

## THE DESCRIPTION OF PHD'S CORPUS

Thesis topic: *Glottalization, tonal contrast, and intonation: an experimental study of Kim Thuong dialect of Muong*

### 1. Data collection:

- Speech materials: Minimal sets/ pairs of real words: **12 minimal sets of 5 smooth tones** (8 complete minimal sets + 4 near minimal sets) + **3 checked minimal pairs**. See details in the tables below (section 4).
- Method of collection: each target word is required to speak four times: two times in isolation and two times in a carrier sentence.
- The carrier sentence is a question including 4 words:

/ja<sup>2</sup>    măt<sup>6</sup>                      cǎŋ<sup>3</sup>/

2sg    to\_know            <target\_item>            INTERROGT

“Do you know \_\_\_\_\_?”

Not only the target words but also all 3 frame words of the carrier sentence is annotated and processed with Peakdet. Hence, for each speaker, we have a total maximum of 660 tokens.

#### The total corpus (per speaker):

$$((5 \times 12) + (2 \times 3) + (5 \times 12 \times 4) + (2 \times 3 \times 4)) \times 2 = 660 \text{ (items)}$$

Elements	Meaning
5	5 tones in smooth syllables
12	12 minimal sets (of 5 smooth tones)
2	2 tones in checked syllables
3	3 minimal pairs (of 2 checked tones)
4	4 syllables of the carrier sentence (1 target word + 3 frame words)
2	2 repetitions

Total corpus: 660 items		
Target syllables: 264 items		Frame syllables: 396 items
<b>In isolation:</b> 132 items	<b>In carrier sentence:</b> 132 items	
- 24 tokens each smooth tone (x 5 tones)	- 24 tokens each smooth tone (x 5 tones)	- /ja <sup>2</sup> / : 132 tokens
- 6 tokens each checked tone (x 2 tones)	- 6 tokens each checked tone (x 2 tones)	- /măt <sup>6</sup> / : 132 tokens
		- /cǎŋ <sup>3</sup> / : 132 tokens

### 2. Annotation

**Annotation scheme:** four digits. Structure: AABC, detailed below.

Position	Meaning	Number code
AA	The UID of the target syllable, in two figures. Numbers from 1 to 9 are padded up with a zero. These numbers corresponding to the UID numbers (second column) in table of minimal sets underneath.	01 – 40: words of 8 minimal sets (8*5) 41 – 60: words of 4 near – minimal sets (4*5) 61 – 66: words of 3 checked minimal pairs (3*2)
B	The nature of the annotated rhyme: B = 0: target syllable said in isolation $1 \leq B \leq 4$ : token from the carrier sentence. B indicates the order of the token inside carrier sentence, 3 means target syllable.	0: target syllable (in isolation) 1: 1 <sup>st</sup> syllable of frame - /ja2/ 2: 2 <sup>nd</sup> syllable of frame - /măt6/ 3: 3 <sup>rd</sup> syllable of frame - target syllable (in carried sentence) 4: 4 <sup>th</sup> syllable of frame - /căn3/
C	repetition	1: first time 2: second time

For example:

- Item with UID 6801 means: the check syllable /kap7/ in isolation, first time.
- Item with UID 3632 means: the target syllable /ma5/ in carrier sentence, second time.
- Item with UID 3611 means: the first syllable /ja2/ of the carrier sentence in case /ma5/ is target syllable, first time.

As a consequence, based on two first digits (AA), we can classify the tokens into 7 sorts of tone:

T1	T2	T3	T4	T5	T6	T7
04	03	02	05	01	61	62
09	08	07	10	06	63	64
14	13	12	15	11	65	66
19	18	17	20	16		
24	23	22	25	21		
29	28	27	30	26		
34	33	32	35	31		
39	38	37	40	36		
44	43	42	45	41		
49	48	47	50	46		
54	53	52	55	51		
59	58	57	60	56		

### 3. Data processing using Peakdet: input and output

#### 3.1. Input:

After annotation, two input files are prepared for data processing with Peaket:

- EGG file (.WAV): the EGG signal was recorded simultaneously with acoustic signals, now are extracted from the annotated audio file.
- REGIONS file (.TXT): this is a text file containing the list of regions created at annotation (i.e. the time codes of the target items annotated in SoundForge).

### 3.2. Output:

As final output, Peakdet produces a “.mat” file. The result is placed in a matrix, containing one sheet per token analyzed. It means the maximum data of each speaker will contain 660 sheets corresponding to 660 rhymes. However, this absolute quantity is not met in many cases. Because the carrier sentence are repeated many time, the speakers tend to reduce the first syllable. The specific status of each data will be provided in the table of speakers in section 5.

The structure of a single matrix:

Each sheet presents the main results of a token in a single matrix including 10 columns x 100 line.

Column	Value	Unit
1 <sup>st</sup>	the beginning of cycle	ms
2 <sup>nd</sup>	the end of cycle	ms
3 <sup>th</sup>	f <sub>0</sub>	Hz
4 <sup>th</sup>	DECPA: Derivative-EGG Closure Peak Amplitude	
5 <sup>th</sup>	Oq determined from raw maximum without smoothing	%
6 <sup>th</sup>	DEOPA	
7 <sup>th</sup>	Oq determined from maximum after smoothing	%
8 <sup>th</sup>	Oq determined from peak detection without smoothing	%
9 <sup>th</sup>	Oq determined from peak detection after smoothing	%
10 <sup>th</sup>	The open quotient values retained after user’s verification	%

## 4. Speech materials: tables of minimal sets/pairs

12 minimal sets (8 complete minimal sets + 4 near minimal sets)

N <sub>o</sub>	UID	IPA	PS	Vietnamese	English
1	01	<b>paj5 (t<sup>h</sup>ăj1)</b>	n	Sải tay	Armspan
	02	<b>paj3</b>	n	(Cái) vại	A kind of porcelain jar to salt vegetables or tubes
	03	<b>(t<sup>h</sup>yp6) paj2</b>	n	Đập nước	Barrage
	04	<b>paj1</b>	n	Vải	Cloth
	05	<b>paj4</b>	n	Quả, trái	Fruit
2	06	<b>rɔ5</b>	n	(Con) rùa	Tortoise
	07	<b>rɔ3</b>	v	Mò (cua, cá)	To find crab or fish (by hand) in a swamp
	08	<b>rɔ2</b>	adj	No	To be sated
	09	<b>(ʔɔ5) rɔ1</b>	adj	Rảnh rỗi	Idle
	10	<b>(paj4) rɔ4</b>	n	Hoa chuối	Banana flower
3	11	<b>pa5</b>	n	Bà	Grand-mother/ an old woman
	12	<b>pa3</b>	v	Bám, trạm vào vai	To cling on one’s shoulder

	13	<b>pa2</b>	n	Ba	Three
	14	<b>pa1 (tʰiən5)</b>	v	Trả (tiền)	To pay
	15	<b>pa4</b>	v	Vá (xăm)	To patch
4	16	<b>laj5</b>	n	Lưỡi	Tongue
	17	<b>(pɤ2) laj3</b>	v	Trở lại,	Return
	18	<b>laj2</b>	v	Lai	To take someone up by bike/ motobike?
	19	<b>laj1 (tʰm5)</b>	n	(Cái) rào/ chắn ao	A fence to prevent fish go out from lake
	20	<b>laj4</b>	v	Lái	To drive
5	21	<b>taj5</b>	v	Đãi gạo	To wash rice
	22	<b>taj3</b>	v	Lôi (núa, bương bằng tay hoặc xe máy)	To drag bamboo by hand, or by motorbike
	23	<b>taj2 (nan3)</b>	n	Tai nạn	Accident
	24	<b>taj1</b>	n	Thác (nước)	Waterfall
	25	<b>taj4</b>	n	Đái	To pee
6	26	<b>kɔ5</b>	n	(Con) cò	Stork
	27	<b>kɔ3</b>	v	Nói	To speak
	28	<b>(kiɛw4) kɔ2</b>	v	Kéo co	Tug of war
	29	<b>kɔ1</b>	n	Cỏ	Grass
	30	<b>kɔ4</b>	v	Có	To have
7	31	<b>kiɛŋ5</b>	n	Cái chum nhỏ	A jar to
	32	<b>kiɛŋ3</b>	pre	(Ổ) cạnh	Next to
	33	<b>kiɛŋ2</b>	n	Canh	Soup
	34	<b>kiɛŋ1</b>	n	(Cái) keng	A kind of gong to gather people
	35	<b>kiɛŋ4</b>	n	Cánh	Wing
8	36	<b>ma5</b>	n	Lỗ cua, lươn	Cave hole of eel or crab
	37	<b>ma3</b>	n	Mạ	Rice seedling
	38	<b>ma2</b>	n	(Con) ma	Ghost
	39	<b>ma1</b>	n	(Mồ) mả	Tomb
	40	<b>ma4</b>	n	Má	Cheek
9	41	<b>ŋa5</b>	v	Ngã	to fall
	42	<b>ŋa3</b>	v	Ngứa	To itch
	43	<b>ŋa2</b>	n	Chói mắt	Dazzle
	44	<b>ŋa1</b>	v	(Nằm) ngả	To recline
	45	<b>na4</b>	n	Cung tên	Archery
10	46	<b>ka5</b>	n	Cà (tím)	African eggplant
	47	<b>ta3</b>	n	Tạ	Dumbbells
	48	<b>ka2</b>	n	(Con) gà	Chicken
	49	<b>ka1</b>	adj	To	Big
	50	<b>ka4</b>	n	Cá	Fish
11	51	<b>kaj5</b>	v	Cài ( cúc áo)	To fasten buttons of clothing
	52	<b>paj3</b>	n	(Cái) vại	A kind of porcelain jar to salt vegetables or tubes
	53	<b>kaj2</b>	n	(Cái) gai	Thorn
	54	<b>kaj1</b>	n	(Rau) cải	Borecole
	55	<b>kaj4</b>	adj	(Con) cái	Female
12	56	<b>ku5</b>	adj	Cũ	Old, ancient
	57	<b>tu3</b>	v	Tụ (máu)	To converge (blood)

	58	<b>ku2</b>	n	Trâu	Buffalo
	59	<b>ku1</b>	n	Củ (khoai)	Tubes
	60	<b>ku4</b>	n	Cú mèo	Owl

### 3 Checked minimal pairs

N <sub>o</sub>	UID	IPA	PS	Vietnamese	English
1	61	<b>pat6</b>	n	Sàn (nhà)	Floor made by a kind of bamboo
	62	<b>pat7</b>	n	Bát	Bowl
2	63	<b>rɔc6</b>	n	Ruột	Intestine
	64	<b>rɔc7</b>	v	Rót (nước)	To pour (water)
3	65	<b>lak6</b>	n	(Củ) lạc	Peanut
	66	<b>lak7</b>	adj	Lác	Squint eye

	Phonemes	Appearance	Minimal set(s)
<b>Smooth syllables</b>			
Initial consonants	p	2	1, 3
	t	1	5
	k	5	6, 7, 10, 11, 12
	m	1	8
	ŋ	1	9
	l	1	4
	r	1	2
Vowels	a	8	1, 3, 4, 5, 8, 9, 10, 11
	ɔ	2	2, 6
	u	1	12
	ie	1	7
Final consonants			
	j	4	1, 4, 5, 11
	ŋ	1	7
<b>Checked syllables</b>			
Initial consonants	p	1	1
	l	1	3
	r	1	2
Vowels	a	2	1, 3
	ɔ	1	2
Final consonants	t	1	1
	c	1	2
	k	1	3

### 5. Information of speakers and their data

- M = male, F = female
- Number of data that can be analyzed: 10 man, 11 woman
- Number of data has been processed: 10 man, 10 woman

- Total: 26 participant, 28 files, 20 files have been process

ID	Birth year	Rec. time	EGG quality	Status of data	Size of .mat file
F1	1983	2018	Bad Crackling noise	No analysis	No analysis
<b>F3</b>	1984	2018	Good	First time: 421/460 tokens - Missing 35 first syllables (position AA1C): vowels reduced - Missing 4 tokens of a carrier sentence. One target syllable (in carrier sentence) is problematic. Second time: 660/660 tokens	100x10x 660
<b>F7</b>	1988	2018	OK	660/660 tokens	119x10x660
F8	1973	2018	Weak EGG	No analysis	No analysis
<b>F9</b>	1991	2018	Good	660/660 tokens	119x10x660
<b>F10</b>	1973	2018	Good	585/660 tokens (Missing 75 tokens) Missing minimal set 11 in carrier sentence, minimal set 12 and 3 minimal pairs both in isolation and carrier sentence (in first repetition) - 11 AA01 - 16 AA31 - 16 AA11, 16 AA21, 16 AA41	178x10x660
F11	1982	2018	Very weak EGG	No analysis	No analysis
<b>F12</b>	1972	2018	Good	660/660 tokens	111x10x660
<b>F13</b>	1967	2018	Good	646/660 tokens Missing 14 tokens (AA1C)	133x10x646
F14	1965	2018	OK But there are a few flat segments	No analysis	No analysis
F16	1956	2019	Weak EGG	No analysis	No analysis
<b>F17</b>	1955	2019	Good	660/660 tokens	100x10x660
F18	1988	2019	Weak EGG	No analysis	No analysis

<b>F19</b>	1975	2019	Good	660/660 tokens	100x10x660
<b>F20</b>	1989	2019	OK	660/660 tokens	100x10x660
<b>F21</b>	1988	2019	OK	660/660 tokens	111x10x660
<b>M1</b>	1979	2018	Good	660/660 tokens	100x10x660
<b>M5</b>	1955	2018	OK	660/660 tokens	100x10x660
<b>M7</b>	1984	2018	Good	660/660 tokens	100x10x660
<b>M8</b>	1975	2018	OK	660/660 tokens	100x10x660
<b>M9</b>	1990	2018	Good	660/660 tokens	100x10x660
<b>M10</b>	1988	2018	Good	660/660 tokens	100x10x660
<b>M11</b>	1987	2018	Good	660/660 tokens	100x10x660
<b>M12</b>		2018	Signal out of range	No analysis	No analysis
		2019	Good	660/660 tokens	100x10x660
<b>M13</b>	1962	2018	Good	660/660 tokens	100x10x660
<b>M14</b>	1959	2018	OK	656/660 tokens Missing 4 tokens: 3 AA1C, 1 AA4C	100x10x656