

# Toward a Phonetic Approach for Multi-Dialect Speech Recognition in Vietnamese

Anonymous ACL submission

## Abstract

Vietnamese automatic speech recognition (ASR) remains challenging due to systematic dialectal variation across Northern, Central, and Southern regions, where identical lexical items often exhibit substantially different pronunciations. Most existing approaches address this variability primarily at the word level, relying on vocabularies that implicitly assume dialect-invariant mappings between orthography and pronunciation, which is linguistically inappropriate for Vietnamese. In this work, we propose a dialect-aware phonetic framework that explicitly models Vietnamese phonological structure and dialectal variation at both the vocabulary and decoding levels. We introduce a phonetic vocabulary that decomposes each syllable into structured phonetic components and maps them to dialect-specific IPA representations. Building on this representation, we design a phonetic-structure decoder that jointly predicts these components, enabling consistent and interpretable modeling. Experiments on the ViMD dataset demonstrate that the proposed approach consistently outperforms or matches strong pretrained baselines across dialects, achieving a WER of 13.35%, a PER of 8.45%, and dialect identification accuracy exceeding 95%, while using fewer parameters and no external pretraining. We will release code and phonetic resources for experimental reproducibility upon the acceptance of this paper.

## 1 Introduction

Automatic Speech Recognition (ASR) for Vietnamese multi-dialect speech remains challenging due to the language’s extensive dialectal variation. Speakers across the Northern, Central, and Southern regions exhibit systematic differences in pronunciation for the same lexical items, resulting in substantial acoustic-lexical mismatches (Phạm and McLeod, 2016; Nga et al., 2021; Dinh et al., 2024) (Figure 1). These mismatches often lead to degraded performance, as current ASR systems as-

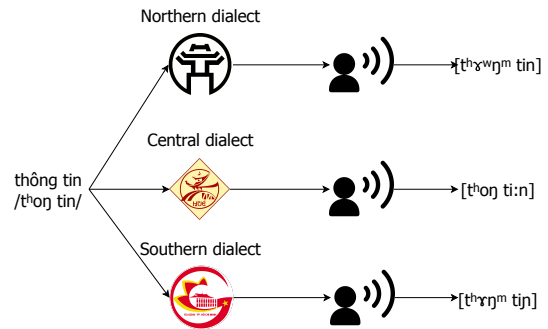


Figure 1: Examples of Vietnamese multi-dialect speech variations.

sume stable mappings between acoustic features and vocabulary units.

Prior work has primarily addressed dialectal variation from the acoustic perspective, through multi-task training (Elfeky et al., 2016; Dan et al., 2022), dialect-specific transfer learning (Ta et al., 2024; Suwanbandit et al., 2023), or multi-stage training combined with augmentation techniques that simulate pronunciation diversity (Park et al., 2019; Wang et al., 2023). While effective to some degree, these approaches implicitly treat dialect differences as variability to be normalized. Crucially, these methods do not address a core limitation at the vocabulary level: word-based and subword-based tokenizers presuppose a dialect-invariant correspondence between orthography and pronunciation. This assumption is particularly problematic for Vietnamese, where dialect-specific phenomena such as vowel shifts, coda neutralization, and tone contour variation are systematic and phonetically rule-based.

Vietnamese further motivates rethinking vocabulary design. As a monosyllabic, tonal language, each word corresponds to a structured combination of an initial consonant, a rhyme, and a tone. These components vary predictably across dialect families

070 and encode distinctions that are essential for both  
071 human perception and machine recognition (Pham  
072 and McLeod, 2016). Conventional Vietnamese  
073 ASR systems that operate on orthographic words  
074 or subwords overlook this structure, limiting their  
075 ability to generalize across dialects and exacerbating  
076 data sparsity and alignment inconsistencies,  
077 particularly for low-resource dialects.

078 In this work, we propose a dialect-aware pho-  
079 netic vocabulary representation for Vietnamese  
080 ASR, grounded in linguistic structure rather than  
081 surface orthography, explicitly encoding Viet-  
082 namese phonology and dialectal variation. Our  
083 method decomposes each syllable into linguisti-  
084 cally meaningful components, then maps them  
085 to dialect-specific International Phonetic Alpha-  
086 bet (IPA) patterns, yielding a compact and inter-  
087 pretable token set that represents systematic cross-  
088 dialectal differences. Beyond vocabulary design,  
089 we introduce a multi-token decoder that jointly pre-  
090 dict the initial, rhyme, and tone of each syllable.  
091 Unlike conventional single-token auto-regressive  
092 decoders, our architecture leverages the internal  
093 structure of Vietnamese phonology to model cross-  
094 component dependencies, thereby improving pho-  
095 netic consistency, reducing ambiguity, and stabiliz-  
096 ing alignment during training. This paper makes  
097 three main contributions:

- 098 1. We develop the first dialect-aware phonetic  
099 vocabulary for Vietnamese ASR that inte-  
100 grates linguistic structure with dialect-specific  
101 pronunciation patterns.
- 102 2. We propose a structured multi-token decod-  
103 ing mechanism that describes syllable compo-  
104 nents jointly to enhance phonetic discrimina-  
105 tion and decoding stability.
- 106 3. We demonstrate through extensive experi-  
107 ments that integrating dialectal phonetics at  
108 both the vocabulary and decoding levels sub-  
109 stantially improves word recognition accuracy,  
110 phonetic consistency, and vocabulary cov-  
111 erage compared with standard methods and  
112 strong pretrained models.

## 113 2 Related Works

114 Research on dialectal ASR has primarily focused  
115 on mitigating performance degradation caused by  
116 pronunciation variations. A prominent direction  
117 centers on multi-task and multi-model architectures  
118 that explicitly incorporate dialect labels during

119 training. Multi-task learning frameworks jointly  
120 optimize dialect classification and ASR to better  
121 align dialect-sensitive acoustic cues with recogni-  
122 tion objectives (Dan et al., 2022). Complement-  
123 arily, systems equipped with dialect identification  
124 modules route utterances to dialect-specific ASR  
125 models, reducing cross-dialect interference through  
126 specialized decoding paths (Elfeky et al., 2016).  
127 While these methods effectively utilize dialect la-  
128 bels, they typically require large-scale annotated  
129 multi-dialect corpora and remain limited to model-  
130 ing differences primarily at the acoustic level.

131 Another line of research addresses data scarcity  
132 and adaptation challenges using transfer learning  
133 or multi-stage training strategies. Models pre-  
134 trained on high-resource standard-accent corpora  
135 are fine-tuned on small dialectal datasets to im-  
136 prove robustness (Ta et al., 2024; Suwanbandit  
137 et al., 2023). More recent architectures, such as  
138 Aformer, adopt progressive, multi-stage optimiza-  
139 tion schedules to capture diverse acoustic patterns  
140 (Wang et al., 2023). Mixture-of-Experts (MoE)  
141 frameworks further enhance cross-dialect general-  
142 ization by combining dialect-sensitive experts with  
143 general experts under dynamic routing strategies  
144 (Zhou et al., 2024). Despite their effectiveness,  
145 these approaches largely conceptualize dialectal  
146 variation as acoustic variability to be normalized or  
147 compensated, thereby underexploiting the deeper  
148 lexical and phonological structure.

149 For Vietnamese, research has increasingly fo-  
150 cused on multi-dialect corpora (Tran et al., 2024;  
151 Dinh et al., 2024), allowing systematic analysis  
152 of pronunciation variation and syllable structure.  
153 Resources such as ViMD (Dinh et al., 2024) pro-  
154 vide valuable data for examining these linguistic  
155 patterns and evaluating dialect-aware ASR models  
156 at the provincial level. Nevertheless, most Viet-  
157 namese ASR systems continue to operate on ortho-  
158 graphic word or subword units (Nga et al., 2021).  
159 Such units assume a dialect-invariant mapping be-  
160 tween text and pronunciation, a problematic as-  
161 sumption for Vietnamese, where systematic cross-  
162 dialect phonological shifts are pervasive. As a re-  
163 sult, the lexical and phonological levels remain  
164 comparatively underexplored in Vietnamese ASR.

165 Overall, existing work predominantly addresses  
166 dialect variation through acoustic modeling and  
167 adaptation techniques. In contrast, vocabulary-level  
168 modeling that explicitly captures dialect-dependent  
169 phonological structure remains limited. This gap  
170 directly motivates our approach: we introduce a

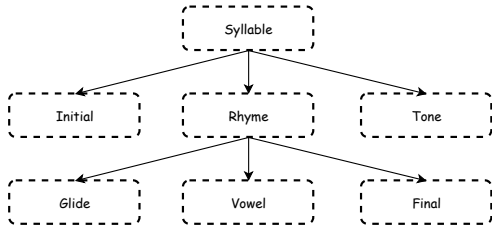


Figure 2: Vietnamese syllable structure.

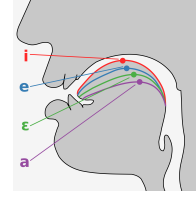


Figure 3: Highest tongue positions of front vowels.

171 phonologically grounded, dialect-aware vocabulary  
 172 representation and a structured decoding mecha-  
 173 nism that jointly model syllable components, en-  
 174 abling Vietnamese ASR systems to more faithfully  
 175 represent, generalize across, and ultimately benefit  
 176 from systematic dialectal variation.

### 3 Preliminaries

#### 3.1 Vietnamese Syllable Structure

179 Vietnamese syllables exhibit a highly regular struc-  
 180 ture consisting of three core components: ini-  
 181 tial, rhyme, and tone (Đoàn Thiện Thuật, 2016;  
 182 Hạo, 1998). Phonetic studies further decompose  
 183 the rhyme into glide, vowel, and final, where the  
 184 vowel is mandatory, and the remaining compo-  
 185 nents are optional (Đoàn Thiện Thuật, 2016; Hạo,  
 186 1998). Consequently, the rhyme serves as the cen-  
 187 tral carrier of phonological information in Viet-  
 188 namese syllables. Linguistic (Giáp, 2008, 2011;  
 189 Hạo, 1998) and natural language processing stud-  
 190 ies (Hieu Nguyen et al., 2025) consistently char-  
 191 acterize Vietnamese as a monosyllabic language,  
 192 in which each word contains at most one syllable.  
 193 Coupled with its phonemic orthography and  
 194 strong grapheme–phoneme correspondence, this  
 195 property enables a straightforward alignment be-  
 196 tween orthographic units and phonetic components.  
 197 These characteristics form the basis of our pro-  
 198 posed dialect-aware phonetic vocabulary, described  
 199 in Section 4.1. In this work, phonemes are deno-  
 200 ted using slashes, while dialect-specific phones are  
 201 denoted using square brackets. Figure 1 illustrates a  
 202 valid example of dialectal variation.

#### 3.2 Phonetic Variation in Multi-Dialect Vietnamese at the Phone Level

203 Multi-dialect variation in Vietnamese primarily  
 204 manifests at the phonetic level rather than the word  
 205 level. For example, a salient characteristic of the  
 206 Northern dialect is the interaction between front  
 207 vowels /i, e, ε/ and velar consonants /k, ŋ/. In  
 208  
 209

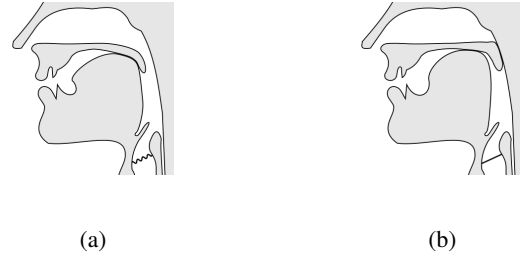


Figure 4: Mouth configurations during production of the voiced velar nasal [ŋ] (a) and the voiceless velar plosive [k] (b)

210 rhymes such as /iŋ, eŋ, εŋ, ik, ek, εk/, front vow-  
 211 els require a high front tongue position (Figure 3),  
 212 while velar consonants require contact with the  
 213 velum (Figure 4), resulting in articulatory incom-  
 214 patibility. This incompatibility leads to vowel cen-  
 215 tralization (/i, e, ε/ → [i, ə, ɜ]) and consonant  
 216 palatalization (/k, ŋ/ → [c, ɲ]), often accom-  
 217 panied by a transitional glide [j]. These patterns re-  
 218 flect systematic, phone-level variation rather than  
 219 changes in lexical identity, providing strong mo-  
 220 tivation for phone-based modeling in Vietnamese  
 221 multi-dialect ASR. A detailed dialect-specific anal-  
 222 ysis is in Appendix C.

### 4 Methodology

223 Our goal is to develop a vocabulary representation  
 224 and decoding framework that explicitly encodes  
 225 Vietnamese phonological structure and dialectal  
 226 regularities. This section presents: (i) the construc-  
 227 tion of a dialect-aware phonetic vocabulary, (ii) a  
 228 structured multi-dialect transcript representation,  
 229 and (iii) a phonetic-structure decoder that jointly  
 230 predicts internal syllable components.  
 231

#### 4.1 Dialect-aware Phonetic Vocabulary

232 Vietnamese syllables follow a regular phonolog-  
 233 ical template consisting of an *initial consonant*,  
 234 *glide*, *vowel*, *coda*, and *tone*. Because dialectal  
 235 variation systematically affects each of these com-  
 236 ponents, our vocabulary construction begins with a  
 237

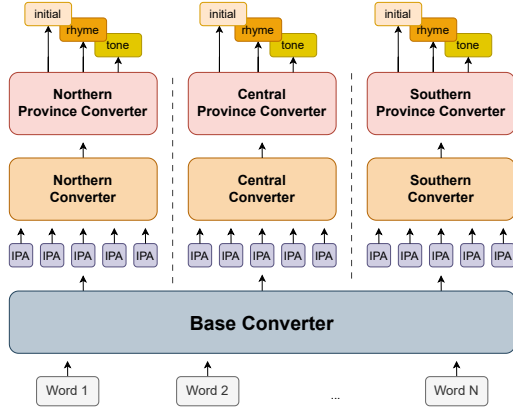


Figure 5: Pipeline for constructing a dialect-aware phonetic vocabulary for Vietnamese.

deterministic syllable parser that decomposes every orthographic word into these five constituent units. The parser enforces phonotactic constraints, including initial–vowel compatibility, coda admissibility, and some orthographic exceptions, ensuring that only phonotactically valid Vietnamese syllables enter the vocabulary. Each component is then mapped to its canonical IPA representation via a *Base Converter* encoding the standard Vietnamese phonology. These IPA arrays serve as a dialect-invariant reference layer that abstracts away from orthographic inconsistencies (Figure 5).

Dialect-specific regularities are flexibly integrated through a modular extension of the converter. A *Dialectal Converter* applies rule sets capturing systematic Northern, Central, and Southern shifts, such as initial mergers (e.g., /v/ → [j]), vowel raising/lowering, diphthong simplification, and coda neutralization. A further *Province Converter* refines these transformations to model a few province-level sociophonetic differences. This hierarchical design allows fine-grained dialect modeling without entanglement with surface orthography. The final representation for each word yields a compact, interpretable, and dialect-aware phonetic vocabulary suitable for Vietnamese ASR.

## 4.2 Multi-dialect Transcript Representation

Rhymes, comprising the medial, nucleus, and coda, constitute the core and consistently stable obligatory phonological unit of Vietnamese and encode the majority of dialect-induced variation in vowel quality and coda realization (Phạm and McLeod, 2016). To balance linguistic granularity with decoding efficiency, we adopt a structured three-part

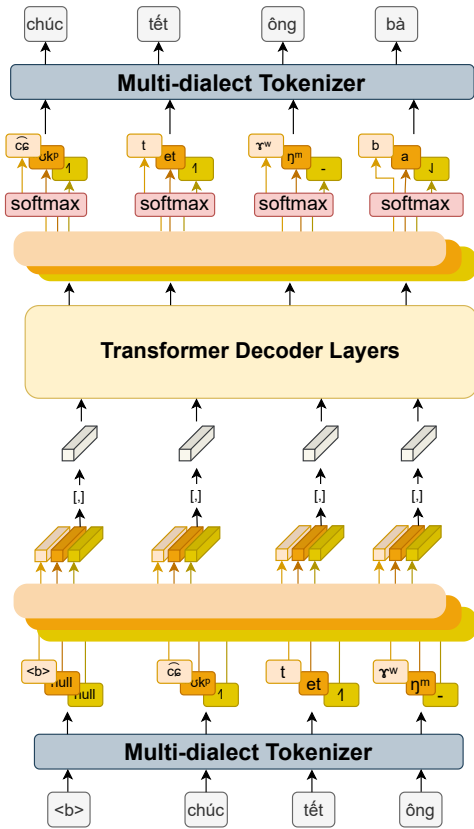


Figure 6: Vietnamese phonetic-structure decoder. The English text shown corresponds to Tet (Vietnamese New Year) wishes addressed to grandparents.

representation (initial, rhyme, tone) rather than using all five phonological components. The resulting vocabulary consists of three disjoint subsets:  $V_{\text{initial}}$  (27 categories),  $V_{\text{rhyme}}$  (~200 categories), and  $V_{\text{tone}}$  (6 categories). During training, output is a sequence of structured triplets. At inference time, predicted phone sequences are deterministically mapped back to orthographic forms using a reverse lexicon, which relies on the high consistency between phonetics and orthography in Vietnamese.

## 4.3 Vietnamese Phonetic-structure Decoder

To leverage the internal structure of Vietnamese syllables, we introduce a *multi-token autoregressive decoder* that predicts the initial, rhyme, and tone in parallel at each timestep, while maintaining autoregression across syllables (Figure 6). **Component-Specific Embeddings:** At timestep  $t$ , the decoder receives three syllabic components ( $\text{initial}_{t-1}$ ,  $\text{rhyme}_{t-1}$ ,  $\text{tone}_{t-1}$ ). Each component is embedded in its own learned subspace, preserving

the distinct phonological roles of each one. **Unification Layer:** The component embeddings are fused using a linear learned unification function, producing a unified representation that captures inter-component phonotactic dependencies.

$$h_t = f_{\text{unify}}(e_t^{\text{initial}}, e_t^{\text{rhyme}}, e_t^{\text{tone}}).$$

**Attention over Acoustic Context:** This step is similar to the standard decoding mechanism (Vaswani et al., 2017), where the unified representation attends to encoder outputs via multi-head attention, thereby integrating acoustic context with phonetic structure and enabling the decoder to capture dialect-sensitive cues and long-range temporal dependencies. **Parallel Multi-Token Prediction:** A shared projection layer produces a single probability vector that is partitioned into three distributions.

$$p(\text{initial}_t), p(\text{rhyme}_t), p(\text{tone}_t).$$

Joint prediction enforces internal consistency, avoids invalid syllables, and aligns naturally with the phonological structure of Vietnamese.

## 5 Experiments

### 5.1 Baselines, Datasets, and Metrics

**Baselines.** We compare the proposed dialect-aware phonetic vocabulary and structured decoding mechanism against a diverse set of strong baselines that reflect common design choices in Vietnamese ASR. *Zero-shot* baselines include publicly available large-scale Vietnamese ASR models, namely PhoWhisper (Le et al., 2024) and wav2vec 2.0 variants (Baevski et al., 2020) (wav2vec2-base-vietnamese, wav2vec2-base-vietnamese-160h, and wav2vec2-base-vietnamese-250h), which serve as off-the-shelf references. For fully supervised evaluations, we implement the Speech Transformer and Conformer (Dong et al., 2018; Gulati et al., 2020) architectures trained with standard *word-based* tokenization, where the output units correspond to orthographic Vietnamese words. To isolate the effect of vocabulary design, we construct *phone-based* counterparts that replace word tokens with phonetic units (initial, rhyme, tone) derived from the proposed dialect-aware vocabulary, while keeping the encoder–decoder architecture, parameter budget, and training configuration fixed. We additionally include a multi-task (MT) baseline that jointly predicts dialect labels and transcriptions, representing prior approaches

that inject dialect information at the acoustic modeling level (Dan et al., 2022). Finally, our full phonetic method, *Vietnamese phonetic-structure decoder* models, combine the phonetic vocabulary with the proposed Vietnamese phonetic-structure multi-token decoder, allowing us to disentangle the contributions of phonetic representation and decoding structure.

**Datasets.** All experiments are conducted on the publicly available Multi-Dialect Vietnamese (ViMD) corpus (Dinh et al., 2024), which comprises speech from all Vietnamese provinces. We adopt the official train–validation–test splits, enforcing strict speaker disjointness while preserving the original dialectal distribution. Only minimal preprocessing is applied: utterances containing non-Vietnamese tokens are filtered out, and audio is kept in its original mono 16kHz waveform to retain natural pronunciation and recording conditions. As a result, ViMD constitutes a challenging and realistic benchmark for evaluating the robustness of ASR systems across Vietnamese dialects.

**Metrics.** We evaluate models from complementary lexical and phonetic perspectives. At the orthographic level, Word Error Rate (WER) is reported as the primary metric, together with dialect-wise results. Northern provinces use the Northern dialect, provinces from Nghe An–Ha Tinh to Thua Thien Hue use the Central dialect, and the remaining provinces use the Southern dialect. For phone-based systems, we assess phonetic accuracy using the Phone Error Rate (PER) over the predicted phonetic sequences. We also report component-wise error rates for syllable initials, rhymes, and tones to analyze the phonological modeling behavior. Together, these metrics capture both transcription accuracy, linguistic fidelity, and dialectal robustness of the proposed framework.

### 5.2 Configuration

For all supervised experiments, models were trained from scratch on the ViMD corpus without using any external language models on a single NVIDIA H100 GPU. All experiments were implemented using a unified pipeline based on the SpeechBrain end-to-end speech processing toolkit (Ravanelli et al., 2021). Speech signals were represented by 80-dimensional log-Mel filterbank (Fbank) features, with a window size of 25ms and a step size of 10ms. To improve robustness, SpecAugment (Park et al., 2019) was applied

Model	Setup	Params	WER $\downarrow$	PER $\downarrow$	WER <sub>N</sub>	WER <sub>C</sub>	WER <sub>S</sub>	Time(s) $\downarrow$
<i>Published Zero-shot Baselines<math>\dagger</math></i>								
PhoWhisper-base	–	74M	21.78	–	18.17	27.18	23.80	0.73
PhoWhisper-small	–	244M	17.11	–	15.25	18.50	18.38	1.21
wav2vec2-base-vi	–	95M	22.48	–	19.87	27.15	23.82	0.03
wav2vec2-base-vi-160h	–	95M	31.12	–	27.09	40.49	32.83	0.03
wav2vec2-base-vi-250h	–	95M	16.83	–	14.51	19.73	18.23	0.07
<i>Published Fine-tuned Baselines<math>\dagger</math></i>								
PhoWhisper-base	FT	74M	16.30	–	–	–	–	–
wav2vec2-base-vi	FT	95M	15.80	–	–	–	–	–
wav2vec2-base-vi-160h	FT	95M	17.49	–	–	–	–	–
wav2vec2-base-vi-250h	FT	95M	13.56	–	–	–	–	–
<i>Our Implementations</i>								
Transformer	–	29M	16.66	–	14.68	18.64	17.95	0.37
Conformer	–	31M	17.56	–	16.05	19.92	18.39	0.29
Transformer	P	26M	15.13	12.52	13.55	<b>15.99</b>	16.27	1.25
Conformer	P	28M	13.96	11.29	12.01	17.20	14.99	0.93
MT-Transformer	P+M	26M	15.93	12.99	14.22	17.49	17.05	1.58
MT-Conformer	P+M	28M	16.75	13.58	14.71	19.56	17.93	1.21
<b>Our-Transformer</b>	<b>P+V</b>	26M	<b>13.35</b>	<b>8.45</b>	<b>11.37</b>	16.06	<b>14.51</b>	0.33
Our-Conformer	P+V	28M	13.80	8.91	11.78	17.07	14.89	0.28

Table 1: Performance comparison on Vietnamese multi-dialect ASR. Metrics include word error rate (WER, %) and phone error rate (PER, %). WER<sub>N</sub>, WER<sub>C</sub>, and WER<sub>S</sub> denote Northern, Central, and Southern dialects. Setup: phonetic supervision (P), multi-task learning (M), Vietnamese phonetic-structure decoder (V).  $\dagger$  Results reported from (Dinh et al., 2024).

during training, including both time masking and frequency masking with fixed mask widths.

For encoder–decoder architectures, Transformer employed 12 encoder and 6 decoder layers. Conformer consisted of 8 encoders and 4 decoders, each with an attention dimension of 256, 4 self-attention heads, and a feed-forward dimension of 2048. All models were optimized using the Adam optimizer (Kingma and Ba, 2015) with the Noam learning-rate schedule. Warm-up steps (Gotmare et al., 2019) were set to 40k for the Transformer and 20k for the Conformer, with initial learning rates of  $10^{-3}$  and  $4 \times 10^{-4}$ , respectively. Label smoothing (Szegedy et al., 2016) and dropout were both set to 0.1 for regularization. Training used a joint CTC–Attention objective (Kim et al., 2017), with the CTC loss weight set to 0.30 (Transformer) and 0.15 (Conformer). This joint objective promotes monotonic alignment through CTC while preserving the modeling flexibility of attention-based decoding, resulting in more stable training.

### 5.3 Main Results

Table 1 presents the main ASR results on the ViMD benchmark, comparing our proposed dialect-aware phonetic framework against strong zero-shot, fine-tuned, and supervised baselines. Overall, the re-

sults demonstrate that explicitly modeling Vietnamese phonological structure and dialectal variation at both the vocabulary and decoding levels yields consistent improvements across model architectures and dialect regions.

Large pretrained models, such as PhoWhisper and wav2vec 2.0, achieve competitive performance in both zero-shot and fine-tuned settings by leveraging large-scale pretraining. However, even after fine-tuning, these models remain constrained by word- or subword-based vocabularies that implicitly assume dialect-invariant pronunciations. In contrast, our models, trained from scratch on ViMD, achieve comparable or lower WERs with substantially fewer parameters, highlighting that linguistic inductive bias, specifically phonetic and dialect-aware representations, can partially compensate for the absence of massive pretraining when dialectal variation is systematic and structured.

**Effect of the phonetic vocabulary design.** Replacing word-level outputs with the proposed phonetic units yields clear improvements over standard Transformer and Conformer baselines. For example, Transformer-P reduces the WER from 16.66% to 15.13%, while Conformer-P achieves a larger reduction, from 17.56% to 13.96%. These results

Design	Transformer				Conformer			
	PER↓	Initial	Rhyme	Tone	PER↓	Initial	Rhyme	Tone
Phonetic baseline	12.52	12.29	14.53	10.52	11.29	11.07	13.33	9.26
+ Multitask	12.99	12.66	15.15	10.90	13.58	13.01	15.94	11.47
+ <b>Phonetic-structure decoder (our)</b>	<b>8.45</b>	<b>8.36</b>	<b>10.07</b>	<b>6.88</b>	8.91	8.82	10.63	7.24

Table 2: Phonetic error analysis of phone-based ASR models. Metrics include overall phone error rate (PER, %) and component-wise PERs for initial, rhyme, and tone units.

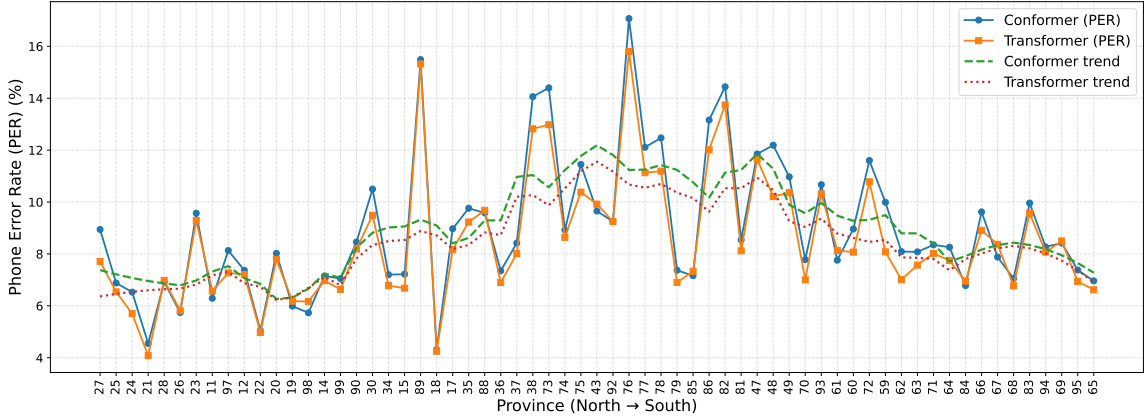


Figure 7: Phone error rate (PER, %) by province for Transformer and Conformer models using the proposed Vietnamese phonetics method.

suggest that explicit phone-level modeling reduces acoustic-lexical mismatch and improves generalization across dialects, with more pronounced gains for Conformer architectures.

**Contribution of the Vietnamese phonetic-structure decoder.** The largest gains are obtained with the full phonetic method, which further reduces WER to 13.35% (Transformer) and 13.80% (Conformer) without increasing model capacity or introducing additional supervision. By jointly predicting initial, rhyme, and tone, the decoder enforces internal phonological consistency and reduces invalid syllable hypotheses, leading to a more compact token representation and more stable decoding.

Dialect-specific analysis shows consistent improvements across Northern, Central, and Southern regions. Central dialect provinces, which exhibit greater internal phonetic variability, also benefit from the proposed approach. Figure 7 further corroborates this trend at the provincial level, where phone error rates are generally lower and more stable across regions.

## 5.4 Results Analysis

**Phonetic Structure and Dialect Modeling** Table 2 provides deeper insight into the main results.

While phonetic baselines already improve over word-based systems, they still exhibit relatively high initial, rhyme, and tone error rates, reflecting ambiguity when these components are predicted independently without structural constraints. The phonetic-structure decoder substantially reduces overall PER, from 12.52% to 8.45% for Transformer and from 11.29% to 8.91% for Conformer, with consistent improvements across all three phonetic components.

By comparison, multi-task learning with dialect classification improves dialect identification accuracy (Table 4), but this does not consistently translate into better ASR performance. In some cases, PER can even degrade relative to phonetic-only baselines, suggesting that injecting dialect information solely at the acoustic level is insufficient and may introduce competing optimization objectives, particularly when training on a modest-sized dataset such as ViMD. In contrast, our method integrates dialectal knowledge directly into the vocabulary and decoding process, resulting in more consistent and stable performance gains.

**Word Level vs. Phonetic Level Modeling** Table 3 compares lexical diversity and frequency bias between word-level and phonetic-level ASR models. Phonetic-level models correctly recognize a

Modeling Unit	Architecture	Unique Correct Words $\uparrow$	Pearson $r$	Spearman $\rho$
Word-level	Transformer	1,696	0.79	0.76
	Conformer	1,718	0.77	0.74
Phonetic-level	Transformer	<b>1,950</b>	0.57	0.45
	Conformer	<b>1,943</b>	0.58	0.47

Table 3: Lexical diversity and frequency bias analysis for word-level and phonetic-level ASR models. Metrics include the number of unique word types correctly recognized at least once in the test set, as well as Pearson/Spearman correlations between word frequency in the training data and per-word recall.

Model	Accuracy	Macro-F1
<i>Published baselines<math>\dagger</math></i>		
PhoWhisper-base	-	86.97
Whisper-base	-	85.59
wav2vec2-base-vi	-	91.02
wav2vec2-large-vi	-	91.47
wav2vec2-xlsr-300m	-	89.01
wav2vec2-large-xlsr-35	-	87.36
<i>Our implementations</i>		
<b>Our-Transformer</b>	<b>95.75</b>	<b>96.11</b>
Our-Conformer	95.47	95.80

Table 4: Classification performance (%) of multi-task phonetic models. Metrics include macro F1 and accuracy.  $\dagger$ Results reported from (Dinh et al., 2024).

substantially larger number of unique word types in the test set (from 1696 to 1950 for Transformer and from 1718 to 1943 for Conformer), indicating improved generalization to infrequent lexical items. Moreover, both Pearson and Spearman correlations between training-set word frequency and test-time recall are markedly lower for phonetic-level models, suggesting reduced reliance on frequency-driven memorization. These results suggest that phonetic-level modeling promotes structural generalization by capturing systematic phonological patterns across dialects, rather than exploiting surface-level lexical statistics.

**Dialect Classification Performance** Table 4 reports dialect classification results for the proposed multi-task phonetic models. Both Transformer and Conformer architectures achieve high overall accuracy (95.8% and 95.5%) and macro-F1 scores (96.1% and 95.8%), consistently matching fine-tuned PhoWhisper-base and wav2vec2-base-vi baselines reported in (Dinh et al., 2024). These results suggest that the learned phonetic representations encode dialect-discriminative information, despite dialect identification being treated as an auxiliary task. This finding provides empirical evidence that explicitly modeling Vietnamese phonological structure facilitates robust dialect characterization within a unified ASR and dialect classifica-

tion framework through multi-task learning.

## 6 Conclusion

In this work, we address a key limitation in Vietnamese multi-dialect ASR: the assumption of a dialect-invariant correspondence between orthography and pronunciation. Motivated by linguistic evidence that Vietnamese dialectal variation is systematic and primarily realized at the phone level, we propose a dialect-aware phonetic framework that integrates phonological structure into both vocabulary design and decoding. We introduce a phonetic vocabulary grounded in Vietnamese syllable structure, representing each syllable as an initial, rhyme, and tone mapped to canonical and dialect-specific IPA forms, together with a phonetic-structure decoder that jointly predicts these components while explicitly modeling their internal dependencies. Experimental results demonstrate that incorporating phonological structure and dialectal regularities provides an effective inductive bias for multi-dialect ASR and complements large-scale pretraining for languages with structured and systematic phonology.

## Limitations

Despite its effectiveness, the proposed approach has several limitations. The structured decoder jointly predicts initials, rhymes, and tones at the syllable level; it does not explicitly model cross-syllabic or cross-word phonological and prosodic phenomena, which remain open research directions. Moreover, while the framework is linguistically grounded, extending it to other languages would require language-specific phonological analysis and rule design, and the current implementation is limited to Vietnamese. Future work includes combining rule-based and data-driven phonetic modeling, integrating large pretrained acoustic models, and extending the framework to other low-resource, multi-dialect languages.

562  
563  
564  
565  
566  
567  
568  
569  
  
570  
571  
  
572  
573  
574  
575  
  
576  
577  
  
578  
579  
580  
581  
582  
583  
584  
585  
  
586  
587  
588  
589  
590  
591  
  
592  
593  
594  
595  
596  
597  
  
598  
599  
  
600  
601  
  
602  
603  
604  
605  
606  
607  
  
608  
609  
610  
611  
612  
613  
614  
615  
616

## References

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. [wav2vec 2.0: A framework for self-supervised learning of speech representations](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Hoàng Thị Châu. 2002. *Phương ngữ học tiếng Việt*. Nhà xuất bản Đại học Quốc gia Hà Nội.

Zhengjia Dan, Yue Zhao, Xiaojun Bi, Licheng Wu, and Qiang Ji. 2022. [Multi-task transformer with adaptive cross-entropy loss for multi-dialect speech recognition](#). *Entropy*, 24(10):1429.

Alexandre de Rhodes. 1651. *Dictionarium Anna-miticum Lusitanum et Latinum*.

Nguyen Dinh, Thanh Dang, Luan Thanh Nguyen, and Kiet Van Nguyen. 2024. [Multi-dialect vietnamese: Task, dataset, baseline models and challenges](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 7476–7498. Association for Computational Linguistics.

Linhao Dong, Shuang Xu, and Bo Xu. 2018. [Speech-transformer: A no-recurrence sequence-to-sequence model for speech recognition](#). In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2018, Calgary, AB, Canada, April 15-20, 2018*, pages 5884–5888. IEEE.

Mohamed Elfeky, Meysam Bastani, Xavier Velez, Pedro J. Moreno, and Austin Waters. 2016. [Towards acoustic model unification across dialects](#). In *2016 IEEE Spoken Language Technology Workshop, SLT 2016, San Diego, CA, USA, December 13-16, 2016*, pages 624–628. IEEE.

Nguyễn Thiện Giáp. 2008. *Từ vựng học tiếng Việt*. Nhà xuất bản Giáo dục Việt Nam.

Nguyễn Thiện Giáp. 2011. *Vấn đề "từ" trong tiếng Việt*. Nhà xuất bản Giáo dục Việt Nam.

Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2019. [A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. 2020. [Conformer: Convolution-augmented transformer for speech recognition](#). In *21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020*, pages 5036–5040. ISCA.

Nghia Hieu Nguyen, Dat Tien Nguyen, and Ngan Luu-Thuy Nguyen. 2025. [Vietnamese words are not constructed from syllables: Rethinking the role of word segmentation in natural language processing for vietnamese texts](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(22):24069–24077. 617  
618  
619  
620  
621  
622

Cao Xuân Hạo. 1998. *Tiếng Việt mấy vấn đề ngữ âm - ngữ pháp - ngữ nghĩa*. Nhà xuất bản Giáo dục Việt Nam. 623  
624  
625

Suyoun Kim, Takaaki Hori, and Shinji Watanabe. 2017. [Joint ctc-attention based end-to-end speech recognition using multi-task learning](#). In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2017, New Orleans, LA, USA, March 5-9, 2017*, pages 4835–4839. IEEE. 626  
627  
628  
629  
630  
631

Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. 632  
633  
634  
635  
636

Thanh-Thien Le, Linh The Nguyen, and Dat Quoc Nguyen. 2024. [Phowhisper: Automatic speech recognition for vietnamese](#). In *The Second Tiny Papers Track at ICLR 2024, Tiny Papers @ ICLR 2024, Vienna, Austria, May 11, 2024*. OpenReview.net. 637  
638  
639  
640  
641

Cao Hong Nga, Chung-Ting Li, Yung-Hui Li, and Jia-Ching Wang. 2021. [A survey of vietnamese automatic speech recognition](#). In *2021 9th International Conference on Orange Technology (ICOT)*, pages 1–4. 642  
643  
644  
645  
646

Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. 2019. [SpecAugment: A simple data augmentation method for automatic speech recognition](#). In *20th Annual Conference of the International Speech Communication Association, Interspeech 2019, Graz, Austria, September 15-19, 2019*, pages 2613–2617. ISCA. 647  
648  
649  
650  
651  
652  
653  
654

Ben Phạm and Sharynne McLeod. 2016. [Consonants, vowels and tones across vietnamese dialects](#). *International Journal of Speech-Language Pathology*, 18(2):122–134. 655  
656  
657  
658

Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, Loren Lugosch, Cem Subakan, Nauman Dawalatabad, Abdelwahab Heba, Jianyuan Zhong, Ju-Chieh Chou, Sung-Lin Yeh, Szu-Wei Fu, Chien-Feng Liao, Elena Rastorgueva, François Grondin, William Aris, Hwidong Na, Yan Gao, and 2 others. 2021. [Speechbrain: A general-purpose speech toolkit](#). abs/2106.04624. 659  
660  
661  
662  
663  
664  
665  
666

Artit Suwanbandit, Burin Naowarat, Orathai Sangpetch, and Ekapol Chuangsuwanich. 2023. [Thai dialect corpus and transfer-based curriculum learning investigation for dialect automatic speech recognition](#). In *24th Annual Conference of the International Speech Communication Association, Interspeech 2023, Dublin, Ireland, August 20-24, 2023*, pages 4069–4073. ISCA. 667  
668  
669  
670  
671  
672  
673  
674

Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. [Re-thinking the inception architecture for computer vision](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2818–2826. IEEE Computer Society.

Bao Thang Ta, Nhat Minh Le, and Van Hai Do. 2024. [Transfer learning methods for low-resource speech accent recognition: A case study on vietnamese language](#). *Eng. Appl. Artif. Intell.*, 132:107895.

Linh Thi Thuc Tran, Han-Gyu Kim, Hoang Minh La, and Su Van Pham. 2024. Automatic speech recognition of vietnamese for a new large-scale corpus. *Electronics*, 13(5).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Xuefei Wang, Yanhua Long, Yijie Li, and Haoran Wei. 2023. [Multi-pass training and cross-information fusion for low-resource end-to-end accented speech recognition](#). In *24th Annual Conference of the International Speech Communication Association, Interspeech 2023, Dublin, Ireland, August 20-24, 2023*, pages 2923–2927. ISCA.

Jie Zhou, Shengxiang Gao, Zhengtao Yu, Ling Dong, and Wenjun Wang. 2024. [Dialectmoe: An end-to-end multi-dialect speech recognition model with mixture-of-experts](#). In *Chinese Computational Linguistics - 23rd China National Conference, CCL 2024, Taiyuan, China, July 25-28, 2024, Proceedings*, volume 14761 of *Lecture Notes in Computer Science*, pages 243–258. Springer.

Đoàn Thiện Thuật. 2016. *Ngữ âm tiếng Việt*. Nhà xuất bản Đại học Quốc gia Hà Nội.

## A Appendix: Experimental Details

### A.1 Dataset Details

Split	Duration (h)	#Utterances	#Unique Words
Train	69.58	13,061	3,859
Dev	8.65	1,631	2,381
Test	9.45	1,787	2,441

Table 5: Statistics of the ViMD dataset across training, development, and test splits.

Table 5 presents the detailed statistics of the ViMD dataset used in our experiments, including the duration, number of utterances, and number of unique words for each split (training, development, and test).

### A.2 Metrics

We evaluate model performance using complementary metrics that capture orthographic transcription accuracy, phonetic fidelity, dialectal robustness, and lexical generalization behavior.

**Word Error Rate (WER).** End-to-end transcription quality is primarily measured using *Word Error Rate* (WER), defined as

$$\text{WER} = \frac{S + D + I}{N}, \quad (1)$$

where  $S$ ,  $D$ , and  $I$  denote the numbers of word-level substitutions, deletions, and insertions obtained via minimum-edit-distance alignment, and  $N$  is the number of words in the reference transcription. In addition to overall WER, we report dialect-wise WER for Northern, Central, and Southern Vietnamese, following the regional grouping defined in the ViMD benchmark.

**Phone Error Rate (PER).** For phonetic-level systems, we report *Phone Error Rate* (PER), computed analogously at the phone level:

$$\text{PER} = \frac{S_p + D_p + I_p}{N_p}, \quad (2)$$

where  $S_p$ ,  $D_p$ , and  $I_p$  are phone-level substitutions, deletions, and insertions, and  $N_p$  is the number of reference phones. PER provides a linguistically grounded assessment of pronunciation modeling, which is particularly relevant for Vietnamese due to systematic dialectal variation at the phone level.

**Component-wise Phonetic Error Rates.** To analyze phonological modeling behavior with respect to Vietnamese syllable structure, we further compute component-wise error rates for syllable *initials*, *rhymes*, and *tones*. For a component  $c \in \{\text{initial, rhyme, tone}\}$ , the error rate is defined as

$$\text{ER}_c = \frac{S_c + D_c + I_c}{N_c}, \quad (3)$$

where  $S_c$ ,  $D_c$ ,  $I_c$ , and  $N_c$  are computed over the aligned component sequence.

**Dialect Classification Metrics.** For multi-task models that jointly perform ASR and dialect identification, we report *classification accuracy* and *macro-averaged F1-score*. Given a set of dialect classes  $C$ , macro-F1 is computed as

$$\text{Macro-F1} = \frac{1}{|C|} \sum_{c \in C} \text{F1}_c, \quad (4)$$

ensuring equal weighting across dialects regardless of class imbalance.

**Lexical Diversity and Frequency Bias.** To assess lexical generalization beyond aggregate error rates, we report the number of *unique word types* correctly recognized at least once in the test set. We further analyze frequency bias by measuring the relationship between training-set word frequency and per-word recall.

For a word type  $w$ , its training-set frequency is defined as

$$f_{\text{train}}(w) = \sum_{u \in \mathcal{D}_{\text{train}}} \sum_{i=1}^{|u|} \mathbf{1}(u_i = w), \quad (5)$$

where  $\mathcal{D}_{\text{train}}$  denotes the training corpus and  $\mathbf{1}(\cdot)$  is the indicator function. To mitigate the heavy-tailed frequency distribution, we apply a logarithmic transformation:

$$\tilde{f}_{\text{train}}(w) = \log(1 + f_{\text{train}}(w)). \quad (6)$$

Per-word recall is defined as

$$\text{recall}(w) = \frac{\sum_{(x,y) \in \mathcal{D}_{\text{test}}} \sum_{i=1}^{|\hat{y}|} \mathbf{1}(\hat{y}_i = w)}{\sum_{(x,y) \in \mathcal{D}_{\text{test}}} \sum_{i=1}^{|y|} \mathbf{1}(y_i = w)}, \quad (7)$$

where  $(x, y)$  is a test utterance–reference pair and  $\hat{y}$  is the predicted word sequence after optimal alignment.

The relationship between training frequency and recognition performance is quantified using Pearson and Spearman correlation coefficients:

$$r = \text{Pearson}(\tilde{f}_{\text{train}}(w), \text{recall}(w)), \quad (8)$$

$$\rho = \text{Spearman}(\tilde{f}_{\text{train}}(w), \text{recall}(w)), \quad (9)$$

computed over all word types appearing at least once in the test set. Lower correlation values indicate reduced dependence on frequency-driven memorization and stronger structural generalization.

### A.3 Additional Dialect Classification Results

Table 6 shows that the Central dialect achieves the highest classification accuracy among the three dialects for both Transformer and Conformer models, as reflected by the strong diagonal entries and minimal cross-dialect confusions. In contrast, most classification errors arise from confusion between the Northern and Southern dialects,

which are mutually misclassified more frequently than either is with the Central dialect. This pattern is consistent across architectures and suggests that Central Vietnamese exhibits more distinctive acoustic-phonetic characteristics, while Northern and Southern dialects share partially overlapping phonological properties that make them harder to separate. Overall, the confusion matrices provide fine-grained evidence that complements the aggregate performance metrics in Table 4.

## B Appendix: Phonetics and Orthography in Vietnamese

Traditionally, the Vietnamese had the Nom alphabet as the main orthography system. However, the Nom alphabet was developed based on the ancient Han alphabet. This means we have to be fluent in the ancient Han to have the ability to use the Nom alphabet for reading and writing. Later on, (de Rhodes, 1651) used the Latin alphabet to describe the speech sound of this language. This orthography system is simple and effective for describing almost all phonetic phenomena in Vietnamese. With its advantages that counterbalance its disadvantages, this Latin alphabet is gradually improved over time and finally becomes the national alphabet system of modern Vietnamese.

Although having the Nom alphabet or the Latin alphabet, these orthography systems all reflect two consistent characteristics of Vietnamese:

1. Vietnamese is a monosyllabic language.
2. The correspondence between graphemes and phonemes in Vietnamese is consistent.

That is, in this language, we do not have the linking pronunciation as in English, and every phoneme has persistent writing forms. In Vietnamese, each syllable has three components: initials, rhymes, and tones. Rhyme has smaller components, which are glide, vowel, and final. We provide the list of all phonemes according to the syllabic structure of Vietnamese:

- 22 phonemes as the initials:
  - Plosive consonants: /b, t, t<sup>h</sup>, k/.
  - Fricative consonants: /f, d, x, z, j, s, ʃ, ç̥, t̚, ŋ, x, v/.
  - Nasal consonants: /n, m, ŋ, p/.
  - Vibrant consonants: /r, l/.
- 01 phonemes as the glide: /ɰ/.

True	Predicted Dialect			True	Predicted Dialect		
	Northern	Central	Southern		Northern	Central	Southern
	Conformer			Transformer			
Northern	711	0	13	Northern	714	0	10
Central	2	148	5	Central	0	152	3
Southern	58	3	847	Southern	57	6	845

Table 6: Dialect confusion matrices for multi-task Transformer and Conformer models on the ViMD test set. Rows denote ground-truth dialects and columns denote predicted dialects.

852	• 15 phonemes as the vowels:	phonemes. These examples might include	886
853	– Diphthongs: /ie, uo, uə/.	the tone mark above or below the graphemes.	887
854	– Monophthongs: /a, ă, ơ, i, ɛ, e, u, ư, o,	Readers can discard these marks to see the true	888
855	ơ, ơ, ə/.	writing form of the vowels in Vietnamese. For	889
856	• 10 phonemes as the finals:	instance, the dot below <b>iê</b> /iə/ in word <b>kiêm</b>	890
857	– Nasal consonants: /n, t/.	/kiem-ɿ̯/ denotes the mid glottalized-raising	891
858	– Labial consonants: /m, p/.	tone /ɿ̯/.	892
859	– Velar consonants: /ŋ, k/.	• 22 initials have 26 writing forms:	893
860	– Palatal consonants: /j, c/.	– b /b/. Eg: <b>ba</b> mẹ, <b>bánh</b> kẹo, <b>buôn</b> bán.	894
861	– Semivowels: /ɥ, ɨ/.	– t /t/. Eg: <b>tâm</b> tư, <b>tĩnh</b> tiến, <b>tính</b> cách.	895
862	• 6 phonemes denote tones:	– th /tʰ/. Eg: <b>thách</b> thức, <b>thành</b> thạo.	896
863	– Flat tone is denoted by nothing.	– k, c, or q /k/. Eg: <b>cách</b> mạng, <b>quan</b> hệ,	897
864	– Low falling tone: /H̄/.	<b>hiện</b> kim.	898
865	– Mid raising tone: /H̆/.	– ph /f/. Eg: <b>phụ</b> huynh, <b>phong</b> cách, <b>phân</b>	899
866	– Mid falling tone: /Ḣ/.	<b>định</b> .	900
867	– Mid glottalized-falling tone: /H̆̚/.	– đ /d/. Eg: <b>đưa</b> đón, <b>đậm</b> đà.	901
868	– Mid glottalized-raising tone: /H̆̆̚/.	– gh or g /ɣ/. Eg: <b>ga</b> tàu, <b>gánh</b> hát, <b>ganh</b>	902
869	The writing form of these phonemes is consistent	<b>ghét</b> .	903
870	regardless of grammar. In particular:	– gi /z/. Eg: <b>giếng</b> nước, <b>giống</b> loài.	904
871	• 6 tones are denoted by a mark above or below	– đ /d/. Eg: <b>đứng</b> đắn, <b>đêm</b> tối, <b>đèn</b> đuốc.	905
872	the graphemes of vowels:	– d /j/. Eg: <b>da</b> dẻ, tiêu <b>dùng</b> , <b>dụng</b> cụ.	906
873	– Flat tone is denoted by nothing (a).	– x /s/. Eg: sản <b>xuất</b> , xuất <b>xứ</b> , xe <b>cộ</b> .	907
874	– Low falling tone /H̄/ is denoted by a	– s /ʃ/. Eg: xác <b>suất</b> , so <b>sánh</b> , sao <b>chép</b> .	908
875	grave accent (à).	– ch /ç/. Eg: <b>chứa</b> <b>chan</b> , <b>che</b> <b>chở</b> .	909
876	– Mid raising tone /H̆/ is denoted by an	– tr /ʈ/. Eg: <b>tranh</b> chấp, tiết <b>trùng</b> .	910
877	acute accent (á).	– ng or ngh /ŋ/. Eg: mong <b>ngóng</b> , tình	911
878	– Mid falling tone /Ḣ/ is denoted by a hook	<b>nghĩa</b> .	912
879	above (â).	– nh /ɲ/. Eg: <b>nhà</b> cửa, nổi <b>nhớ</b> , <b>nhung</b> lụa.	913
880	– Mid glottalized-falling tone /H̆̚/ is de-	– l /l/. Eg: <b>lắm</b> lem, <b>lung</b> linh, <b>lối</b> về.	914
881	noted by a tilde above (ã).	– r /r/. Eg: <b>rậm</b> rạp, <b>rón</b> rén, <b>rực</b> rỡ.	915
882	– Mid glottalized-raising tone /H̆̆̚/ is de-	– kh /x/. Eg: <b>khó</b> khăn, <b>khởi</b> sắc, <b>khâm</b>	916
883	noted by a dot below (ạ).	<b>khá</b> .	917
884	In the following texts, we provide writing	– v /v/. Eg: <b>vui</b> vẻ, <b>vương</b> vấn, <b>vấy</b> vùng.	918
885	forms of vowels regarding the mentioned	– m /m/. Eg: <b>mong</b> <b>mỏi</b> , <b>máy</b> <b>mắn</b> , <b>mênh</b>	919
		<b>mông</b> .	920
		– n /n/. Eg: đất <b>nước</b> , <b>núi</b> non, <b>nông</b> cạn.	921
		• The glide /ɥ/ has two writing forms: u or o.	922
		Eg: <b>quê</b> nhà, <b>hoa</b> cỏ, <b>khuyến</b> khích.	923

- 924 • 03 diphthongs have 8 writing forms:
- 925 – iê, yê, ia, or ya /ie/. Eg: **khiếm** thị, **yên**
- 926 **ắng**, **chia** sẻ, **khuya** khoắt.
- 927 – uô or ua /uo/. Eg: **khuôn** khổ, **mua** bán.
- 928 – ươ or ưa /uə/. Eg: **khướu** giác, **dây** **dưa**.

- 929 • 12 monophthongs have 13 writing forms:
- 930 – a /a/. Eg: **ba** mẹ, **tranh** vẽ, **ngã** **ba**, **mả**
- 931 **miết**, **làng** **chài**.
- 932 – ă or a /ă/. Eg: ánh **ắ**ng, **nắ**m tay, **nắ**m
- 933 **thắ**ng, **á**y **ná**y, **chạ**y **nhắ**y.
- 934 – â /â/. Eg: **nâng** niu, **ắ**n định, **ắ**n vang.
- 935 – i or y /i/. Eg: **thi** cử, **trữ** nặng, **bữ** môi.
- 936 – ê /e/. Eg: **kết** quả, **thể** hiện, **mân** **mê**.
- 937 – e /ɛ/. Eg: mùa **hè**, **xe** **cộ**, **té** **ngã**.
- 938 – u /u/. Eg: **thu** mua, **mủ**m **mỉ**m, **lung** lay,
- 939 **trung** thành.
- 940 – ư /ɯ/. Eg: **trư**ng cầu, xây **dự**ng, **ứ**ng ý.
- 941 – o /ɔ/. Eg: **chăm** sóc, **mong** **ngố**ng, **trong**
- 942 **veo**.
- 943 – oo /ɔː/. Eg: **xoong** chào.
- 944 – ô /o/. Eg: **trồ**ng **đồ**ng, **ồ**ng hút, **cổ** **nhân**.
- 945 – ơ /ə/. Eg: **mơ** **mộng**, **cơ** **nhữ**, **chơi** **bờ**i.

- 946 • 10 final consonants have 12 writing forms:
- 947 – i or y /i/. Eg: **làng** **chài**, **mỗ**i **mệ**t, **chạ**y
- 948 **đũa**, **bay** **nhắ**y.
- 949 – m /m/. Eg: **ê**m **ắ**m, **nhiệ**m màu, **mâm** **cổ**.
- 950 – n /n/. Eg: **nan** giải, **non** **nót**, **tản** **mạn**.
- 951 – ng /ŋ/. Eg: **sang** **trọng**, **trồ**ng **trắ**i, **sung**
- 952 **túc**.
- 953 – nh /ɲ/. Eg: **nhanh** **nhẹn**, **bê**nh **vực**, **binh**
- 954 **quyền**.
- 955 – p /p/. Eg: **phậ**p **phồ**ng, **thắ**p **thỏ**m, **thắ**p
- 956 **tùng**, **gượ**ng **é**p, **ú**c **hié**p, **tắm** **ướ**p.
- 957 – t /t/. Eg: **lấn** **át**, **bát** **đĩa**, **kết** quả, **hó**t **hả**.
- 958 – c /k/. Eg: **cúc** áo, **chực** chờ, **bốc** **vác**.
- 959 – ch /c/. Eg: **cách** thức, **chích** **ngừa**.
- 960 – u or o /ɯ/. Eg: **trau** **chu**ốt, **chị**u **đự**ng,
- 961 **xoong** **chắ**o, **xiêu** **vẹ**o.

962 It is important to note that the writing form of

963 phonemes in Vietnamese is consistent in all use-

964 age cases, regardless of grammar (tense, aspect,

965 mood). This feature defines Vietnamese as an iso-

966 lating language whose morphology is an aspect of

967 phonetics rather than grammar, as in inflectional

968 languages. To this end, given a Vietnamese word,

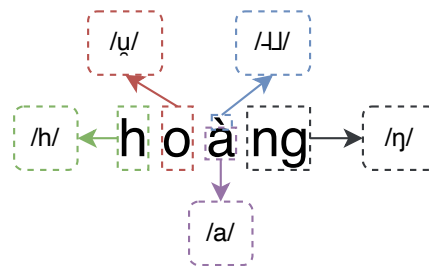


Figure 8: Example for the consistency between graphemes and phonemes in Vietnamese.

969 native speakers can find no difficulty in pronounc-

970 ing it. On the other hand, Vietnamese natives can

971 write down any Vietnamese word by listening to

972 its speech sound without knowing how to spell it.

973 For instance, given the word **hoàng**, its phonetic

974 representation is /huan̩-ɰ/. We can determine the

975 grapheme of the initial /h/ is **h**, of the glide /u/ is

976 **o**, of the vowel /a/ is **a**, of the /ŋ/ is **ng**, and tone is

977 denoted by the grave accent above **à** (Figure 8).

978 However, although the conversion from

979 grapheme to phonemes is straightforward (that is,

980 many graphemes correspond to unique phonemes),

981 the inversion is not all ways many-to-one mapping.

982 There are some phonemes having one-to-many

983 mapping with graphemes, such as the initial /ɣ/

984 can be written as **g** or **gh**, or the diphthong /ie/

985 can be written as **iê**, **yê**, **ia** or **ya**. Actually, a particular

986 writing form of such phonemes is determined

987 consistently via the neighbor phonemes. We

988 detail here rules for determining writing forms of

989 phonemes having multiple respective graphemes:

- 990 • The initial /k/ is written as:
- 991 – **k** if followed by **i** /i/, **ê** /e/, **e** /ɛ/, **iê** /ie/.
- 992 Eg: **kíp** thời, **kiểu** cách, **kèm** cặp.
- 993 – **q** if followed by **u** /u/ as the glide. Eg:
- 994 **quê** hương, **quà** cáp.
- 995 – **c** otherwise. Eg: **cứng** **cỏi**, **cỏi** **mỏ**, **hạt**
- 996 **cát**, **rau** **củ**.
- 997 • The initial /ɣ/ is written as:
- 998 – **gh** if followed by **i** /i/, **ê** /e/, **e** /ɛ/, **iê** /ie/.
- 999 Eg: **bàn** **ghế**, **ghi** chép, **ghe** tàu.
- 1000 – **g** otherwise. Eg: **thanh** **gươm**, **gồng**
- 1001 **gánh**, **gọi** **điện**.
- 1002 • The initial /ŋ/ is written as:

1003	– <b>ng</b> if followed by textbf /i/, ê /e/, e /ɛ/,	– <b>y</b> if the rhyme has /ă, ɔ̃/ as the vowel. Eg:	1046
1004	<b>iê</b> /ie/. Eg: <b>nghe</b> ngóng, <b>nghiêm</b> nghị,	<b>bay</b> lượn, <b>chạy</b> nhảy, <b>cây</b> cày.	1047
1005	<b>ng</b> hê sĩ.		
1006	– <b>ng</b> otherwise. Eg: <b>ngành</b> nghề, <b>ngõ</b>	Following these orthographic rules in Viet-	1048
1007	ngịch, <b>ngọt</b> ngào.	namese, there is no ambiguity in converting	1049
1008	• The diphthong /ie/ is written as:	phonemes to graphemes and vice versa. Analysis	1050
1009	– <b>iê</b> if the rhyme has a final consonant and	in this Section together with the analysis in the fol-	1051
1010	no glide. Eg: <b>kiến</b> thức, tiết <b>kiệm</b> .	lowing Appendix C serve as the fundamentals for	1052
1011	– <b>yê</b> if the rhyme has a final consonant and	our Dialect-aware Tokenization algorithm, which	1053
1012	the glide written as <b>u</b> . Eg: <b>khuyên</b> bảo,	is described comprehensively in Appendix D.	1054
1013	<b>uẩn</b> chuyển, câu <b>chuyện</b> .		
1014	– <b>ya</b> if the rhyme has no final consonant	<b>C Appendix: Multi-dialect Speech in</b>	1055
1015	and the glide written as <b>u</b> . Eg: <b>đêm</b>	<b>Vietnamese</b>	1056
1016	<b>khuya</b> .	Most studies in multi-dialect speech in Vietnamese	1057
1017	– <b>ia</b> if the rhyme has no final consonant	(Hạo, 1998; Đoàn Thiện Thuật, 2016; Châu, 2002)	1058
1018	and no glide. Eg: <b>bìa</b> sách, <b>chia</b> sẻ.	agree that Vietnamese has three dialects of speech,	1059
1019	• The diphthong /uo/ is written as:	which are the Northern dialect, the Central dialect,	1060
1020	– <b>uô</b> if the rhyme has a final consonant.	and the Southern dialect. These dialects were stud-	1061
1021	Eg: nổi <b>buồn</b> , <b>muối</b> biển, <b>muộn</b> màn,	ied and grouped mostly depending on how the resi-	1062
1022	<b>chuồn</b> chuồn.	dents deal with the rhymes. However, when diving	1063
1023	– <b>ua</b> if the rhyme has no final consonant.	into each particular dialect, there are more com-	1064
1024	Eg: <b>chùa</b> chiền, nhảy <b>múa</b> .	plexed phonetic phenomena, and more variations	1065
1025	• The diphthong /uə/ is written as:	are explored. Section 3.2 above briefly gives an	1066
1026	– <b>uơ</b> if the rhyme has a final consonant.	overview of the main differences among the three	1067
1027	Eg: bia <b>rượu</b> , <b>hưởng</b> thụ, chiêm <b>ngưỡng</b> .	main dialects in Vietnam. In this section, we pro-	1068
1028	– <b>ua</b> if the rhyme has no final consonant.	vide the reader with a more detailed discussion and	1069
1029	Eg: <b>chùa</b> chiền, nhảy <b>múa</b> .	analysis of each dialect. These discussions form	1070
1030	• The monophthong /ă/ is written as:	the fundamentals for our multi-dialect tokeniza-	1071
1031	– <b>a</b> if is is followed by character <b>y</b> . Eg:	tion method, which is described in the following	1072
1032	<b>máy</b> bay, <b>cay</b> nông, <b>tay</b> chân.	sections.	1073
1033	– <b>ă</b> otherwise. Eg: <b>bắt</b> tay, <b>bằng</b> lòng, may		
1034	<b>mắn</b> .	<b>C.1 Northern Dialect</b>	1074
1035	• The monophthong /i/ is written as:	According to (Hạo, 1998), Northern dialect speech	1075
1036	– <b>i</b> if the rhyme has a final consonant. Eg:	can be found in provinces from Thanh Hóa to	1076
1037	<b>lúu</b> rít.	the northern borders of Vietnam. In particular,	1077
1038	– <b>y</b> if the rhyme has no consonant. Eg: <b>kỷ</b>	the Northern dialect is distributed in Thanh Hóa,	1078
1039	<b>luật</b> , <b>lý</b> do.	Hà Giang, Cao Bằng, Bắc Kạn, Lạng Sơn, Tuyên	1079
1040	• The final /i/ is written as:	Quang, Thái Nguyên, Phú Thọ, Bắc Giang, Quảng	1080
1041	– <b>i</b> if the rhyme has the front vowel /a/, the	Ninh, Lào Cai, Lai Châu, Yên Bái, Điện Biên, Sơn	1081
1042	central vowel /u, uə/ or the back vowels	La, Hòa Bình, Hà Nội, Hà Nam, Bắc Ninh, Hải	1082
1043	/u, o, ɔ, uo/ as the vowel. Eg: <b>mãi</b> mê,	Dương, Hải Phòng, Hưng Yên, Nam Định, Ninh	1083
1044	<b>mui</b> thuyền, <b>hỏi</b> han, <b>mỗi</b> chài, <b>gửi</b> gắm,	Bình, Thái Bình, and Vĩnh Phúc provinces.	1084
1045	<b>lười</b> biếng, <b>nuôi</b> nấng.	In initials, Northern dialect does not give distinct	1085
		pronunciation on the following phonemes (Hạo,	1086
		1998; Châu, 2002) (the writing forms are provided	1087
		in the parentheses):	1088
		• /s/ (x) and /ʃ/ (s) are pronounced as [s].	1089
		• /ç/ (ch) and /ʈ/ (tr) are pronounced as [ç].	1090
		• /z/ (gi), /r/ (r), and /j/ (d) are pronounced as	1091
		[z].	1092

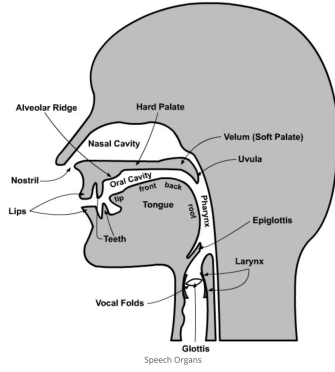


Figure 9: The speech organs.

- In Hà Nội, Vĩnh Phúc, Bắc Ninh, Hà Đông, and Hưng Yên, there is a mispronunciation between /l/ (l) and /n/ (n).

In rhymes, this dialect has the general following features (Hạo, 1998; Châu, 2002):

- The glide /ɥ/ does not follow /u, o, ɔ/.
- The semivowel /j/ does not follow the front vowels /i, e, ɛ/ except for the cases where the final consonants are velar.
- The back diphthongs /uɪ, uə/ when being followed by the glide /ɥ/ they are transmitted to the central diphthongs [i, iə]. That is the reason why the following words **con hương, ly rượu, cấp cứu, âm mưu** are pronounced approximately the same as **con hiêu, ly riệu, cấp kúu, âm miu**, respectively.

In addition, this dialect exhibits two phonetic features on rhymes, including front vowels /i, e, ɛ/ or rounded back vowels /u, o, ɔ/ followed by velar consonants /ŋ, k/. The front vowels require the tongue to be raised and forward to the teeth in the mouth, while the velar consonants require the tongue to be raised and backward to the velum. This results in the placement of the tongue in the middle of the mouth, not forward to the teeth or backward to the velum (Figure 9). From that on, the velar consonants are palatalized to become [j, c] while the front vowels are transmitted to the central vowels [i, ə, ɜ]. The gliding between the central vowels and the palatalized consonants introduces a glide phone [j] in the middle. To this end, we have the following variations for the rhymes having front vowels /i, e, ɛ/ followed by velar consonants /ŋ, k/

(the writing forms are provided in the parentheses) (Hạo, 1998; Châu, 2002):

- /iŋ/ → [iʲŋ] (inh).
- /eŋ/ → [əʲŋ] (ênh).
- /ɛŋ/ → [ɜʲŋ] (anh).
- /ik/ → [iʲc] (ich).
- /ek/ → [əʲc] (êch).
- /ɛk/ → [ɜʲc] (ach).

In the cases of the rounded back vowels /u, o, ɔ/, the velar finals are labialized because of the effectiveness of labial vowels. These vowels are shortened according to the existence of the velar consonants, hence they lost the labial factor at the beginning. To this end, we have the following variations for the rhymes having rounded back vowels /u, o, ɔ/ followed by velar consonants /ŋ, k/ (the writing forms are provided in the parentheses) (Hạo, 1998; Châu, 2002):

- /uŋ/ → [ʷuŋᵐ] (ung).
- /oŋ/ → [ʷoŋᵐ] (ông).
- /ɔŋ/ → [ʷɔŋᵐ] (ong).
- /uk/ → [ʷʊkᵖ] (uc).
- /ok/ → [ʷʊkᵖ] (ôc).
- /ɔk/ → [ʷɔkᵖ] (oc).

## C.2 Central Dialect

The Central dialect is distributed in more limited regions than the other two dialects. This dialect can be found in Nghệ An, Hà Tĩnh, Quảng Bình, Quảng Trị, and Huế provinces. Moreover, this dialect has a particular minor dialect for each province, which makes the Central dialect the most varied one in Vietnam. All minor variations of the Central dialect share the same features in initials and tones, while the rhymes are the only factor that makes them different from each other.

In general, this dialect has five tones:

- In Quảng Bình, Quảng Trị, and Huế, there is no distinction between tone /Hᵛ/ and /Hʔᵛ/.
- In Nghệ Tĩnh, there is no distinction between tone /Hʔᵛ/ and /Hʔᵛ/.

1166	However, this dialect shares the same features relevant to initials in all provinces:	<b>C.2.3 Minor Dialect in Quảng Bình</b>	1203
1167		In Quảng Bình, for rhymes having the nasal consonants /n, t/ following the front vowels /i, e, ε/, these vowels are pronounced longer than usual:	1204
1168	• /j, z/ → [z].		1205
1169	• /tʃ/ → [t].	• /in/ → [i:n].	1206
1170	• /ç/ → [c].	• /en/ → [e:n].	1207
1171	The most significant factors that separate the minor dialects of the Central dialect are rhymes having front vowels /i, e, ε/ followed by the velar consonants /ŋ, k/. Rhymes having the rounded back vowels /u, o, ɔ/ followed by the velar consonants in the Central dialect share the same color as those in the Northern dialect.	• /εn/ → [ε:n].	1208
1172		• /it/ → [i:t].	1209
1173		• /et/ → [ε:t].	1210
1174		• /εt/ → [ε:t].	1211
1175			1212
1176		However, these vowels (i:, e:, ε:) become shorter (i, e, ε) when being followed by the velar consonants /ŋ, k/. From that on, the distinction between these pairs /in/ - /iŋ/, /en/ - /eŋ/, /εn/ - /εŋ/, /it/ - /ik/, /et/ - /ek/, /εt/ - /εk/ is not necessarily dependent on the final but on the length of the vowel. To this end, this dialect replaces the velar consonants /ŋ, k/ by the nasal consonants /n, t/. These result in the following variations:	1213
1177			1214
1178	<b>C.2.1 Minor Dialect in Nghệ An and Hà Tĩnh</b>		1215
1179	Rhymes in Nghệ An and Hà Tĩnh reflect the same rules as the Northern dialect for those having front vowels followed by velar consonants. In these rhymes, the velar consonants /ŋ, k/ are palatalized to become [j, c], but the front vowels keep their original form. To this end, we have:		1216
1180		• /iŋ/ → [iɲ].	1217
1181		• /eŋ/ → [eɲ].	1218
1182		• /εŋ/ → [εɲ].	1219
1183		• /ik/ → [ic].	1220
1184		• /ek/ → [ec].	1221
1185		• /εk/ → [εc].	1222
1186			1223
1187			1224
1188			1225
1189			1226
1190			1227
1191	<b>C.2.2 Minor Dialect in Quảng Trị</b>	<b>C.2.4 Minor Dialect in Huế</b>	1228
1192	However, the rhymes mentioned in the above Section C.2.1 are pronounced slightly differently from those in Quảng Trị. In this province, both front vowels and velar consonants keep their original form:	In this minor dialect, for rhymes having front vowels /i, e, ε/ followed by velar consonants /ŋ, k/, these vowels are transmitted to the central vowels [i, ə, ɜ]. However, as the distinction between these pairs /iŋ/ - /in/, /ik/ - /it/, /eŋ/ - /en/, /ek/ - /et/, /εŋ/ - /εn/, /εk/ - /εt/ largely depends on the difference of the vowels, hence having the velar consonants for these rhymes /iŋ, ik, eŋ, ek, εŋ, εk/ is not necessary. To this end, these velar consonants are replaced by the respective nasal consonants /n, t/. From then on, we have:	1229
1193			1230
1194			1231
1195			1232
1196			1233
1197	• /iŋ/ → [iɲ].	• /iŋ/ → [in].	1234
1198	• /eŋ/ → [eɲ].	• /ik/ → [it].	1235
1199	• /εŋ/ → [εɲ].	• /eŋ/ → [əɲ].	1236
1200	• /ik/ → [ik].		1237
1201	• /ek/ → [ek].		1238
1202	• /εk/ → [εk].		1239

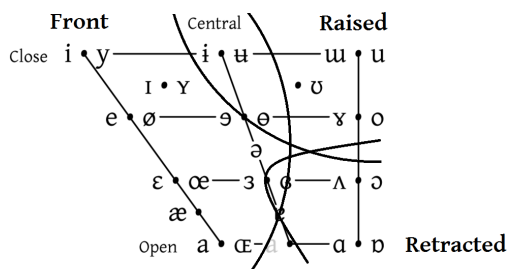


Figure 10: Vowel diagram.

- /ek/ → [ɛt].
- /ɛŋ/ → [ɛn].
- /ɛk/ → [ɛt].

The front vowels /i, e, ɛ/ preceding the nasal consonants /n, t/ become longer. That is:

- /in/ → [i:n].
- /it/ → [i:t].
- /en/ → [e:n].
- /et/ → [e:t].
- /ɛn/ → [ɛ:n].
- /ɛt/ → [ɛ:t].

From these analyses, we can see that the Huế dialect exhibits the same behavior as the dialect in Quảng Bình. However, the way of pronouncing rhymes having the front vowels followed by the velar consonants is different in these two minor dialects. These contribute to the huge diversity of dialects in Vietnamese, which results in a significant challenge in modeling multi-dialect speech ASR in this language.

### C.3 Southern Dialect

(Hạo, 1998; Châu, 2002) determined that the southern dialect can be found from the Đà Nẵng province to the southern border of Vietnam. In particular, the Southern dialect is distributed in Đà Nẵng, Quảng Nam, Quảng Ngãi, Kom Tum, Bình Định, Gia Lai, Phú Yên, Đắk Lắk, Khánh Hòa, Lâm Đồng, Ninh Thuận, Bình Thuận, Bình Phước, Đồng Nai, Bình Dương, Bà Rịa - Vũng Tàu, Hồ Chí Minh, Long An, Đồng Tháp, Tiền Giang, Bến Tre, Vĩnh Long, An Giang, Trà Vinh, Cần Thơ, Sóc Trăng, Kiên Giang, Bạc Liêu, and Cà Mau.

This dialect has five tones rather than six tones as the Northern dialect: there is no distinction between tone /H/ and /H?/. In initials, Southern dialect does not give any difference between the following phonemes (Châu, 2002; Hạo, 1998):

- /z/ (gi), /j/ (d), and /v/ are all pronounced as [j].
- /k<sup>w</sup>/ (qu) is pronounced as [w].
- This dialect also exhibits the same phonetic rule as the Northern dialect, where the velar consonants /ŋ, k/ follow the front vowels /i, e, ɛ/. These velar consonants are palatalized to become [ɲ, c], but the front vowels /i, e, ɛ/ keep their original.
- In provinces Quảng Nam, Quảng Ngãi, Kon Tum, Gia Lai, Đắk Lắk, Bình Định, Phú Yên, Khánh Hòa, Ninh Thuận, Bình Thuận, and Lâm Đồng: /t̚/ (tr) is pronounced as [t̚]. While in other provinces, /t̚/ (tr) and /c̚/ (ch) is pronounced as [c], /s̚/ (s) and /s̚/ (x) are all pronounced as [s].

In rhymes, when these diphthongs /ie, uə/ followed by the semivowel /j/ or labial consonants /m, p, ɸ/, they are converted to the monophthongs in the same row of the vowel diagram (Figure 10). In particular, in this dialect we have:

- /iep/ → [ip].
- /uəp/ → [ɸp].
- /iem/ → [im].
- /uəm/ → [ɸm].
- /ieu/ → [iu].
- /uəɸ/ → [ɸɸ].

This is the reason why in Southern dialect, the following words *tiếp tục*, *quả mướp*, *đánh chiếm*, *ươn mềm*, *tình yêu*, *con hươu* are approximately pronounced as *típ tục*, *quả múp*, *đánh chím*, *ưm mềm*, *tình iu*, *con hutu*, respectively.

In addition to the general features of Southern dialect, two minor dialects share the same characteristics of Southern dialect but have their own particular features. The first one is the dialect distributed in provinces Quảng Nam and Quảng Ngãi, and the second one is the dialect found in the Mekong Delta.

### C.3.1 Minor Dialect of Southern Dialect in Quảng Nam and Quảng Ngãi

In provinces Quảng Nam and Quảng Ngãi, the residents additionally pronounce the vowels differently from the Southern dialect. In rhymes without final consonants, the close monophthongs /i, u, ɯ/ are converted into diphthongs whose the first components are the respective more-open vowels (the row below of the vowel diagram, Figure 10), while the second components are the respective glides. In particular, we have:

- /i/ → [ij].
- /u/ → [ɯw].
- /ɯ/ → [ɣɯ].

In rhymes having /i/ or /u/ as the final consonants, the diphthongs /uo, ɯə, ie/ are transmitted into the monophthongs on the same row of the vowel diagram, but these monophthongs are longer to match the length of the original diphthongs:

- /uo/ → [u:].
- /ɯə/ → [ɯ:].
- /ie/ → [i:].

Moreover, the following two rhymes are pronounced totally differently:

- /oi/ → [uə].
- /ai/ → [ae].

This is the reason why people in these two provinces pronounce the following words *nói* and *hai* approximately the same as *núa* and *he*, respectively.

Furthermore, in rhymes having labial finals /m, p/ and the rounded semivowel /ɯ/, the vowels in this dialect are pronounced as follows:

- /a/ → [ɔ].
- /ɔ/ → [o].
- /o/ → [ɣ].
- /u/ → [ɯ].
- /ǎ, ǝ/ → [a].
- In rhymes having /m/ as the final, /uo/ → [ɯ].

In rhymes having the velar consonants as the final, while the front vowels /i, e, ɛ/ do not vary, the following vowel changed as follows:

- /a/ → [ɑ:].
- /ɔ/ → [ɑ:<sup>m</sup>].
- /o/ → [Λ].
- /ǎ/ → [a].
- /uo/ → [u:].
- /ie/ → [i:].
- /ɯə/ → [ɯ:].

### C.3.2 Minor Dialect of Southern Dialect in Mekong Delta

Mekong Delta includes Bình Dương, Bình Phước, Đồng Nai, Bà Rịa - Vũng Tàu, Tây Ninh, Hồ Chí Minh, Long An, Tiền Giang, Hậu Giang, Long An, Bến Tre, Đồng Tháp, Vĩnh Long, Trà Vinh, Cần Thơ, Sóc Trăng, An Giang, Kiên Giang, Bạc Liêu, and Cà Mau. These provinces have a unique dialect that shares many features with the Southern dialect, with their own distinctions.

In particular, the front vowels /i, e/ followed by velar consonants /ŋ, k/ do not transmitted to the respective central vowels, but the velar consonants are replaced by the nasal consonants /ŋ, k/ → /n, t/, which is similar to the Central dialect without changing in length of the front vowels. That is:

- /iŋ/ → [in].
- /ik/ → [it].

In addition, the front vowel /e/ followed by the velar consonants /ŋ, k/ is converted to the respective central vowel because of the velar consonants, but as the velar consonants /ŋ, k/ are replaced by the nasal consonants [n, t], which results in this vowel is longer. From then on, we have:

- /eŋ/ → [ə:n].
- /ek/ → [ə:t].

However, the alveolar vowels /n, t/ following vowel /a/, and the central vowels /ə, ɛ/ become the velar consonants /ŋ, k/. This is the reason why in this dialect, these words *buôn bán*, *ít ỏi*, *ân cần*, *chen lán* are pronounced approximately the same as *buôn bán*, *ích ỏi*, *âng cần*, *cheng lán*, respectively.

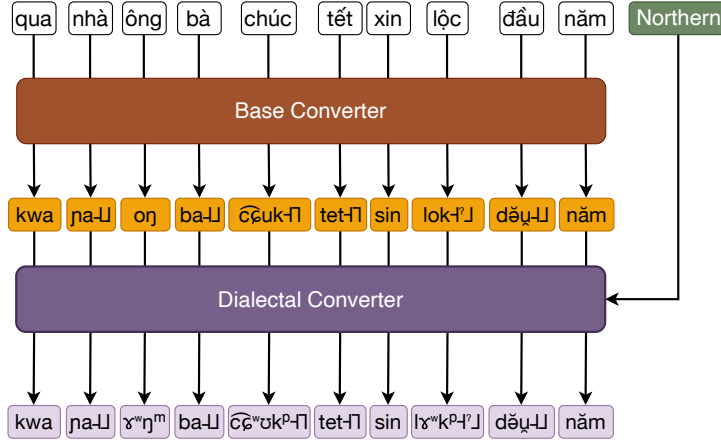


Figure 11: Dialect-aware tokenization algorithm.

On the other hand, the rounded back ones /u, o, ɔ/) have variations depending on the following consonants. If these vowels are followed by the nasal consonants /n, t/, they keep their rounded features while their length is longer /u, o, ɔ/ → [uː, oː, ɔː]. In the case these vowels are followed by velar consonants /ŋ, k/, they lose the rounded characteristics as in the Northern dialect with a little bit difference, while these velar consonants are labialized to become [ŋᵐ, kᵑ]. Moreover, the nasal consonants [n, t] followed by these vowels are replaced by the velar-labialized consonants [ŋᵐ, kᵑ], respectively. In particular, we have:

- /un, uŋ/ → [ʷuŋᵐ].
- /ut, uk/ → [ʷʊkᵑ].
- /on, oŋ/ → [ɣŋᵐ].
- /ɔt, ɔk/ → [ʌkᵑ].

## D Appendix: Dialect-Aware Tokenization Algorithm

We have analyzed in detail the phonetic and orthographic features of Vietnamese in Section B and comprehensively describe the multi-dialect variations of this language in Section C. From these analysis, we develop an Dialect-aware Tokenization algorithm (DiaToken). DiaToken receives a line of text which is the transcription of the given multi-dialect voice and the province ID, then processes and returns the sequence of IPA (International Phonetic Alphabet) characters describing the dialect speech sound of the given audio correspond to the dialect. DiaToken has two phases (Figure 11): (1) Retrieving phonemes from graphemes given the

text of the transcript (the Base Converter module) and (2) Converting phonemes to IPA characters respective to dialectal phones (the Dialectual Converter module).

### D.1 Base Converter Module

**Algorithm 1:** The algorithm for converting text to phonemes.

**Data:** Transcript of the audio

$$w = (w_1, w_2, \dots, w_n).$$

**Result:** A sequence of syllables

$p = (p_1, p_2, \dots, p_n)$  of the given input transcript

$w = (w_1, w_2, \dots, w_n)$ . Each

phoneme  $p_i =$

$$(p_i^{init}, p_i^{glide}, p_i^{vowel}, p_i^{final}, p_i^{tone})$$

is a triplet of IPA for the initial, rhyme, and tone.

```

1 phonemes ← an empty list [];
2 for  $W$  in  $w$  do
3    $p_{tone}, W \leftarrow \text{get\_tone}(W)$ ;
4    $p_{initial}, W \leftarrow \text{get\_initial}(W)$ ;
5    $p_{glide}, W \leftarrow \text{get\_glide}(W)$ ;
6    $p_{vowel}, W \leftarrow \text{get\_vowel}(W)$ ;
7    $p_{final} \leftarrow \text{get\_final}(W)$ ;
8   phonemes ← Append
   ( $p_{init}, p_{rhyme}, p_{tone}$ );
9 end
10 return phonemes;
```

We describe in Alg. 1 a overview of algorithm of the Base Converter module for converting orthographic form of the transcript to the sequence of phonemes. In this algorithm, we repre-

---

**Algorithm 2:** Algorithm for determining initial of a word.

---

```

1 Function get_initial(W):
  Input: A word W in Vietnamese
  Output: The initial i of W
2  onsets  $\leftarrow$  [ngh, tr, th, ph, nh, ng, kh, gi,
   gh, ch, q, đ, x, v, t, s, r, n, m, l, k, h, g,
   d, c, b];
3  i  $\leftarrow$  None;
4  for onset in onsets do
5    if W starts with onset then
6      /* Words starting with "qu"
7       are kept the onset for
8       later process. */
9      if onset  $\neq$  k then
10     | Remove onset from W;
11     end
12     i  $\leftarrow$  IPA character of onset
13     according to Appendix B;
14   end
15   break;
16 end

```

---

1442 sent each word as a vector of five syllabic compo-  
1443 nents: initial, glide, vowel, final, and tone. From  
1444 that on, these five phonemes of each word will  
1445 be mapped to the dialectal phone according to the  
1446 given province of the audio.

---

**Algorithm 3:** Algorithm for determining final of a word.

---

```

1 Function get_final(W):
  Input: A word W in Vietnamese
  Output: The final f of W
2  codas  $\leftarrow$  [ng, nh, ch, u, n, o, p, c, m, y,
   i, t];
3  if W in codas then
4    | return W;
5  end
6  return None;
7 end

```

---

1447 Moreover, as described in Alg. 5, Alg. 2, Alg.  
1448 6, Alg. 4, and Alg. 3, the computational complex-  
1449 ity of the Base Convert (Figure 11) is linear  $\mathcal{O}(n)$   
1450 while using a fix sized of vocabulary of phonemes  
1451 but tokenizing unlimited number of Vietnamese  
1452 words.

---

**Algorithm 4:** Algorithm for determining vowel of a word.

---

```

1 Function get_vowel(W):
  Input: A word W in Vietnamese
  Output: The vowel v of W
2  nuclei  $\leftarrow$  [oo, uo, ua, uo, ua, ie, ye, ia,
   ya, e, e, u, u, o, i, y, o, a, o, a];
3  for nucleus in nuclei do
4    if W starts with nucleus then
5      | v  $\leftarrow$  IPA character of nucleus
6      according to Appendix B;
7      Remove nucleus from W;
8      return v, W;
9    end
10 end
11 return None, W;

```

---



---

**Algorithm 5:** Algorithm for determining tone of a word.

---

```

1 Function get_tone(W):
  Input: A word W in Vietnamese
  Output: The tone t of W
2 end
3 t  $\leftarrow$  tone of word according to Appendix B;
4 Remove tone from W;
5 return t, W;

```

---

---

**Algorithm 6:** Algorithm for determining glide of a word.

---

```

1 Function get_glide(W):
   Input: A word W in Vietnamese
   Output: The glide g of W
2 if W starts with qu then
3   | Remove qu from W;
4   | return u, W;
5 end
6 for case in [oa, oă, oe] do
7   | if W starts with case then
8   | | Remove o from W;
9   | | return u, W;
10  | end
11 end
12 for case in [uê, uy, uơ, ua, uâ, uya] do
13  | if W starts with case then
14  | | Remove u from W;
15  | | return u, W
16  | end
17 end
18 return None, W;
19 end

```

---

## 1453 D.2 Dialectal Converter Module

1454 Having the phoneme representation for the given  
1455 transcript, the DiaToken continues to retrieve the  
1456 dialectal representation of the given transcript ac-  
1457 cording to the name of the province. Vietnam has  
1458 63 provinces and their speech sound are grouped  
1459 into three large dialects as described in Appendix  
1460 C. DiaToken receives the name of the province then  
1461 determines the respective dialect and finally con-  
1462 vert every phonemes to the corresponding phones.  
1463 The overall algorithm is described in Alg. 7.

1464 From Alg. 8, Alg. 9, and Alg. 10, the compu-  
1465 tational complexity is linear  $\mathcal{O}(n)$ , which means  
1466 the complexity of the Dialectal Converter Mod-  
1467 ule is  $\mathcal{O}(n)$ . Finally, the overall complexity of the  
1468 DiaToken algorithm is  $\mathcal{O}(n)$  which is efficient for  
1469 tokenizing transcript in the context of multi-dialect.

---

**Algorithm 7:** The algorithm for converting phonemes to dialectal phones

---

**Data:**

- The name of the province **ProID**.
- A sequence  $p = (p_1, p_2, \dots, p_n)$  representing phonemic syllables of the transcript  $w = (w_1, w_2, \dots, w_n)$ . Each  $p_i = (p_i^{init}, p_i^{glide}, p_i^{vowel}, p_i^{final}, p_i^{tone})$  is a vector of five syllabic components of word  $w_i$ .

**Result:** A sequence

$$ph = (ph_1, ph_2, \dots, ph_n)$$

representing dialectal syllables of the transcript

$$w = (w_1, w_2, \dots, w_n).$$

Each  $ph_i = (ph_i^{init}, ph_i^{rhyme}, ph_i^{tone})$  is a vector of three dialectal phonetic components of word  $w_i$ .

```

1 northern  $\leftarrow$  set of provinces having the
   Northern dialect according to Appendix C;
2 middle  $\leftarrow$  set of provinces having the
   Central dialect according to Appendix C;
3 southern  $\leftarrow$  set of provinces having the
   Southern dialect according to Appendix C;
4 if ProID  $\in$  northern then
5   | pth  $\leftarrow$  get_northern(p);
6 end
7 if ProID  $\in$  middle then
8   | pth  $\leftarrow$  get_middle(p);
9 end
10 if ProID  $\in$  southern then
11  | pth  $\leftarrow$  get_southern(p);
12 end
13 return ph;

```

---

---

**Algorithm 8:** Algorithm for determining Northern dialect of a syllable.

---

```
1 Function get_nothern(s):  
   Input: A sequence of phonemic syllables  $s = (s_1, s_2, \dots, s_n)$   
   Output: A sequence of Northern phonetic syllables  $s = (s_1, s_2, \dots, s_n)$   
2   for  $s_i$  in s do  
3      $s_i^{init} \leftarrow s_i$ ;  
4      $s_i^{init} \leftarrow$  IPA character for Northern dialect of the initial  $s_i^{init}$  according to Appendix C.1;  
5      $s_i^{glide} \leftarrow$  IPA character for Northern dialect of the glide  $s_i^{glide}$  according to Appendix C.1;  
6      $s_i^{vowel} \leftarrow$  IPA character for Northern dialect of the vowel  $s_i^{vowel}$  according to Appendix C.1;  
7      $s_i^{final} \leftarrow$  IPA character for Northern dialect of the final  $s_i^{final}$  according to Appendix C.1;  
8      $s_i^{tone} \leftarrow$  IPA character for Northern dialect of the tone  $s_i^{tone}$  according to Appendix C.1;  
9      $s_i^{rhyme} \leftarrow s_i^{glide} \oplus s_i^{vowel} \oplus s_i^{final}$ ; /*  $\oplus$  represents the concatenation operator */  
10     $s_i \leftarrow (s_i^{init}, s_i^{rhyme}, s_i^{tone})$ ;  
11  end  
12 end
```

---

---

**Algorithm 9:** Algorithm for determining Central dialect of a syllable.

---

```
1 Function get_nothern(s):  
   Input: A sequence of phonemic syllables  $s = (s_1, s_2, \dots, s_n)$   
   Output: A sequence of Central phonetic syllables  $s = (s_1, s_2, \dots, s_n)$   
2   for  $s_i$  in s do  
3      $s_i^{init} \leftarrow s_i$ ;  
4      $s_i^{init} \leftarrow$  IPA character for Central dialect of the initial  $s_i^{init}$  according to Appendix C.2;  
5      $s_i^{glide} \leftarrow$  IPA character for Central dialect of the glide  $s_i^{glide}$  according to Appendix C.2;  
6      $s_i^{vowel} \leftarrow$  IPA character for Central dialect of the vowel  $s_i^{vowel}$  according to Appendix C.2;  
7      $s_i^{final} \leftarrow$  IPA character for Central dialect of the final  $s_i^{final}$  according to Appendix C.2;  
8      $s_i^{tone} \leftarrow$  IPA character for Central dialect of the tone  $s_i^{tone}$  according to Appendix C.2;  
9      $s_i^{rhyme} \leftarrow s_i^{glide} \oplus s_i^{vowel} \oplus s_i^{final}$ ; /*  $\oplus$  represents the concatenation operator */  
10     $s_i \leftarrow (s_i^{init}, s_i^{rhyme}, s_i^{tone})$ ;  
11  end  
12 end
```

---

---

**Algorithm 10:** Algorithm for determining Southern dialect of a syllable.

---

```
1 Function get_nothern(s):  
   Input: A sequence of phonemic syllables  $s = (s_1, s_2, \dots, s_n)$   
   Output: A sequence of Southern phonetic syllables  $s = (s_1, s_2, \dots, s_n)$   
2   for  $s_i$  in s do  
3      $s_i^{init} \leftarrow s_i$ ;  
4      $s_i^{init} \leftarrow$  IPA character for Southern dialect of the initial  $s_i^{init}$  according to Appendix C.3;  
5      $s_i^{glide} \leftarrow$  IPA character for Southern dialect of the glide  $s_i^{glide}$  according to Appendix C.3;  
6      $s_i^{vowel} \leftarrow$  IPA character for Southern dialect of the vowel  $s_i^{vowel}$  according to Appendix C.3;  
7      $s_i^{final} \leftarrow$  IPA character for Southern dialect of the final  $s_i^{final}$  according to Appendix C.3;  
8      $s_i^{tone} \leftarrow$  IPA character for Southern dialect of the tone  $s_i^{tone}$  according to Appendix C.3;  
9      $s_i^{rhyme} \leftarrow s_i^{glide} \oplus s_i^{vowel} \oplus s_i^{final}$ ; /*  $\oplus$  represents the concatenation operator */  
10     $s_i \leftarrow (s_i^{init}, s_i^{rhyme}, s_i^{tone})$ ;  
11  end  
12 end
```

---