

Highlighting Utterances in Chinese Spoken Discourse

Shu-Chuan Tseng
Institute of Linguistics
Academia Sinica, Nankang 115
Taipei, Taiwan
tsengsc@gate.sinica.edu.tw

Abstract

This paper presents results of an empirical analysis on the structuring of spoken discourse focusing upon how some particular utterance components in Chinese spoken dialogues are highlighted. A restricted number of words frequently and regularly found within utterances structure the dialogues by marking certain significant locations. Furthermore, a variety of signals of monitoring and repairing in conversation are also analysed and discussed. This includes discourse particles, speech disfluency as well as their prosodic representation. In this paper, they are considered a kind of “highlighting-means” in spoken language, because their function is to strengthen the structuring of discourse as well as to emphasise important functions and positions within utterances in order to support the coordination and the communication between interlocutors in conversation.

1 Introduction

In the field of discourse analysis, many researchers such as Sacks, Schegloff and Jefferson (Sacks et al. 1974; Schegloff et al. 1977), Taylor and Cameron (Taylor & Cameron 1987), Hovy and Scott (Hovy & Scott 1991) and the others have extensively explored how conversation is organised. Various topics have been investigated, ranging from turn taking, self-corrections, functionality of discourse components, intention and information delivery to coherence and segmentation of spoken utterances. The methodology has also accordingly changed from commenting on fragments of conversations, to theoretical considerations based on more materials, then to statistical and computational modelling of conversation. Apparently, we have experienced the emerging interests and importance on the structure of spoken discourse.

Recently, Clark and Brennan (Clark & Brennan 1991) proposed that production-related as well as perception-related activities have to be coordinated in conversation. They should serve the purpose of a speaker to get an addressee's attention, to plan and produce utterances, to recognise when the addressee does not understand, to initiate and manage repairs and to display or acknowledge understanding. Since interlocutors in conversation usually communicate freely and spontaneously, utterances may be interrupted without being completed and turn taking may take place unpredictably. Besides, there can be erroneous or unknown words and phrases, when speakers are not able to express their thoughts properly or when they spontaneously create new words or compound words for their thoughts. While monitoring and correcting their speech, speakers make repairs and they may need time for re-planning and editing their speech. Discourse particles and pauses (silent or filled) are often used for this purpose. Especially discourse particles are usually located in related positions in a given

discourse. They may possibly signal turn-initial positions or they may express special pragmatic functions such as surprise or hesitation.

With regard to the internal structure of utterances, speech disfluency no doubt results in serious problems for natural language processing systems. But, not for humans! Why? One of the reasons could be that speakers of a common language share similar knowledge and competence on how to express emotion pragmatically and how to perform monitoring and repairing in a reasonable way, whenever there is a need. Therefore, the communication partners can easily decode the meaning and the function of speech disfluency in conversation, because they would encode it by using similar sequence combinations. Moreover, speech disfluency also contributes to the segmentation of discourse by pointing out problematic speech sequences.

In the production of speech, all of the phenomena mentioned above are at the same time accompanied by prosody. No doubt, prosody is the most powerful highlighting means in spoken language. A grammatically ill-formed sentence may just sound like a perfectly correct sentence, when it is produced with a professional and fluent intonation. However, it is not the intention of interlocutors in conversation to impress each other in this way. In contrast, the speaker would rather let the addressee know that he/she made a mistake previously and what he/she is now saying is correct. Thus, it is very likely that discourse particles and speech disfluency are highlighted by prosodic means to emphasise the semantic and syntactic inadequacy.

This paper mainly deals with means emphasising particular speech sequences and four possibilities of observing these activities are 1) frequently used words, 2) discourse particles and markers, 3) speech disfluency and 4) prosodic marking. In the coming sections, they will be investigated on the basis of results of a corpus analysis. Before the results are presented, speech data used for the analyses are introduced first.

2 Chinese Spoken Dialogues

Three dialogues from Taiwan Putonghua Corpus (TWPTH) have been analysed. Putonghua refers to Mandarin. All subjects were born in Taiwan and their first language is Taiwanese. In other words, this set of data is on Taiwanese accented Mandarin. The subjects were given the instruction to talk on any topic they wanted to. The data obtained are therefore to a great extent spontaneous and natural. Each dialogue is about 20 minutes long. Major topics in the dialogues are family, work and study, although the communication partners in dialogues did not stick to a specific topic all the time. Detailed information about the dialogues is summarised in Table 1 including the gender of the subjects and the number of turns, words and characters found in the dialogues. We used the transcripts of the speech data with syntactic annotations specifically defined for this analysis.

Table 1: Statistics of the Speech Data

DIALOGUES	SUBJECTS: SEX	URNS	WORDS	CHARAC- TERS	TOTAL WORDS	TOTAL CHARAC- TERS
Dialogue-01	D1-A: F	183	1688	7410	3625	15518
	D1-B: F	189	1937	8108		
Dialogue-02	D2-A: F	163	1247	5500	2622	11556
	D2-B: M	168	1375	6056		
Dialogue-03	D3-A: M	144	1782	7936	2878	12896
	D3-B: M	142	1096	4960		

In this study, we segment the conversations into turns, instead of utterances. The reason is that utterance boundaries in spoken discourse are difficult to determine. In most of the cases, the initial and the final locations of utterances are not clearly indicated as those of sentences in written texts. Gross et al. and Traum and Heeman have recently suggested various principles of determining utterance boundaries (Gross et al. 1993, Traum & Heeman 1997). Nevertheless, they need further elaborations. Utterances are ambiguous. But turns are relatively simple to recognise. Besides, segmenting the dialogues into turns completely serve the purpose of this study.

Next, the dialogues have been tagged using the tagging system designed for the Academia Sinica Balanced Corpus (hereafter Sinica Corpus) (tag set cf. Chen et al. 1996, CKIP 1995). It should be noted that the majority of the Sinica Corpus data are written texts and the tags have been developed restrictively for well-formed written sentences. Hence we added extra-tags for ill-formed sequences and spontaneous speech phenomena. They are for instance pauses (pause), unidentifiable sounds (NSS), speech errors (POS-se), partial or full speech repetitions (POS-pr and POS-fr, respectively) and discourse particles (DP), where POS stands for the part of speech tags defined for the Sinica Corpus. For instance, the utterance segment 所以久而久而久之小孩子的模式就會跟我們一樣¹ has been tagged as in Example 1.

Example 1: 所以(Cbb) 久而(D-pr) 久而久之(D) 小孩子(Na) 的(DE) 模式(Na)
 suǒ yǐ jiǔ ér jiǔ ér jiǔ zhī xiǎo hái zi de mó shì
 就(D) 會(D) 跟(P) 我們(Nh) 一樣(VH)
 jiù huì gēn wǒ mén yí yàng (dialogue-01)

3 Discourse Particles: Type and Position

Discourse markers usually include particles such as "oh", "well", "now" and "then" and connectives such as "so", "because", "and", "but" and "or" in English (Schiffirin 1987). These can also be called cue phrases, cue words, or discourse particles in general (Hirschberg & Litman 1993). Since there is no consistent definition for distinguishing discourse particles and discourse markers in the literature, a primitive distinction is proposed in this paper. Discourse particles express particular pragmatic indications such as doubt or anger. They verbalise emotion, give listeners time to process problematic utterances, but they do not have lexicalised meaning in lexicon. As typically accepted in the literature, discourse particles are only used in a pragmatic way without having any syntactic or lexical participation in the discourse. On the other hand, discourse markers can be lexical items listed in a lexicon such as verbs, nouns and connectives. But when used in spoken language, they may have similar functions as discourse particles. They mark special locations in utterances and they stress particular speaker intentions in discourse as discourse particles usually do.

Table 2: Distribution of Discourse Particles and Particle-Like Words in Dialogues

	URNS	DISCOURSE PARTICLES	NÀ- AND ZHÈ- WORDS	AVERAGE PER TURN
Dialogue-01	372	23 types, 479 tokens	7 types, 112 tokens	1.6 particle
Dialogue-02	331	18 types, 477 tokens	7 types, 113 tokens	1.6 particle
Dialogue-03	286	19 types, 295 tokens	6 types, 169 tokens	1.6 particle

¹ Therefore, long time, after a long time the behaviour patterns of the children will be the same as ours.

The goal of this part of study is to determine how often and which kinds of discourse particles are preferably used in our spoken dialogues. Two groups of words are under consideration. They are the typical discourse particles and determiners/adverbials having similar pragmatic function as discourse particles. As illustrated in Table 2, 23 typical and different discourse particles were found in the dialogues. On average, every turn contains at least one discourse particle. More than ten discourse particles were produced in all three dialogues with high frequency. This result clearly supports the notion that speakers of a common language tend to use similar particles for similar purposes and they all use them very frequently. Across all three dialogues, given a specific discourse particle, the sentential position where it is located is usually the same. In other words, there is a regular mapping between the types of discourse particles and the positions where they are to be found. The following examples should illustrate the phenomenon. 嘛 (mǎ) is normally utterance-final, whereas 嗯 (en) is more likely to be found in the utterance-initial position.

Example 2: 那 我 在 想 這個 到底 因為 他 八歲 讀書 嘛²
 nà wǒ zài xiǎng zhè ge dào dī yīn wèi tā bā suì dú shū mǎ
 (dialogue-01)
 而且 是 比較 新 新鮮 一點 你 在 鄉下 嘛³
 ér qiě shì bǐ jiào xīn xīnxiān yì diǎn nǐ zài xiāngxià mǎ
 (dialogue-02)
 有 這個 意思 是 啊 你 可以 自己 學 嘛⁴
 yǒu zhè ge yì si shì a nǐ kě yǐ zì jǐ xué mǎ
 (dialogue-03)

Example 3: 嗯嗯 小涵 就 噉 好像 沒辦法 控制
 en en xiǎo hán jiù ou hǎo xiàng méi bàn fǎ kòng zhì
 自己 的 噉 行為 噉 那 一些 的⁵
 zì jǐ de ou xíng wéi ou nà yì xiē de (dialogue-01)
 嗯 對 噉 你 你 像 我 來到 這邊 那
 en en duì en nǐ nǐ xiàng wǒ lái dào zhè biān nà
 呵 除了 工作 就是 工作 啦⁶
 he chú le gōng zuò jiù shì gōng zuò la (dialogue-02)
 噉 那 博士 博士 那 進去 等級 就 比較 高 噉
 en nà bó shì bó shì nà jìn qù děng jí jiù bǐ jiào gāo mǎ
 還是 跟 一般人 一樣⁷
 hái shì gēn yì bān rén yí yàng (dialogue-03)

Determiners/adverbials 那 (nà), 那個 (nà ge), 那麼 (nà me), 那樣 (nà yàng), 這 (zhè), 這個 (zhè ge), 這麼 (zhè me) and 這樣 (zhè yàng)⁸, serving similar pragmatic purposes as discourse particles, have been investigated, too. Used in conversation, these words function more like connectives such as *so*, *therefore*, or *then*. This kind of pragmatic implication is characteristic of

² Then I'm thinking this on earth because he went to school at the age of eight MA.

³ In addition, it is a little bit fresh, fresher you are in the countryside MA.

⁴ Having this intention yeah A you can learn it yourself MA.

⁵ EN EN Xiao-Han seems OU that she could not control her OU behaviour herself that kind of.

⁶ EN EN yeah EN you you like me came here NA HE after work is still work LA.

⁷ EN NA PhD PhD NA enter the level then higher MA or it is the same as the ordinary people.

⁸ The original semantic meaning of 那 (nà), 那個 (nà ge), 那麼 (nà me), 那樣 (nà yàng), 這 (zhè), 這個 (zhè ge), 這麼 (zhè me) and 這樣 (zhè yàng) are "that", "that one", "that way, then", "that way", "this", "this one", "this way, so" and "this way, so".

discourse markers. Thus, counting discourse particles and 那 (nà) and 這 (zhè) words together, we found 1.6 particle words in each turn on average, as illustrated in Examples 4 and 5.

Similar to the typical discourse particles, 那 (nà) and 這 (zhè) words are more likely to be found in certain syntactic positions than in other positions. For instance, 那 (nà) as a connective appears mostly in the utterance-initial position, while 那個 (nà ge) seems to retain more characteristics of determiners and is often located in the mid-utterance position. As shown in Table 2, 那 (nà) and 這 (zhè) words were very often used as particle-like discourse markers in the dialogues. This result demonstrates the fact that words can be used differently in written and spoken language.

Example 4: 那 邱玉芬 他們 家 謝玉姣 雖然 二 三樓
 nà qiū yù fēn tā mén jiā xiè yù jiāo suī rán èr sān lóu
 不過 她 還要 拿 拿出 一百萬 來 那個 補 一
 bú guò tā hái yào ná ná chū yì bǎi wàn lái nà ge bǔ yì
 貼 一百萬 貼 出來
 tiē yì bǎi wàn tiē chū lái⁹ (dialogue-01)
 那 我 還有 那 小孩子 呵 他們 差不多
 nà wǒ hái yǒu nà xiǎo hái zi he tā mén chà bù duō
 兩三歲 的 時候 我 就 一直 帶出 我 就
 liǎng sān suì de shí hòu wǒ jiù yì zhí dài chū wǒ jiù
 好像說 給 他 出去 看看 因為 這邊 的
 hǎo xiàng shuō gěi tā chū qù kàn kàn yīn wèi zhè biān de
 話 真的 甚麼 都 沒有 真的 要 玩 的 要 呵
 huà zhēn de shén me dōu méi yǒu zhēn de yào wán de yào he
 要 穿的 啊 這樣 我 老大 呵 哎
 yào chuān de a zhè yàng wǒ lǎo dà he ai
 這邊 來 的 剛好 待 這邊 嗒 不行 那
 zhè biān lái de gāng hǎo dāi zhè biān nuò bù xíng nà
 甚麼 東西 都 還要 到 台中 去 買¹⁰
 shé me dōng xi dōu hái yào dào tái zhōng qù mǎi (dialogue-02)
 那 你 那 論文 作 的 怎樣 差不多了吧¹¹
 nà nǐ nà lùn wén zuò de zě yàng chà bù duō le ba (dialogue-03)

Example 5: 嘿 對 像 他們 我 家 隔壁 那個 棟 噉
 hei duì xiàng tā mén wǒ jiā gé bì nà ge dòng ou
 可能 是 光是 老大 和 老二 的 三樓 已經 賣掉 了¹²
 kě néng shì guāng shì lǎo dà hé lǎo èr de sān lóu yǐ jīng mai diào le
 (dialogue-01)
 沒有 啦 現在 才 大三 那 在 高雄 噠 那個
 méi yǒu la xiàn zài cái dà sān nà zài kāo xióng en nà ge

⁹ NA Qiu-Yu-Fen their family Xie-Yu-Jiao although the second the third floor but she still prov provided one million to NAGE add one add one million.

¹⁰ NA I still have NA children HE they almost at the age of two three I kept bringing out I then just like let him go out to look because here there is really nothing, if really want to have fun, have to HE have clothes A this kind my oldest HE AI from here, just staying here NUO it doesn't work NA you have to go to Tai-Zhong to buy everything.

¹¹ NA your NA paper is, how, almost done BA.

¹² HE1 yeah like them next to us, NAGE building OU it is probably the third floor owned by the oldest and the second has already sold.

高雄	醫學院	嘛	呵	生物系	現在	已經	三年級了			
kāo xióngyī xué yuàn		mǎ	he	shēng wù xì	xiàn zài	yǐ jīng	sān nián jí le ¹³			
(dialogue-02)										
有	你	如果	有	高考	的	話	他	<u>那個</u>	好像	會
yǒu	nǐ	rú guǒ	yǒu	gāo kǎo	de	huà	tā	nà ge	hǎo xiàng huì	
升級	<u>那個</u>	加級	會	比較	比較	快	哦	拿	薪水	也
shēng jí	nà ge	jiā jí	huì	bǐ jiào	bǐ jiào	kuài	o	ná	xīn shuǐ yě	
是	當然	是	比較	好一點						
shì	dān rán	shì	bǐ jiào	hǎo yì diǎn (dialogue-03)						

Regarding the usage of discourse particles and the particle-like 那 (nà) and 這 (zhè) words in the dialogues, we clearly found that given a discourse particle or a particle-like word, its sentential position is in many cases predictable. This typical characterisation of discourse particles may be applied to design an efficient parsing strategy for the natural language processing systems (Carbonell & Hayes 1983, Fischer & Brandt-Pook 1998), especially for spoken language systems to recognise the turn structure in conversation.

4 Frequent Words in Discourse

As mentioned earlier, in addition to discourse particles, lexical words may also be discourse markers, when they are of discourse use. In this section, we look at frequently used words in discourse to identify words highlighting important positions within utterances. Two groups of them are analysed: frequently used words in all three dialogues and frequently used words in turn-initial positions.

4.1 Frequent Words in Dialogues

The first 100 most frequently produced words in the dialogues were collected. Results show that 36 words among them were found in all three dialogues, although the subjects have different education background and they have spoken on diverse topics. These include 7 verbs: 在 (zài), 是 (shì), 就是 (jiù shì), 說 (shuō), 去 (qù), 要 (yào), 有 (yǒu)¹⁵, 6 particles: 哦 (o), 嗯 (en), 哎 (ai), 啦 (la), 啊 (a), 嘛 (mǎ), 5 adverbials: 也 (yě), 就 (jiù), 都 (dōu), 很 (hěn), 對 (duì)¹⁶, 4 grammatical particles: 呢 (ne), 嗎 (ma), 了 (le), 的 (de), 4 nouns: 話 (huà), 時候 (shí hòu), 人 (rén), 小孩子 (xiǎo hái zi)¹⁷, 3 nà and zhè words: 這樣 (zhè yàng), 那個 (nà ge), 那 (nà)), 3 pronouns: 他 (tā), 我 (wǒ), 你 (nǐ)¹⁸, 2 negation: 不 (bù), 沒有 (méi yǒu), 1 adjective: 好 (hǎo) and 1 connective: 所以 (suǒ yǐ)¹⁹.

Six particles, four adverbials, three 那 (nà) and 這 (zhè) words and one connective can evidently be characterised as discourse markers, because the original semantic meaning of these words seems to decrease. Their use becomes more pragmatic and their function as indicator of structuring

¹³ Nothing LA now he is junior NA in Kao-Xiong EN NAGE Kao-Xiong medical college MA HE department of biology now already the third year.

¹⁴ Have if you have passed the exam for officials he NAGE would be promoted NAGE promote would more more fast O the pay is of course is a little bit better.

¹⁵ 在, 是, 就是, 說, 去, 要 and 有 mean "is located", "is", "that is", "say", "go", "want" and "have", respectively.

¹⁶ 也, 就, 都, 很 and 對 mean "too", "then", "all", "very" and "correct", respectively.

¹⁷ 話, 時候, 人 and 小孩子 mean "words, or cases", "time", "man, or people" and "kids", respectively.

¹⁸ 他, 我 and 你 mean "he", "I" and "you", respectively.

¹⁹ 不, 沒有, 好, 所以 mean "no", "not", "good, or well" and "so", respectively.

utterances in conversation is also strengthened. As shown in Example 6, the adjective 所以 (so, therefore) is used to begin a conclusion or a confirmation.

Example 6: 所以 久而 久而久之 小孩子 的 模式
 suǒ yī jiǔ ér jiǔ ér jiǔ zhī xiǎo hái zi de mó shì
 就 會 跟 我們 一樣 (translation cf. Example 1)
 jiù huì gēn wǒ mén yí yàng (dialogue-01)
所以 說 鄉下 呵 就是 這樣 很 純樸 啦²⁰
 suǒ yī shuō xiāng xià he jiù shì zhè yàng hěn chún pú la (dialogue-02)
 那 所以 我 就是 在 那邊 先 寫一寫 然後 再
 nà suǒ yī wǒ jiù shì zài nà biān xiān xiě yì xiě rán hòu zài
 回到 這邊 來 debugging²¹
 huí dào zhè biān lái (dialogue-03)

4.2 Frequent Words in Turn-Initial Position

Words used to initiate turns mark the beginning position of all turns and at the same time they are also beginnings of a number of utterances. We have listed all turn-initial words with word frequency ranking. The ten most frequently used turn-initial words are given in Table 3, in terms of our six subjects respectively. It shows a clear group pattern of word types. Astonishingly, among the most frequently used turn-initial words, all six speakers have preferred 嗯 (en), 哎 (ai), 對 (duì) and 噉 (ou). 那 (nà) and 我 (wǒ) were also preferably used by five speakers. It seems that the syntactic positions of these words are related to particular pragmatic purposes. Chui (Chui 2000) interestingly discussed the case of 對 (duì) in conversation with regard to the ritualization of its pragmatic functions.

Table 3: Turn-Initial Words with Word Frequency

D1-A		D1-B		D2-A		D2-B		D3-A		D3-B	
嗯(en)	37	嗯(en)	45	嗯(en)	18	噉(en)	9	噉(hài)	22	噉(en)	18
他(tā)	10	那(nà)	9	噉(ou)	8	對(duì)	8	那(nà)	18	那(nà)	17
哎(ai)	9	哎(ai)	8	噉(hài)	7	是(shì)	8	噉(ou)	7	啊(a)	9
對(duì)	7	對(duì)	7	這樣子	7	哎(ai)	8	我(wǒ)	7	噉(ou)	6
噉(ou)	6	噉(ou)	7	(zhè yàng zi)		噉(o)	7	對啊(duì a)	6	噉(ai)	5
她(tā)	5	他(tā)	6	噉(o)	6	呵(he)	7	哎(ai)	6	我(wǒ)	4
那(nà)	4	我(wǒ)	5	那(nà)	6	噉(ou)	5	對(duì)	4	呢(e)	4
我(wǒ)	4	她(tā)	5	對(duì)	5	啊(a)	5	噉(en)	4	對啊	3
有(yǒu)	4	就(jiù)	4	噉噉(en en)	5	我(wǒ)	5	這樣	4	(duì a)	
好像(hǎo xiàng)	4	噉(o)	4	真的	5	呢(e)	5	(zhè yàng)	4	就(jiù)	3
				(zhēn de)				噉(o)	4	噉(o)	3
				哎(ai)	5						

By analysing all turn-initial words in the dialogues, we found that 30% of the overall turns were initiated by discourse particles 噉 (en), 噉 (ou), 噉 (hai), 哎 (ai), 噉 (o), 啊 (a) and 呢 (e). 對 (duì)

²⁰ So in the countryside HE its so very plain LA.

²¹ NA so I just write something there then come back here debugging.

and 是 (shì)²² make up 5%, whereas *nà-zhè* words 那 (nà), 這樣子 (zhè yàng zi)²³ and 這樣 (zhè yàng) are 6.6%. In other words, these 12 words highlight the beginning of more than 40% of turns in our data. This is an amazing result. However, this does not necessarily lead to the conclusion that the detection of these words means the detection of turn beginning, because in addition to turn-initial positions, these words are also quite frequently used in other sentential positions. We need to carry out further normalisation analyses to see if the occurrences of these words are in fact more frequent in the turn-initial positions than elsewhere in discourse. Besides, it would be interesting to see if turn-initial markers are also to be frequently found in utterance-initial positions. This is not dealt with in this paper. If this turns out to be true, it would help us execute the segmentation of spoken discourse to a great extent.

5 Chinese Speech Disfluency

More and more researchers have started to work on Chinese repairs regarding the organization of conversation and the natural language processing (Chui 1996, Lee & Chen 1997). Among our data, 373 overt immediate speech repairs were identified. This does not include simple hesitations and partial repetitions. 27 speech repairs among them contain an editing term. These are 哎 (ai), 啊 (a), 呵 (he), 呃 (e), 噉 (ou), 呢 (ne), 嘿 (hei) and 啦 (la). It makes up 7.2 % of the overall speech repairs. Both Labov (Labov 1966) and Hindle (Hindle 1983) approve the edit signal hypothesis: speech disfluency, “non-fluency” as called by them, is usually accompanied by editing terms. Editing terms functionally signal the location of speech repairs. But this empirical result on Mandarin Chinese data does not seem to support the edit signal hypothesis.

It is also to be noted that some discourse particles are specifically used under certain circumstances, i.e. they appear only in particular positions within utterances. For instance, among the eight distinctive editing terms found in speech repairs, the most frequent discourse particle found in the turn-initial position 嗯 is not included. 嗯 is not used to indicate doubt or the self-monitoring of the speaker, but to signal the intention of speaker to be ready to take over the turn or to be used as “thinking pause” of the speaker. This strongly supports the notion that discourse particles do have their independent pragmatic function and “lexicalised meaning” (Fischer & Johantokrax 1995), when they are used to strengthen specific pragmatic implication in discourse.

With regard to the number of words, 10.4% of the overall words in TWPTH are involved in repairing sequences (detailed definitions of speech repairs cf. Tseng 1999). Analysing phrases containing repairs, we obtained the following result (Tseng 2000). Chinese speech repairs are most likely to be found in verb and noun phrases. They make up 37.7% and 41.2% of the total speech repairs, respectively. The most likely position for the Chinese speech repairs to be initiated is the phrasal boundary, where the second most likely position is the morpheme with the central semantic content of the problem word involved.

This clearly shows an important relationship between spoken utterances and syntax. To be more specific, syntactic features play a role for the speakers to process the organization of their speech. When they have to interrupt their utterances or to resume their speech, they prefer certain syntactic locations. Speech disfluency interacts with the syntactic organisation of utterances to the extent that speech disfluency gives cues to important positions in spoken utterances which are especially difficult to deal with in natural language processing systems.

²²對 and 是 can also be used as adjective *correct* and verb *is*. But when they appear at the turn-initial position, they usually confirm prior utterances and they should be translated into *right* and *yes*.

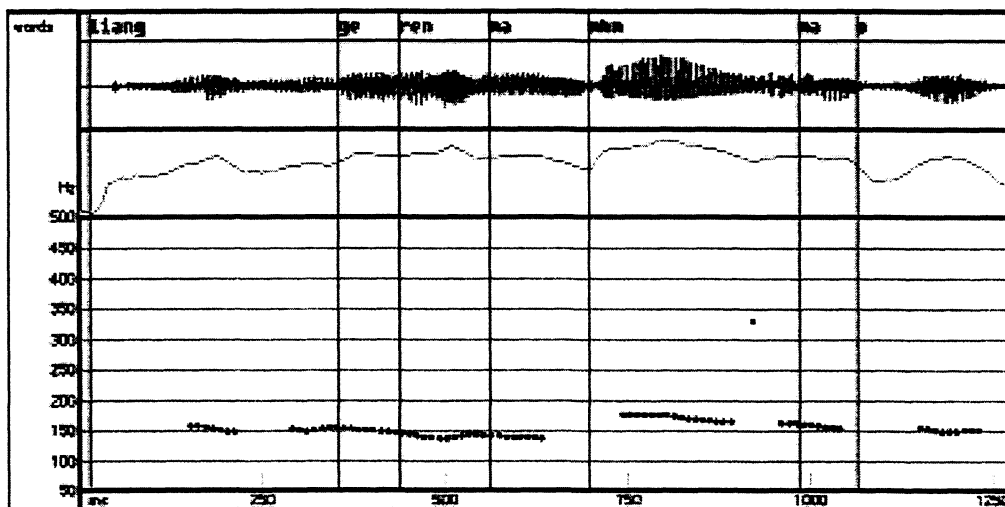
²³這樣子 and 這樣 have similar meaning “this way”, or “so”.

6 Pitch Contour of Discourse Particles and Disfluency

A crucial means for highlighting particular speech sequences is prosody. Especially in spoken language, pitch height, amplitude and melodic as well as intonational variants are all employed to emphasise intentions of the speakers (Levelt & Cutler 1983). Here, we are not tempting to cover all prosodic aspects in conversation. Because pitch height obviously plays the most important role in determining the prosodic properties of speech data we preliminarily restrict the prosodic analysis to pitch contour. The TWPTH data were digitally recorded at 44.1 kHz sampling rate. The prosodic analysis was carried out using PitchWorks developed by SCICON R&D in California.

We first look at the utterance segment 兩個人嘛噯嘛哦²⁴ in Figure 1, where there are three discourse particles occurring four times: 嘛 (mǎ) twice, 哦 (o) and the hesitation discourse particle 噯 (en, mhm) once. Regarding the fundamental frequency (F0) contour of the discourse particles 嘛 (mǎ), 哦 (o) and 噯 (en, mhm), we found a variety of falling and rising combinations. Even the well-known hesitation discourse particle mhm, typically having extremely flat F0 contour, has falls and rises. Because of the scarcity of data presented here, more prosodic data need to be obtained to determine the prosodic features of discourse particles in Mandarin.

Figure 1: F0 Contour of Some Discourse Particles.

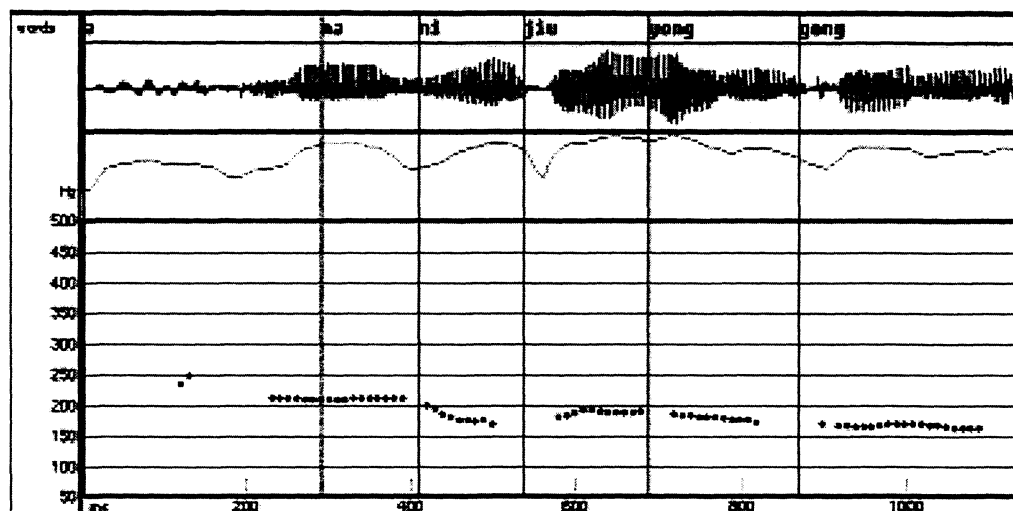


The pitch contour of 那 (nà) in the utterance 啊那你就用功²⁵ as particle-like words (instead of as determiners) is first flat, slightly rising, then at the end of the word slightly falling. This is illustrated in Figure 2. The slightly falling F0 contour at the end of the word could be related to the falling lexical tone of 那 (nà). Or it could possibly have something to do with the pragmatic functions of 那 (nà) to indicate the intention of the speaker to take over the turn. Nevertheless, this preliminary result raises an interesting issue: are there any prosodic distinguishable differences between 那 (nà) and 這 (zhè) words as particle-like words and as determiners? Is the syntactic and pragmatic change of the use of 那 (nà) and 這 (zhè) words marked by prosody?

²⁴ liǎng ge rén mā mǎ (en) mā o. (Two people MA MHM MA O.)

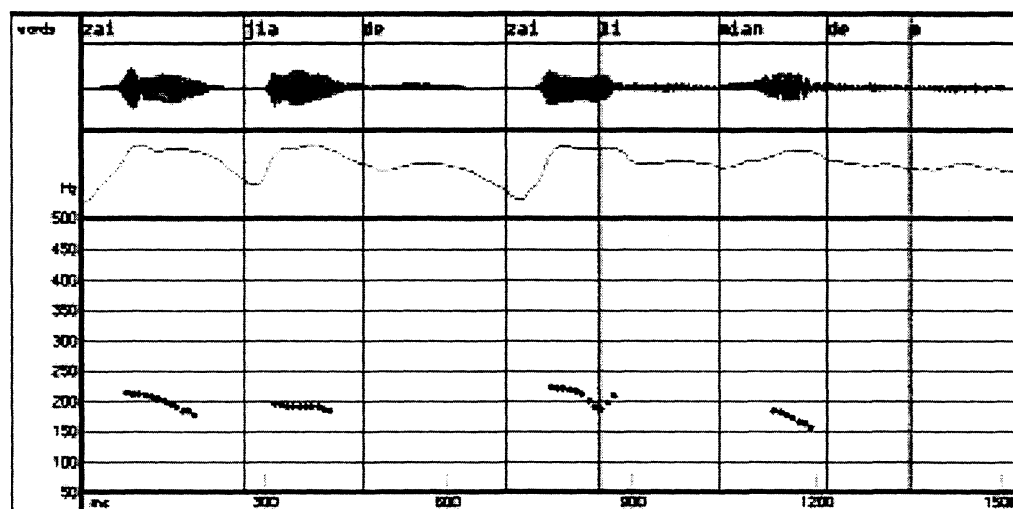
²⁵ a nà nǐ jiù yòng gōng. (A then you should work hard.)

Figure 2: F0 Contour of Particle-like Nà.



Following the baseline declination, the reset hypothesis related to the prosodic characterisation of speech repairs proposes that the F0 value after the interruption is higher than expected. And the previous sentential segment produced with a similar F0 value is very likely to be the corresponding sentential position before the interruption. This phenomenon can be clearly observed in Figure 3. Within the repair in the utterance 在家的在裡面的哦²⁶, the pitch height around the resumption location within the speech repairs, the second 在 (zài), is at about the same height as that around the corresponding position in the reparandum, the first 在 (zài).

Figure 3: F0 Contour around Resumption Location in Speech Disfluency



²⁶ zài jiā de zài lǐ miàn de o. (At home, inside O.)

7 Conclusion

Based on the results of a corpus analysis, the means of highlighting specific utterance components in Chinese spoken dialogues were investigated and the interrelationship between these means were discussed in this paper. Tagged and annotated speech data were used to analyse frequent types of word sequences including their frequency of occurrences, their position within utterances and their prosodic representation. A limited number of discourse particles and words were found in regular utterance positions (utterance-initial, mid-utterance and utterance-final; turn-initial) to highlight the particularity of the internal structure of utterances and dialogues. The interaction between syntax and pragmatics with respect to these phenomena needs to be further examined. Currently, statistical analyses on the prosodic representation of discourse particles and speech repairs are in progress.

Acknowledgements

The study presented in this paper is in part supported by National Science Council (NSC 89-2411-H-001-098). I'd like to thank the Industrial Research Technology Institute (IRTI) for generously providing the TWPTH corpus and Hui-Hsin Tseng for her work on annotating the speech data.

References

- Carbonell, J. and Hayes, P. 1983. Recovery Strategies for Parsing Extragrammatical Language. *American Journal of Computational Linguistics*. (3/4):123-146.
- Chen, K.-J., Huang, C.-R., Chang, L.-P. and Hsu, H.-L. 1996. SINICA CORPUS: Design Methodology for Balanced Corpora. *PACLIC 11*. 167-176.
- Chui, K.-W. 2000. Ritualization in Evolving Pragmatic Functions: A Case Study of DUI. In *Proc. of the 7th International Symposium on Chinese Language and Linguistics*. 177-192. Jia-Yi.
- Chui, K.-W. 1996. Organization of Repair in Chinese Conversation. *Text 16/3*. 343-372.
- CKIP. 1995. Sinica Balanced Corpus. Technical Report no. 95-02/98-04. (in Chinese)
- Clark, H. and Brennan, S. 1991. Grounding in Communication. In *Perspectives on Socially Shared Cognition*. Resnik/Levine/Tealsey (eds.). 127-149.
- Fischer, K. and Brandt-Pook, H. 1998. Automatic Disambiguation of Discourse Particles. In: *Proc. of Coling/ACL '98 Workshop on Discourse Relations and Discourse Markers*, Montreal, Canada, 107-113.
- Fischer, K. and Johanntokrax, M. 1995. Ein linguistisches Merkmalsmodell für die Lexikalisierung von diskursgesteuerten Partikeln. *SFB 360*. No. 8.
- Gross, D., Allen, J. and Traum, D. 1993. The TRAINS 91 Dialogues. Technical Report 92-1, Dept. of Computer Science. University of Rochester.
- Heeman, P. and Allen, J. 1999. Speech Repairs, Intonational Phrases and Discourse Markers: Modelling Speakers' Utterances in Spoken Dialogue. *Computational Linguistics 25/4*. 527-571.
- Hindle, D. 1983. Deterministic Parsing of Syntactic Non-fluencies. In *Proc. of ACL '83*. 123-128.
- Hirschberg, J. and Litman, D. 1993. Empirical Studies on the Disambiguation of Cue Phrases. *Computational Linguistics*, 19(3): 501-530.
- Hovy, E. and Scott, D. (eds.) 1996. *Computational and Conversational Discourse. Burning Issues - An Interdisciplinary Account*. Springer.

- Labov, W. 1966. On the Grammaticality of Every-day Speech. In Annual Meeting of the Linguistic Society of America. 41. New York.
- Lee, Y.-S. and Chen, H.-H. 1997. Using Acoustic and Prosodic Cues to Correct Chinese Speech Repairs. In Proc. of EUROSPEECH '97. 2211-2214. Rhodes.
- Levelt, W. J. 1983. Monitoring and Self-Repair in Speech. *Cognition*. 14. 41-104.
- Levelt, W. J. and Cutler, A. 1983. Prosodic marking in Speech Repair. *Journal of Semantics*. 2(2):205-217.
- Nakatani, C., Hirschberg, J. and Grosz, B. 1995. Discourse Structure in Spoken Language: Studies on Speech Corpora. Presented at the AAA-I-95 Spring Symposium on Empirical Methods in Discourse Interpretation and Generation.
- Sacks, H., Schegloff, E. and Jefferson, G. 1974. A Simplest Systematics for the Organization of Turn-Taking for Conversation. *Language*. 50(4): 696-735.
- Schegloff, E., Jefferson, G. and Sacks, H. 1977. The Preference of Self-Correction in the Organization of Repair in Conversation. *Language*. 53(2): 361-382.
- Schiffirin, D. 1987. *Discourse Particles*. Cambridge University Press.
- Shriberg, E. and Lickley, R. 1992. Intonation of Clause-Internal Filled Pauses. In Proc. of ICSLP 92. Banff. 991-994.
- Taylor, T. and Cameron, D. 1987. *Analysing Conversation*. Pergamon Press.
- Traum, D. and Heeman, P. 1997. Utterance Unit in Spoken Dialogue. In *Dialogue Processing in Spoken Language Systems*. Maier/Mast/Luper Foy (eds.). Lecture Notes in Artificial Intelligence. Springer Verlag.
- Tseng, S.-C. 2000. Repair Patterns in Spontaneous Chinese Dialogs: Morphemes, Words and Phrases. In Proc. of ICSLP 2000. 453-456.
- Tseng, S.-C. 1999. *Grammar, Prosody and Speech Disfluencies in Spoken Dialogues*. PhD Thesis. University of Bielefeld.