-	ts^hanɿ	-	-	ts ^h am ^平
Trans-DPD-BST	ts^ha:mɿ	-	-	ts^hamɿ	-	ts^hanɿ	-	-	dzam ^平
Trans-DPD-IIM-BST	ts^ha:mɿ	-	-	ts^hamɿ	-	ts^hanɿ	-	-	dzam ^平
Reference	je:kɿ	ioʔɿ	ioʔɿ	ioʔɿ	ioʔɿ	iauɿ	hiaʔɿ	ioɿ	jak 入
GRU-SUPV									jak 入
GRU-IIM									jak 入
GRU-BST									jak 入
GRU-IIM-BST									jak 入
GRU-DPD	je:kɿ	ioʔɿ	ioʔɿ	iaʔɿ	ioʔɿ	yeɿ	hiaʔɿ	ieɿ	jek入
GRU-DPD-IIM	je:kɿ	ioʔɿ	ioʔɿ	ioʔɿ	yoʔɿ	yeɿ	hiaʔɿ	ioɿ	jak 入
GRU-DPD-BST	je:kɿ	ioʔɿ	ioʔɿ	ioʔɿ	ioʔɿ	yeɿ	hiaʔɿ	ioɿ	jak 入
GRU-DPD-IIM-BST	je:kɿ	ioʔɿ	ioʔɿ	ioʔɿ	ioʔɿ	yeɿ	hiaʔɿ	ioɿ	jak 入
Trans-SUPV									jek ^v 入
Trans-IIM									jak 入
Trans-BST									jek入
Trans-IIM-BST									jak 入
Trans-DPD	je:kɿ	ioʔɿ	ioʔɿ	ioʔɿ	yoʔɿ	iauɿ	hiɿʔɿ	ioɿ	jew去
Trans-DPD-IIM	je:kɿ	ioʔɿ	ioʔɿ	ioʔɿ	ioʔɿ	ieɿ	hiaʔɿ	ieɿ	jak 入
Trans-DPD-BST	je:kɿ	ioʔɿ	ioʔɿ	ioʔɿ	ioʔɿ	yeɿ	hiaʔɿ	ioɿ	jak 入
Trans-DPD-IIM-BST	je:kɿ	ioʔɿ	ioʔɿ	iaʔɿ	ioʔɿ	yeɿ	hiaʔɿ	ioɿ	jaw去
Reference	fetɿ	-	futɿ	hutɿ	xuaʔɿ	xuɿ	-	-	xwot 入
GRU-SUPV									ywot入
GRU-IIM									yuw去
GRU-BST									ywot入
GRU-IIM-BST									xjut入
GRU-DPD	fetɿ	-	fatɿ	hutɿ	fiyaʔɿ	fuɿ	-	-	k ^h wot入
GRU-DPD-IIM	fetɿ	-	futɿ	hutɿ	xuəʔɿ	xuɿ	-	-	xwot 入
GRU-DPD-BST	fa:tɿ	-	vatɿ	hɣatɿ	xuaʔɿ	xua	-	-	ywot入
GRU-DPD-IIM-BST	fetɿ	-	fitɿ	hutɿ	xuəʔɿ	xuɿ	-	-	xjut入
Trans-SUPV									xjut入
Trans-IIM									xwot 入
Trans-BST									xjwot入
Trans-IIM-BST									xwot 入
Trans-DPD	fetɿ	-	fatɿ	hɣatɿ	xuaʔɿ	xua	-	-	xwat入
Trans-DPD-IIM	fetɿ	-	futɿ	hutɿ	fiəʔɿ	xuɿ	-	-	xwot 入
Trans-DPD-BST	fetɿ	-	fitɿ	hutɿ	xuəʔɿ	xuɿ	-	-	xwot 入
Trans-DPD-IIM-BST	fetɿ	-	fitɿ	hutɿ	euəʔɿ	euɿ	-	-	xut入
Reference	sa:nɿ	-	-	ts^hanɿ	-	ʃs^hanɿ	-	-	dzɛn ^平
GRU-SUPV									ts ^h æn ^平
GRU-IIM									dzæn ^平
GRU-BST									dzæn ^平
GRU-IIM-BST									dzon ^平
GRU-DPD	sa:nɿ	-	-	sanɿ	-	ʃs ^h anɿ	-	-	dzon ^平
GRU-DPD-IIM	ts ^h a:nɿ	-	-	ts^hanɿ	-	ʃs ^h anɿ	-	-	te ^h an ^平
GRU-DPD-BST	ts ^h a:nɿ	-	-	ts^hanɿ	-	ʃs ^h anɿ	-	-	ʂæn ^平
GRU-DPD-IIM-BST	ts ^h a:nɿ	-	-	ts^hanɿ	-	ʃs ^h anɿ	-	-	dzan ^平
Trans-SUPV									ʃs ^h æn ^平
Trans-IIM									ʃs ^h æn ^平
Trans-BST									ʃs ^h æn ^平
Trans-IIM-BST									dzan ^平
Trans-DPD	ts ^h a:nɿ	-	-	ts^hanɿ	-	ʃs ^h anɿ	-	-	ʃs ^h æn ^平
Trans-DPD-IIM	ts ^h a:nɿ	-	-	ts^hanɿ	-	ʃs ^h anɿ	-	-	ʃs ^h æn ^平
Trans-DPD-BST	sa:nɿ	-	-	ts^hanɿ	-	ʃs^hanɿ	-	-	ʃs ^h æn ^平
Trans-DPD-IIM-BST	ts ^h a:nɿ	-	-	ts^hanɿ	-	ʃs ^h anɿ	-	-	dzæn ^平

Table 29: Example outputs on WikiHan for all strategy-architecture combinations. Bold: correct protoform or reflex; ‘-’: the dataset does not contain this reflex; blank: reflex prediction is not applicable for the strategy.

	Romanian	French	Italian	Spanish	Portuguese	Latin
Reference	-	kɔ̃tʁibye	kontribuire	kontribwir	kuɲtʁibuiɾ	kɔ̃ntribɔ̃ere
GRU-SUPV						kɔ̃ntribɔ̃brem
GRU-IIM						kɔ̃ntributɔ̃rem
GRU-BST						kɔ̃ntribure
GRU-IIM-BST						kɔ̃ntribuirem
GRU-DPD	-	kɔ̃tʁibye	kontribuire	kontribwir	kuɲtʁibuiɾ	kɔ̃ntribɔ̃re
GRU-DPD-IIM	-	kɔ̃tʁibyl	kontribuire	kontribuer	kuɲtʁibuiɾ	kɔ̃ntribɔ̃re
GRU-DPD-BST	-	kɔ̃tʁibye	kontribuire	kontribwir	kuɲtʁibuiɾ	kɔ̃ntribuirem
GRU-DPD-IIM-BST	-	kɔ̃tʁibye	kontribuire	kontribwir	kuɲtʁibuiɾ	kɔ̃ntribuirem
Trans-SUPV						kɔ̃ntribire
Trans-IIM						kɔ̃ntribɔ̃re
Trans-BST						kɔ̃ntribire
Trans-IIM-BST						kɔ̃ntribɔ̃re
Trans-DPD	-	kɔ̃tʁibye	kontribuire	kontribwir	kuɲtʁibuiɾ	kɔ̃ntribɔ̃re
Trans-DPD-IIM	-	kɔ̃tʁibyɾ	kontriburre	kontribuir	kuɲtʁibuiɾ	kɔ̃ntribɔ̃re
Trans-DPD-BST	-	kɔ̃tʁibye	kontribuire	kontribwir	kuɲtʁibuiɾ	kɔ̃ntribɔ̃re
Trans-DPD-IIM-BST	-	kɔ̃tʁibye	kontribuire	kontribuir	kuɲtʁibuiɾ	kɔ̃ntribɔ̃re
Reference	tʃenzor	s̄s̄æɾ	tʃensore	θensor	seɲsoɾ	kensɔ̃rem
GRU-SUPV						kɪnsɔ̃res
GRU-IIM						kɪnsɔ̃rɔ̃m
GRU-BST						kensɔ̃rem
GRU-IIM-BST						kensɔ̃rem
GRU-DPD	tʃensor	s̄s̄æɾ	tʃensore	θensor	seɲsoɾ	kensɔ̃re
GRU-DPD-IIM	tʃensor	s̄s̄æɾ	tʃensore	θensor	seɲsuɾ	kɪnsɔ̃rem
GRU-DPD-BST	tʃensor	s̄s̄æɾ	tʃensore	θensor	seɲsoɾ	kensɔ̃rem
GRU-DPD-IIM-BST	tʃensor	s̄s̄æɾ	tʃensore	θensor	seɲsuɾ	kensɔ̃rem
Trans-SUPV						kenssɔ̃rem
Trans-IIM						kɪɲkɔ̃rem
Trans-BST						kensɔ̃rem
Trans-IIM-BST						kensɔ̃rem
Trans-DPD	tʃensor	s̄s̄æɾ	tʃensɔ̃re	θensor	seɲsoɾ	kɪɲserɔ̃rem
Trans-DPD-IIM	tʃensor	s̄s̄æɾ	tʃensore	θensor	seɲsoɾ	kensɔ̃rem
Trans-DPD-BST	tʃensor	s̄s̄æɾ	tʃensore	θensor	seɲsoɾ	kensɔ̃rem
Trans-DPD-IIM-BST	tʃensor	s̄s̄æɾ	tʃensore	θensor	seɲsoɾ	kensɔ̃rem
Reference	insuflatsje	ɛ̃syflasjɔ̃	insufflatsione	insuflaθjon	ɪɲsufləsɛ̃θ	insɔ̃fflatsionem
GRU-SUPV						insɔ̃flatsionem
GRU-IIM						ɪnsɔ̃pʰlatsionem
GRU-BST						ɪnsɔ̃flatsionem
GRU-IIM-BST						ɪnsɔ̃fflatsionem
GRU-DPD	insuflatsje	ɛ̃syflasjɔ̃	insufflatsione	insuflaθjon	ɪɲsufletɛ̃θ	insɔ̃flatsionem
GRU-DPD-IIM	insuflatsje	ɛ̃syflsjɔ̃	insufflatsione	insuflaθjon	ɪɲsufləsɛ̃θ	ɪnsɔ̃pʰlatsionem
GRU-DPD-BST	insuflatsje	ɛ̃syflasjɔ̃	insufflatsione	insuflaθjon	ɪɲsufləsɛ̃θ	ɪnsɔ̃flatsionem
GRU-DPD-IIM-BST	insuflaksje	ɛ̃syflakjɔ̃	insufflatsione	insuflaθjon	ɪɲsufləsɛ̃θ	ɪnsɔ̃flatsionem
Trans-SUPV						insɔ̃flatsionem
Trans-IIM						ɪnsɔ̃pʰlatsionem
Trans-BST						insɔ̃flatsionem
Trans-IIM-BST						insɔ̃flatsionem
Trans-DPD	insuflatsje	ɛ̃syflasjɔ̃	insufflatsione	insuflaθjon	ɪɲsufləsɛ̃θ	insɔ̃flatsionem
Trans-DPD-IIM	insullatsje	ɛ̃syllasjɔ̃	insullatsione	insullaθjon	ɪɲsulleɛ̃θ	insɔ̃latsionem
Trans-DPD-BST	insuflatsje	ɛ̃syflasjɔ̃	insufflatsione	insuflaθjon	ɪɲsufləsɛ̃θ	insɔ̃flatsionem
Trans-DPD-IIM-BST	insuflatsje	ɛ̃syflasjɔ̃	insufflatsione	insuflaθjon	ɪɲsufləsɛ̃θ	insɔ̃flatsionem
Reference	-	-	perdzurio	-	perəzurjɔ̃	perɪurɪɔ̃m
GRU-SUPV						perɔ̃rɪɔ̃m
GRU-IIM						perɔ̃rɪɔ̃m
GRU-BST						perɔ̃urɪɔ̃m
GRU-IIM-BST						perɔ̃urɪɔ̃m
GRU-DPD	-	-	perizurio	-	perizurɪɔ̃	perɔ̃urɪɔ̃m
GRU-DPD-IIM	-	-	perizurio	-	perizurɪɔ̃	perɪurɪɔ̃m
GRU-DPD-BST	-	-	perdzurio	-	perəzurjɔ̃	perɔ̃urɪɔ̃m
GRU-DPD-IIM-BST	-	-	pergzurio	-	perəgurjɔ̃	perɔ̃urɪɔ̃m
Trans-SUPV						perɔ̃urɪɔ̃m
Trans-IIM						perɔ̃urɪɔ̃m
Trans-BST						perɔ̃urɪɔ̃m
Trans-IIM-BST						perɔ̃urɪɔ̃m
Trans-DPD	-	-	perdzurio	-	perəzurjɔ̃	perɔ̃urɪɔ̃m
Trans-DPD-IIM	-	-	perdzurio	-	perəgurjɔ̃	perɔ̃urɪɔ̃m
Trans-DPD-BST	-	-	perdzurio	-	perəzurjɔ̃	perɔ̃urɪɔ̃m
Trans-DPD-IIM-BST	-	-	perdzurio	-	perəzurjɔ̃	perɔ̃urɪɔ̃m

Table 30: Example outputs on Rom-phon for all strategy-architecture combinations. Bold: correct protoform or reflex; ‘-’: the dataset does not contain this reflex; blank: reflex prediction is not applicable for the strategy.