



Rich Prosodic Information Exploration on Spontaneous Mandarin Speech

Cheng-Hsien Lin¹, Chung-Long You¹, Chen-Yu Chiang², Yih-Ru Wang¹, Sin-Horng Chen¹

¹Dept. of Electrical and Computer Engineering, National Chiao Tung University, Taiwan

²Dept. of Communication Engineering, National Taipei University, Taiwan

Summary

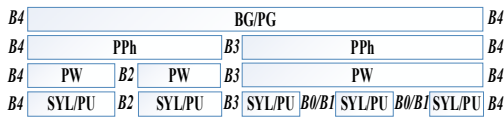
The joint prosody labeling and modeling algorithm for Mandarin read speech is extended for spontaneous speech by additionally considering the affecting patterns:

- Particular Unit (PU) type: particle, marker, uncertain and foreign speech
- Part of speech for normal speech part
- Break types of syllable boundaries
- Syllable contraction and lengthening

- A hierarchical prosodic model (HPM) is constructed for mandarin spontaneous speech using MCDC corpus
- The prosodic characteristics and disfluency events of spontaneous Mandarin speech are explored by investigating the parameters of HPM

The Hierarchical Prosody Modeling

- The four-layer prosodic structure consisting syllable/particular unit (SYL/PU), prosodic word (PW), prosodic phrase (PPh), and breath group/prosodic phrase group (BG/PG) is adopted
- Two types of tags are employed to represent the prosody structure: **BREAK type** (syllable/PU juncture) and **PROSODY STATE** (inside syllable/PU)



The HPM describes the relationships among prosody tags (T), prosody-acoustic (A) and linguistic features (L), and can be expressed by 8 sub-models:

$$P(T|A,L,A) = \prod_{n=1}^N \left(P(sp_n | B_{n-1}, p_n, f_n^{n+1}, pos_n) P(sd_n | B_{n-1}, q_n, t_n, s_n, pos_n, cl_n) \right) \quad 1$$

$$P(P_i | P_j) = \prod_{n=1}^N \left(P(p_n | p_{n-1}, B_{n-1}) P(q_n | q_{n-1}, B_{n-1}) P(r_n | r_{n-1}, B_{n-1}) \right) \quad 2$$

$$P(P_i | P_j) = \prod_{n=1}^N \left(P(pd_n | ed_n, pj_n, dl_n | B_{n-1}) P(B_n | I_n) \right) \quad 3$$

$$P(P_i | P_j) = \prod_{n=1}^N \left(P(se_n | B_{n-1}, r_n, t_n, f_n, pos_n) \right) \quad 4$$

$$P(P_i | P_j) = \prod_{n=1}^N \left(P(p_n | p_{n-1}, B_{n-1}) P(q_n | q_{n-1}, B_{n-1}) P(r_n | r_{n-1}, B_{n-1}) \right) \quad 5$$

$$P(P_i | P_j) = \prod_{n=1}^N \left(P(pd_n | ed_n, pj_n, dl_n | B_{n-1}) P(B_n | I_n) \right) \quad 6$$

$$P(P_i | P_j) = \prod_{n=1}^N \left(P(p_n | p_{n-1}, B_{n-1}) P(q_n | q_{n-1}, B_{n-1}) P(r_n | r_{n-1}, B_{n-1}) \right) \quad 7$$

$$P(P_i | P_j) = \prod_{n=1}^N \left(P(pd_n | ed_n, pj_n, dl_n | B_{n-1}) P(B_n | I_n) \right) \quad 8$$

Prosody tags include Breaks (B) and prosody states of pitch contour (p), duration (q), and energy (r) for syllable/PU

$$T = \{B_n, p_n, q_n, r_n | n = 1 \dots N\}$$

- (1,2,3,7) Prosody-acoustic feature related sub-models, representing within syllable/PU features: pitch (sp), duration(sd), energy(se), and between syllable/PU features: pause duration(pd), energy dip(ed), pitch jump(pj) and duration lengthening(dl)

$$A = \{sp_n, sd_n, se_n, pd_n, ed_n, pj_n, dl_n | n = 1 \dots N\}$$

- (4,5,6) Prosody-state sub-models of pitch, duration and energy
- (8) Break-syntax sub-model

The linguistic feature set (L) includes: reduced linguistic feature(l), tone(t), base-syllable(s), final(f), part-of-speech(pos), and contraction/lengthening tag(cl)

$$L = \{l_n, t_n, s_n, f_n, pos_n, cl_n | n = 1 \dots N\}$$

sp, sd, and se are further expressed:

$$P(sp_n | B_{n-1}, p_n, f_n^{n+1}, pos_n) = \begin{cases} N(sp_n; \beta_{sp} + \beta_{cl}^f + \beta_{cl}^s + \beta_{cl}^p + \beta_{cl}^t + \beta_{cl}^s + \beta_{cl}^p + \beta_{cl}^t, R_{sp}^s), & \text{for SYL} \\ N(sp_n; \beta_{sp} + \beta_{cl}^f + \beta_{cl}^s + \beta_{cl}^p + \beta_{cl}^t + \beta_{cl}^s + \beta_{cl}^p + \beta_{cl}^t, R_{sp}^p), & \text{for PU} \end{cases}$$

$$P(sd_n | B_{n-1}, q_n, t_n, s_n, pos_n, cl_n) = \begin{cases} N(sd_n; \gamma_{cl} + \gamma_{cl}^f + \gamma_{cl}^s + \gamma_{cl}^p + \gamma_{cl}^t + \gamma_{cl}^s + \gamma_{cl}^p + \gamma_{cl}^t + \mu_{sd}, R_{sd}^s), & \text{for SYL} \\ N(sd_n; \gamma_{cl} + \gamma_{cl}^f + \gamma_{cl}^s + \gamma_{cl}^p + \gamma_{cl}^t + \gamma_{cl}^s + \gamma_{cl}^p + \gamma_{cl}^t + \mu_{sd}, R_{sd}^p), & \text{for PU} \end{cases}$$

$$P(se_n | B_{n-1}, r_n, t_n, f_n, pos_n) = \begin{cases} N(se_n; \alpha_{cl} + \alpha_{cl}^f + \alpha_{cl}^s + \alpha_{cl}^p + \alpha_{cl}^t + \alpha_{cl}^s + \alpha_{cl}^p + \alpha_{cl}^t + \mu_{se}, R_{se}^s), & \text{for SYL} \\ N(se_n; \alpha_{cl} + \alpha_{cl}^f + \alpha_{cl}^s + \alpha_{cl}^p + \alpha_{cl}^t + \alpha_{cl}^s + \alpha_{cl}^p + \alpha_{cl}^t + \mu_{se}, R_{se}^p), & \text{for PU} \end{cases}$$

$\beta_x, \gamma_x, \alpha_x$ are the affecting patterns (APs) of affecting factors x for syllable/PU pitch contour, duration, and energy models

$\mu_x, \Sigma_x, R_x, C_x$ denote the global mean vector and covariance matrix of modeling residual

- The unsupervised PLM algorithm is employed to simultaneously train the HPM and label the corpus with prosodic tags.

$$T^*, \Lambda^* = \arg \max_{T, \Lambda} P(T|A, L, \Lambda)$$

Database and Preprocess

Database

- MCDC 8-hour dialogues from 16 speakers
- Texts with phonetic and linguistic tags are transcribed by linguist experts

Spontaneous speech characteristics: repetition/restart/repair, pause, pronunciation variation and sociolinguistic phenomena are also annotated

Syllable alignment

HMM forced alignment with manual error correction

Syllable acoustic features

Utterance level mean-and-variance normalization for duration and energy

Speaker level mean-and-variance normalization for F0 values

Syllable contour is represented by four Legendre Polynomial coefficients

Modeling results of HPM

- The total residual errors (TREs) of modeling for normal syllable

Pitch		Duration		Energy	
AF	TRE	AF	TRE	AF	TRE
-	-	+Contra. & Lengthen (cl.)	81.2%	-	-
+Tone (t _n)	90.7%	+Tone	78.9%	+Tone	94.4%
+Coarticulation (cp _n)	84.1%	+Syllable Type (s _n)	68.6%	+Final (f _n)	85.3%
+POS (pos _n)	77.2%	+POS	65.2%	+POS	82.4%
+Break (b _n)	71.7%	+Break	56.0%	+Break	80.1%
+Prosodic State (ps _n)	8.0%	+Prosodic State (qs _n)	3.6%	+Prosodic State (es _n)	2.3%

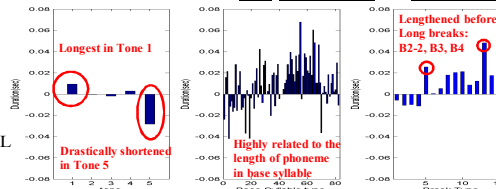
- TREs of modeling for PU

Pitch		Duration		Energy	
AF	TRE	AF	TRE	AF	TRE
-	-	+Contra. & Lengthen	84.4%	-	-
+PU ID (pu _n)	83.3%	+PU ID	72.2%	+PU ID	86.5%
+Break	76.3%	+Break	64.6%	+Break	81.4%
+Prosodic State	9.5%	+Prosodic State	2.9%	+Prosodic State	44.2%

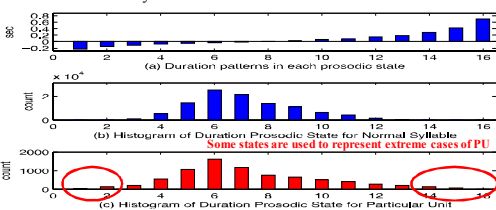
- SYL/PU duration APs of contraction, lengthening and normal speech

	Contraction	Lengthening	Normal
SYL duration (ms)	-28	156	9
PU duration (ms)	-44	143	-19

- SYL duration APs of tone, base-syllable, and adjacent breaks



- SYL duration APs of prosody states and the histograms of normal syllable and PU



Analysis on Disfluency Events

- The utterance examples of Repetition

Type	Start	Interrupt	End	Sentence
1	對對對	我有聽說	所以像我先講	以前有一種說法 *(dui)*(dui)*(dui)wo yo ting shuo suo yi xiang wo xian sheng
2	因為	(他)舉了一個例子		yin wei*(ta)*(ta)#ju le yi ge li zi
3	待一整天	(類似)(類似)#亞歷山大那種		dai yi zheng tian dui*(lei si)*(lei si)#ya li shan da na zhong
4	吹涼風	就覺得(很舒服)(很舒服)所以這個		chui liang feng jiu jue de*(hen shu fu)*(hen shu fu)#suo yi zhe ge

- Break labeling results of repetition at different boundaries

Type	Dominated by Non-pause breaks				Interrupt boundary			
	B0	B1	B2-1	B2-2	B2-2	B2-3	B3	B4
1	28%	43%	9%	3%	15%			<1%
2	7%	18%	25%	16%	9%	20%		4%
3&4	4%	11%	27%	15%	15%	19%		8%

Type	End boundary Dominated by Short-to-long pause breaks				Utterance beginning			
	B0	B1	B2-1	B2-2	B2-3	B3	B4	Bs
1	5%	5%	41%	8%	3%	29%		9%
2	19%	23%	37%	5%	11%	3%		1%
3&4	17%	29%	27%	6%	12%	6%		3%

Type	Short stops for all three types				Start boundary				Utterance beginning							
	B0	B1	B2-1	B2-2	B2-3	B3	B4	Bs	B0	B1	B2-1	B2-2	B2-3	B3	B4	Bs
1	2%	7%	14%	6%	8%	11%	6%		46%							
2	8%	25%	22%	9%	3%	14%	5%		17%							
3&4	9%	18%	23%	5%	5%	14%	6%		19%							

- The utterance examples of Restart and Repair

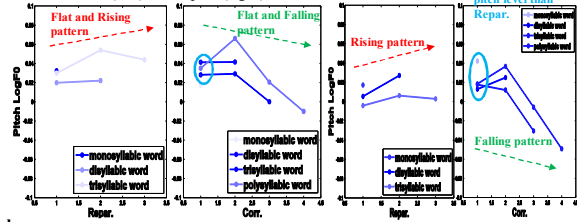
Restart	我我我我我父親#也沒有
Repair	四十分鐘左右#(就回到)#台北
	(wo wo wo wo wo fu qin)#ye mei you
	si shi fen zhong zuo you#(jiu dao)#tai bei

- Break labeling results of restart and repair at different boundaries

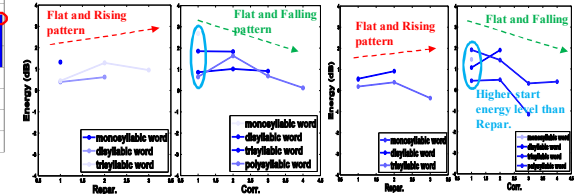
Boundary	Restart						
	B0	B1	B2-1	B2-2	B2-3	B3	B4
Start	10%	29%	24%	8%	10%	14%	5%
Interrupt	2%	9%	34%	25%	9%	16%	5%
End	15%	37%	20%	4%	15%	7%	2%

Boundary	Repair						
	B0	B1	B2-1	B2-2	B2-3	B3	B4
Start	13%	30%	24%	7%	13%	10%	3%
Interrupt	3%	11%	29%	27%	10%	18%	2%
End	13%	36%	16%	5%	18%	10%	1%

- The averaged pitch prosody state APs of reparandum and correction for restart (left) and repair (right)



- The averaged energy prosody state APs of reparandum and correction for restart and repair



- The averaged duration prosody state APs of reparandum and correction for restart and repair

