



# COMPUTATIONAL LINGUISTIC MATERIAL FOR VIETNAMESE SPEECH PROCESSING: APPLYING IN VIETNAMESE TEXT-TO-SPEECH

Pham Van Dong  
Faculty of Information Technology  
Hanoi University of Mining and Geology  
Hanoi, Vietnam

Nguyen Tien Thanh  
Viettel CyberSpace Center  
Viettel  
Hanoi, Vietnam

Mac Dang Khoa  
VinBigdata  
VinGroup  
Hanoi, Vietnam

Tran Do Dat  
MOST  
Ministry of Science and Technology of Vietnam  
Hanoi, Vietnam

Do Thi Ngoc Diep  
MICA  
Hanoi University of Science and Technology  
Hanoi, Vietnam

Vu Thi Hai Ha  
Phonetics Department  
Vietnam Institute of Linguistics  
Hanoi, Vietnam

Dang Thanh Mai  
Foreign language department  
Hanoi University of Mining and Geology  
Hanoi, Vietnam

**Abstract:** The motivation of this paper is to propose a set of best-quality linguistic materials for Vietnamese speech processing, which can be used for Vietnamese TTS and ASR problems. This proposed material includes: (1) a pronunciation dictionary, which adapts from X-SAMPA, (2) a rule-based grapheme to phoneme for Vietnamese. In order to test and evaluate, we have built a Vietnamese TTS system based on the Merlin engine, using the above materials, and evaluating the quality of speech and the accuracy of pronunciation. The results show that the applicability of these materials is favorable for further research and development on Vietnamese speech processing.

**Keywords:** Text-to-speech, Dictionary, Grapheme-to-Phoneme, X-SAMPA, computer coding, speech processing, Vietnamese

## I. INTRODUCTION

Build a speech processing system consisting of many components, of which one of the primary components is linguistic processing, requiring many linguistic materials such as pronunciation dictionaries, etc. This critical part determines pronunciation (in TTS) and correct speech recognition (ASR)[1].

Vietnamese speech processing has been studied since the early 2000s, and so far, it has achieved many results in two main research areas: speech recognition and synthesis. MICA has researched and created some typical application systems such as Hoa Sung (Speech synthesis in the Vietnamese language), and VOVA (High-quality Vietnamese TTS engine) for Android). Leading Vietnamese companies such as Vais (Vietnam AI System), Viettel<sup>1</sup>, Zalo<sup>2</sup>, Fpt<sup>3</sup>, and the world's

leading company Google have created many applications such as Google Text to Speech<sup>4</sup>.

Although there are many pieces of research and products, there is no unified set of public linguistic materials for Vietnamese. Usually, product research and development parties build their independent linguistic materials, so sharing, and inheriting are difficult. Therefore, the motivation of this paper is to propose a set of materials with the best quality for Vietnamese processing, which can be used for Vietnamese TTS and ASR problems.

This paper is organized as follows. After presenting the introduction and the linguistic materials in section 2, section 3 will show the experiment. Section 4 presents the evaluation process. Furthermore, the final section presents the conclusions and future directions of the study. Proposed linguistic materials

This section will present the main contributions, including the Vietnamese syllable list, the Vietnamese phone list, and the Vietnamese computer encoding.

<sup>1</sup> <https://viettelgroup.ai/service/tts>

<sup>2</sup> <https://zalo.ai/products/text-to-audio-converter>

<sup>3</sup> <https://fpt.ai/vi/tts>

<sup>4</sup> <https://cloud.google.com/text-to-speech>

## A. Unified XSAMPA phone list for Vietnamese

### 1) The overview of Vietnamese phonetic

According to A. G. Haudricourt [2], in 1954, the Vietnamese-Muong language group was a language or dialect in the early period of Christianity. Then, through the interaction process with the Chinese language and especially the Tai-Kadai language, the language with a highly developed tone system, the tone system in Vietnamese appeared and today, according to the law of forming tones. The appearance of the tones began around the 6th century (northern part of Vietnamese history) with three tones and steady development around the 12th century (Ly) with six tones. Then some early consonants changed to this day. During the change process, the final consonants change the syllable ends, and the consonant ends change from confusion to tangible.

Vietnamese has divided into three main dialects: North, Central, and South. In fact, because of the very long history of development, the North dialect is considered to be a national standard script. It consists of provinces from Northern Vietnam to Thanh Hoa province. The Central dialect refers to that spoken between Nghe An province and Da Nang province. The South dialect is spoken from the South of Da Nang province to the whole Southern province. A mixture of dialects in an area has a boundary between two dialects. For example, the dialect of Thanh Hoa province is nearly similar to the North and the Central ones, sometimes also arranged into the North-central dialect. The differences between dialects are primarily reflected in the tones and consonants. Three main varieties of Vietnamese, North, Central, and South, which are slightly different from each other, are described below.

Table I. Initials in Northern Vietnamese dialect [3]

		L abial	Den tal/ Alv eolar	P alatal	V elar	G lottal
Nasal		m	n	ɲ	N	
Plosive /Affricate	unaspirated		t	c	k	/
	aspirated	p <sup>h</sup>	t <sup>h</sup>		k <sup>h</sup>	
	implosive	ɓ	ɗ			
Fricative	voiceless	f	s	ʃ	ɣ	h
	voiced	v	z		ʒ	
Approximant			l			

The dialects differ in using tone systems. The Northern dialect uses six tones (level, falling, rising, drop, broken, curve), while the Central and the Southern ones have only five tones, especially some regions of the Central have only four tones. The Central and Southern speakers do not discern clearly between broken and curved. They almost used rising instead of the broken tone. Some North-Central varieties maintain the rising but have a merge of curve and drop. The vowels and consonant pronunciation are also distinct. In the northern dialect, there are six tones.

Northern Vietnamese has no retroflexes /ɳ, ʎ, ʈ/, nor the rhymes /oɰw, ɔɰ↗w/. It is the only area that entirely exists the eight final consonants in spelling, including -p, -t, -k, -m, -n, -ng, -u, and -i. d, gi, and r are all pronounced /z/. ch and tr are pronounced /c/, while x and s are pronounced /s/.

Some rural speakers merge /l/ and /n/ into /l/, although this is not considered standard.

Central Vietnamese is known for preserving archaic characteristics in various ways, such as the most significant number of initial consonants, the high amount of ancient words, sub-local tonal variations within the area, and so on (Hoang TC 2003). In the central and southern dialects, retroflexes /ɳ, ʎ, ʈ/,

but only five tones exist. Five tones system differs from the Northern Vietnamese system in quantity and quality. There is confusion between the curve and the broken, not distinguishing the broken from the drop. The initial consonant system consists of 24 consonants. In many dialects, there are two aspirated consonants [p<sup>h</sup>, k<sup>h</sup>] instead of two fricatives [f, ɣ]. The front vowel [i, e, E], and back vowel [u, o, ɔ] in the Central dialect tend to shift slightly towards the middle than in the Northern dialect.

Table II. Initials in Central Vietnamese dialect [3]

		Labial	Dental /Alveolar	Retroflex	Palatal	Velar	Glottal
Nasal		m	n		ɲ	N	
Plosive /Affricate	unaspirated		t	ɳ	c	k	/
	aspirated	p <sup>h</sup>	t <sup>h</sup>			k <sup>h</sup>	
	implosive	ɓ	ɗ				
Fricative	voiceless	f	s	ʃ		ɣ	h
	voiced	v	z			ʒ	
Approximant			l				

Southern dialect, where there is less variation in dialect boundaries. In general, there are clearly 5 tones, exist retroflexes /ɳ, ʎ, ʈ/. Confusing the finals -n, -t with -N, -k. The initials system of the Southern dialect has some differences from the initials system in other dialects of Vietnamese. First, the appearance of two the initials /j/ and /w/, in which, /j/ corresponds to three initials /v/, /z/, /ʒ/ in other dialects, while /w/ in the Southern dialect corresponds to the combinations of consonant and glide in other dialects /ʃw/, /hw/, /kw/, /Nw/. Second, confuse /s/ with /ʃ/, /c/ with /ɲ/. There is confusion between the curve and the broken tone.

Table III. Initials in Southern Vietnamese dialect [3]

		Labial	Dental /Alveolar	Retroflex	Palatal	Velar	Glottal
Nasal		m	n		ɲ	N	
Plosive /Affricate	unaspirated	p	t	ɳ	c	k	/
	aspirated		t <sup>h</sup>				
	implosive	ɓ	ɗ				
Fricative	voiceless	f	s	ʃ		ɣ	h
	voiced		z				
Approximant			l	r	j	w	
Rhotic							

The tone of the text is written with the sign: "à" (falling tone), "ã" (broken tone), "á" (curve tone), "ạ" (rising tone), "ạ" (drop tone), described in Fig.1.

Tone 1	Tone 2	Tone 3	Tone 4	Tone 5	Tone 6
ngang	huyền ‘\’	ngã ‘~’	hỏi ‘?’	sắc ‘/’	nặng ‘.’

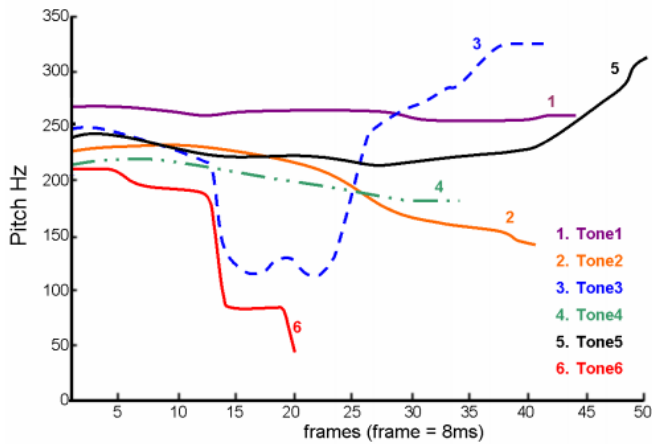


Figure 1 Example of contours of six tones (female subject Pham Ngoc Yen), as described in [4]

Thus, we recommend a listed Vietnamese phone consisting of 19 initial consonants, 12 final consonants, and 13 vowels. Six Vietnamese sounds in our work are recorded in 6 numbers: 0. level tone, 1. falling tone, 2. broken tone, 3. curve tone, 4. rising tone, 5. drop tone.

## 2) Computer encoding

The International Phonetic Alphabet (IPA<sup>5</sup>) is an alphabetic system of phonetic notation based primarily on the Latin script. The International Phonetic Association devised it in the late 19th century as a standardized representation of speech sounds in written form. IPA has the disadvantage that it is challenging to represent and process on computer code. The Speech Assessment Methods Phonetic Alphabet (SAMPA) is a computer-readable phonetic script using 7-bit printable ASCII characters based on the IPA. The Extended Speech Assessment Methods Phonetic Alphabet (X-SAMPA) is a variant of SAMPA developed in 1995 by John C. Wells [5]. It is designed to unify the individual language SAMPA alphabets and extend SAMPA to cover the entire range of characters in the 1993 version of IPA. X-SAMPA is still useful as the basis for an input method for true IPA.

Based on the study of Vietnamese phonetics above, we propose the Vi-XSAMPA phonetic set as an extension of the mapping between IPA and X-SAMPA, serving the coding for the Vietnamese TTS synthesized system due to the following reasons:

- IPA characters have many special characters that make programming difficult, such as  $\text{ŋ}$ ,  $\text{j}$ ,  $\text{ɛ}$ ,  $\text{z}$  ...
- Some X-SAMPA characters will cause errors in the list of questions and Merlin programming, such as the  $\text{@}$  character,  $\text{7}$  ... because it coincides with the programming keyword

The proposal for the Vietnamese phone set was proposed by [6] and aims to apply to speech synthesis. 25 characters for initial/final consonants and 16 for vowels/diphthongs was used in [6]. In this phone set, many characters make programming difficult. They are  $\text{ts}$ ,  $\text{k}_+$ ,  $\text{N}_+$ ,  $\text{7}$ ,  $\text{i@}$ ,  $\text{u@}$ , and  $\text{M@}$ .

We have proposed the VI X-SAMPA table in Table IV. Our phonetic set includes 53 characters, developed from X-SAMPA, and covers the entire phonemic set. For tone, the

numbers 1 - 5 will be used to represent the tone, unlike the Standard X-SAMPA, which uses some symbols which can be challenging to process. We replaced hard-to-program characters with more good characters. They are:  $\text{b}_< \rightarrow \text{b}$ ,  $\text{d}_< \rightarrow \text{d}$ ,  $\text{t}_\text{h} \rightarrow \text{th}$ ,  $\text{a}_\text{X} \rightarrow \text{aX}$ ,  $\text{O}_\text{X} \rightarrow \text{OX}$ ,  $\text{E}_\text{X} \rightarrow \text{EX}$ ,  $\text{7} \rightarrow \text{ow}$ ,  $\text{7}_\text{X} \rightarrow \text{aa}$ ,  $\text{i@} \rightarrow \text{ie}$ ,  $\text{M@} \rightarrow \text{wa}$ ,  $\text{u@} \rightarrow \text{uo}$ ,  $\text{k}_+ \rightarrow \text{kz}$ ,  $\text{k}_\text{p} \rightarrow \text{kp}$ ,  $\text{N}_+ \rightarrow \text{nz}$ ,  $\text{N}_\text{m} \rightarrow \text{Nm}$ ,  $\text{p} \rightarrow \text{pc}$ ,  $\text{t} \rightarrow \text{tc}$ ,  $\text{M} \rightarrow \text{0}$ ,  $\text{L} \rightarrow \text{1}$ ,  $\text{H} \rightarrow \text{4}$ .

Some special symbols in XSAMPA are omitted:  $\text{_<}$ , for example,  $\text{b}_< \Rightarrow \text{b}$  (similar to  $\text{b}$  in English), because these two sounds are close, and also want to be able to later share with the English phone board.

Table IV. VIETNAMESE VI-XSAMPA

NO	IPA	Letter	X-SAMPA	VI-XSAMPA
1	$\text{b}$	$\text{b, p}$	$\text{b}_<$	$\text{b}$
2	$\text{c}$	$\text{ch, tr}$	$\text{c}$	$\text{c}$
3	$\text{d}^f$	$\text{đ}$	$\text{d}_<$	$\text{d}$
4	$\text{f}$	$\text{ph}$	$\text{f}$	$\text{f}$
5	$\text{ɣ}$	$\text{g, gh}$	$\text{G}$	$\text{G}$
6	$\text{h}$	$\text{h}$	$\text{h}$	$\text{h}$
7	$\text{k}$	$\text{c, k, q}$	$\text{k}$	$\text{k}$
8	$\text{l}$	$\text{l}$	$\text{l}$	$\text{l}$
9	$\text{m}$	$\text{m}$	$\text{m}$	$\text{m}$
10	$\text{n}$	$\text{n}$	$\text{n}$	$\text{n}$
11	$\text{j}_\text{n}$	$\text{nh}$	$\text{J}$	$\text{J}$
12	$\text{ŋ}$	$\text{ng, ngh}$	$\text{N}$	$\text{N}$
13	$\text{s}$	$\text{x, s}$	$\text{s}$	$\text{s}$
14	$\text{t}$	$\text{t}$	$\text{t}$	$\text{t}$
15	$\text{t}^h$	$\text{th}$	$\text{t}_\text{h}$	$\text{th}$
16	$\text{v}$	$\text{v}$	$\text{v}$	$\text{v}$
17	$\text{z}$	$\text{d, r, gi}$	$\text{z}$	$\text{z}$
18	$\text{x}$	$\text{kh}$	$\text{X}$	$\text{X}$
19	$\text{w}$	$\text{o, u}$	$\text{w}$	$\text{w}$
20	$\text{a}$	$\text{a}$	$\text{a}$	$\text{a}$
21	$\text{ǣ}$	$\text{a(u, y), ǣ}$	$\text{a}_\text{X}$	$\text{aX}$
22	$\text{o}$	$\text{o, oo}$	$\text{O}$	$\text{O}$
23	$\text{ǫ}$	$\text{o (c, ng)}$	$\text{O}_\text{X}$	$\text{OX}$
24	$\text{e}$	$\text{ê}$	$\text{e}$	$\text{e}$
25	$\text{ɛ}$	$\text{e}$	$\text{E}$	$\text{E}$
26	$\text{ě}$	$\text{a (nh, ch)}$	$\text{E}_\text{X}$	$\text{EX}$
27	$\text{ɤ}$	$\text{σ}$	$\text{7}$	$\text{ow}$
28	$\text{ĩ}$	$\text{â}$	$\text{7}_\text{X}$	$\text{aa}$
29	$\text{i}$	$\text{i, y}$	$\text{i}$	$\text{i}$
30	$\text{u}$	$\text{u}$	$\text{M}$	$\text{M}$
31	$\text{o}$	$\text{ô}$	$\text{o}$	$\text{o}$
32	$\text{u}$	$\text{u}$	$\text{u}$	$\text{u}$
33	$\text{iə}$	$\text{ie}$	$\text{i@}$	$\text{ie}$
34	$\text{uə}$	$\text{ura, uɔ}$	$\text{M@}$	$\text{wa}$



complete phonetic dictionary of more than 18,000 Vietnamese words following the process below, the entire source code and detailed documentation shared on my Github<sup>12</sup>:

Here are some results of our dictionary:

- ba 0 b a
- bai 0 b a j
- bam 0 b a mc
- ban 0 b a nc
- bang 0 b a Nz
- banh 0 b EX Nz
- bao 0 b a wc

The process of creating a dictionary is described in detail in Figure 4

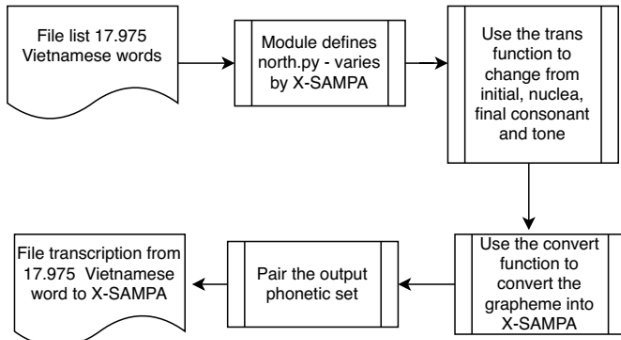


Figure 4. Create dictionary process

## II. EXPERIMENT IN VIETNAMESE TTS

This section aims to applying those materials with one of the fundamental problems of Speech Processing: Text to speech. We also check if the materials can help the TTS system correctly pronounce Vietnamese sounds, especially difficult, rare, and often mispronounced: an, ang, eng, éc, eng, eat, etc. We have prepared the data and trained the TTS system using Merlin Tool.

### A. Data

The training database is a voice database of more than 3 hours of recording, divided into 2.360 sentences, Le Diem recording voice. The recording file is in a studio with a sampling frequency of 16khz.

### B. Merlin Toolkit

The Merlin [9] is a toolkit for neural network-based speech synthesis. Fig.3 is a standard Merlin DNN synthesis architecture. The system takes linguistic features as input and employs neural networks to predict acoustic features, which are then passed to a vocoder to produce the speech waveform. Various neural network architectures are implemented, including a standard feedforward neural network, mixture density neural network, recurrent neural network (RNN), and long short-term memory (LSTM) recurrent neural network, amongst others. The toolkit is Open Source, written in Python, and extensible.

Standard Merlin DNN synthesis architecture is described in detail in.

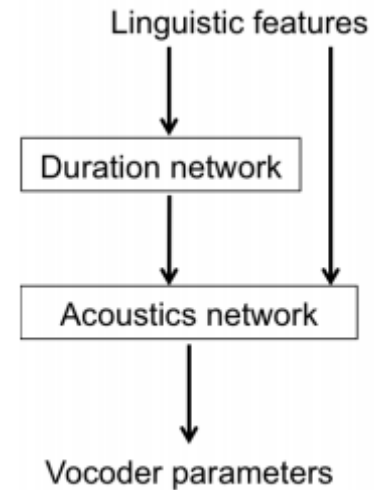


Figure 5. Standard Merlin DNN synthesis architecture [10]

### C. Training TTS

For testing this article, we used Merlin, changed the dictionary and Question List files, and tested and evaluated the results. The system uses a 60-dimensional vector of Mel coefficients containing spectral envelope information, a 25-dimension vector of aperiodicities, and a logarithm of F0 [9]. In the training phase, these vectors are used as the acoustic model deep neural network output. In the synthesis phase, these feature vectors (predicted by the DNN acoustic model) are used as input for the vocoder to synthesize speech signals.

## III. EVALUATION

The evaluation aims to evaluate the quality of output synthesized Vietnamese speech.

The output Vietnamese speech quality of the TTS system was evaluated according to two synthetic speech quality assessment standards. The naturalness of speech was assessed using the MOS (Mean Opinion Score) criterion and rated with five levels (bad-1, poor-2, fair-3, good-4, excellent-5). The intelligibility criterion refers to the ability to fully convey content through synthetic speech, measured as a percentage of the content intelligible ranging from 0% (worst) to 100% (best).

All these assessments for two criteria were conducted through perceptual experiments with listeners. The system was tested in a low-noise environment with 26 Vietnamese people, balanced between men and women, between the ages of 18 and 70, with no hearing or vision impairments or diseases. All test participants do not participate in the training data-building process. The entire testing process will be guided and supervised by technical staff. During the test, each participant will take turns testing ten pre-designed questionnaires. Each questionnaire comprises five Vietnamese sentences selected randomly from an original set of 200 sentences in 10 different fields: culture, society, international, health, law, sport, agriculture, economy, education, tourism, and politics. These sentences were new and did not exist in the training data. Sentences were distributed among listeners. Each sentence in the original will get the same number of evaluations; 7 different people will hear each sentence.

Participants can listen to the voice results once or again if needed. Then participants will rate the two criteria according to their subjective feelings. The final criteria score for the system was defined as the average value of the evaluation results for all sentences, all hearings, and all participants. The results of the evaluation process are summarized in 54.

<sup>11</sup> <http://www.hieuthi.com/blog/2017/03/21/all-vietnamese-syllables.html>

<sup>12</sup> [https://github.com/phamvandong/dictionary\\_xsampa\\_for\\_TTS](https://github.com/phamvandong/dictionary_xsampa_for_TTS)

Table V. TTS system evaluation result

<b>Output speech quality</b>	<b>Vietnamese</b>
<i>MOS (0-5)</i>	4,2
<i>Intelligibility (%)</i>	88%

The MOS scores for Vietnamese were set to 4.2. The high scores indicate that the output speech was almost as natural as human speech. The intelligibility scores of 88% for Vietnamese also show that the output speech was easy to understand and listen to. Both criteria show that Vietnamese's speech's output is of good quality.

The evaluation results show that the Vietnamese TTS system can achieve high results in synthesized speech quality.

#### IV. CONCLUSION

This paper has summarized and proposed Vietnamese speech processing materials, including a list of more than 18.000 Vietnamese words; VI-X-SAMPA transliteration to encode Vietnamese TTS. We have created a Vietnamese question list with an extensive study of phonetics, creating an 808-dimensional matrix in TTS using the Merlin tool.

We also tested the above materials by using Merlin to build a Vietnamese speech synthesis system. The results are evaluated with a MOS score of 4.2 and an accuracy of 88%, which shows that it is highly satisfactory when applying these materials to problems of Vietnamese speech processing.

#### V. REFERENCES

- [1] P. Taylor, "Text-To-Speech Synthesis," *Camb. Univ. Press*, 2009.
- [2] A.-G. Haudricourt, "La place du vietnamien dans les langues austroasiatiques," *Bull. Société Linguist. Paris*, vol. 49, no. 1, pp. 122–128, 1953.
- [3] "Phương Ngữ Học Tiếng Việt (NXB Đại Học Quốc Gia 2009) - Hoàng Thị Châu - 287 Trang | PDF," *Scribd*. <https://www.scribd.com/document/534836166/Ph%C6%B0%C6%BAng-Ng%E1%BB%AF-H%E1%BB%8Dc-Ti%E1%BA%BFng-Vi%E1%BB%87t-NXB-%C4%90%E1%BA%A1i-H%E1%BB%8Dc-Qu%E1%BB%91c-Gia-2009-Hoang-Th%E1%BB%8B-Chau-287-Trang> (accessed Dec. 14, 2022).
- [4] Q. C. Nguyen, "Reconnaissance de la parole en langue Vietnamienne," PhD Thesis, Grenoble INPG, 2002.
- [5] J. C. Wells, "Computer-coding the IPA: a proposed extension of SAMPA," *Revis. Draft*, vol. 4, no. 28, p. 1995, 1995.
- [6] N. T. T. Trang, C. D'Alessandro, A. Rilliard, and T. Do Dat, "HMM-based TTS for Hanoi Vietnamese: issues in design and evaluation," in *14th Annual Conference of the International Speech Communication Association (Interspeech 2013)*, 2013, pp. 2311–2315.
- [7] J. Kirby, "Kirby, James. vPhon: a Vietnamese phonetizer." Nov. 15, 2016. Accessed: Nov. 21, 2019. [Online]. Available: <https://github.com/kirbyj/vPhon>
- [8] T. T. T. Nguyen, "HMM-based Vietnamese Text-To-Speech: Prosodic Phrasing Modeling, Corpus Design System Design, and Evaluation," Paris 11, 2015. Accessed: May 27, 2017. [Online]. Available: <http://www.theses.fr/2015PA112201>
- [9] Z. Wu, O. Watts, and S. King, "Merlin: An Open Source Neural Network Speech Synthesis System.," in *SSW*, 2016, pp. 202–207.
- [10] Z. Malisz, H. Berthelsen, J. Beskow, and J. Gustafson, "Controlling Prominence Realisation in Parametric DNN-Based Speech Synthesis.," in *INTERSPEECH*, 2017, pp. 1079–1083.