

本期要目

壹. ROCLING-2007 CFP

貳. 學術活動預告 (OntoLex 07, PACLIC 21)

參. 專文—語音韻律簡介(江振宇、蕭希群、余秀敏、廖元甫)

肆. 會員個人資料更新通知單

第二頁

第三~四頁

第五~十九頁

第二十頁

九十七年度會費開始繳交

九十七年度「一般個人」及「學生」會員有效期間即將於六月三十日到期，為保障各位會員之權益，敬請各位會員如期繳交會費，若您個人資料有所異動，亦請惠予通知本會，更新個人資料通知單及信用卡繳費單請參閱本訊第 20 頁。

博碩士論文獎 7/1 開始申請

名額及獎項：

1. 博士論文優等獎：一名，獎金二萬元，並頒給學生及指導教授獎狀各乙紙。
2. 博士論文佳作獎：一名，獎金一萬元，並頒給學生及指導教授獎狀各乙紙。
3. 碩士論文優等獎：一名，獎金一萬元，並頒給學生及指導教授獎狀各乙紙。
4. 碩士論文佳作獎：三名，獎金各伍仟元，並頒給學生及指導教授獎狀各乙紙。

申請資格及限制：

1. 國內大專院校博碩士班應屆畢業生從事計算語言學相關研究方向者，由其指導教授推薦。
2. 參賽限制：每位指導教授以推薦一篇博士論文及兩篇碩士論文為限。(含個人指導與共同指導)。

申請期限 7/31 止，申請程序及辦法請參閱本會網頁：<http://www.aclclp.org.tw/doc/shipreg.htm>。

光華雜誌 1976~2005 語料庫簡介

《光華雜誌》為全世界研究「台灣學」的珍貴資料，經由遠流集團智慧藏學習科技公司耗資千萬的數位化工程之後，光華雜誌自 1976~2005 年共三十年的五千萬中英對照文字與三萬張珍貴影像，均完整收錄在《光華雜誌智慧藏》知識庫系統中。其中，五千萬中英對照文字資料庫，是台灣唯一一套「以優美的中英文介紹台灣文化、歷史、生活、人物的中英對照語料庫」。

從台灣的文化（大稻埕/鹿港仔/布袋戲/歌仔戲/陶瓷藝術/篆刻藝術/龍山寺/三峽祖師廟/捏麵人與吹糖人/金馬奔騰/原位民文化/民歌年代...）、台灣的變遷（原住民社會/十大建設/美麗島事件/農村沿革/外銷王國/解嚴時期/電子工業/全民健保/港澳回歸/網路經濟/政黨輪替...）、到台灣的人物（李天祿/林懷民/席德進/劉其偉/黃春明/侯孝賢/馬友友/柏楊/趙耀東/李遠哲/王永慶/吳舜文/鄧麗君/李安/何大一...），這些珍貴的中英文字紀錄，都在光華雜誌語料庫中，不僅是學校推廣本土英語教育的最佳教材，更是研究語言學相關技術的學術單位必備語料庫。

遠流集團智慧藏學習科技公司為推廣自然語言處理相關技術之應用，特別授權本會代為發行「光華雜誌 1976 年~2005 年語料庫」，授權使用價格為 NT\$40,000 元，使用期限三年，語料庫範例及申請辦法請參閱本會網頁：

http://www.aclclp.org.tw/use_gh_c.php。

Rocling 2007 Call for Papers

The 19th Conference on Computational Linguistics and Speech Processing

(<http://nlg12.csie.ntu.edu.tw/~Rocling/>)

Taipei, Taiwan, 6-7 September 2007

The 19th Conference on Computational Linguistics and Speech Processing (Rocling 2007) will be organized by National Taiwan University in Taipei, Taiwan on September 6-7, 2007. ROCLING is the most historied and comprehensive conference focused on computational linguistics, speech processing, and related areas in Taiwan.

ROCLING XIX will be hosted by National Taiwan University. The conference will feature invited lectures, panel discussions, and lecture and poster sessions.

SCOPE

Rocling 2007 invites the submission of original papers in all areas of natural language and speech processing, including, but not limited to:

- | | |
|---|--------------------------------------|
| (a) cognitive/psychological linguistics | (l) semantic web |
| (b) discourse/dialogue modeling | (m) semantics/pragmatics |
| (c) information extraction/text mining | (n) speech analysis/synthesis |
| (d) information retrieval | (o) speech recognition/understanding |
| (e) language understanding/generation | (p) spoken dialog systems |
| (f) lexicon/morphology | (q) spoken language processing |
| (g) machine translation/multilingual processing | (r) syntax/parsing |
| (h) named entity recognition | (s) text summarization |
| (i) NLP applications/tools/resources | (t) web knowledge discovery |
| (j) phonetics/phonology | (u) word segmentation/POS tagging |
| (k) question answering | (v) others |

Best Paper Award:

The best paper will be selected and announced at ROCLING XIX.

DEADLINES:

- **Paper submission deadline:** July 2, 2007
- Notification of acceptance: August 10, 2007
- Camera ready papers due: August 20, 2007

ORGANISING COMMITTEES:

Conference Chair:

Hsin-Hsi Chen, National Taiwan University

Program Co-Chairs:

Kuang-Hua Chen, National Taiwan University

Berlin Chen, National Taiwan Normal University

Organizing Chair:

Lun-Wei Ku, National Taiwan University

1st Call for Papers for OntoLex 07

From Text to Knowledge: The Lexicon/Ontology Interface
Workshop at ISWC07 (6th International Semantic Web Conference)

<http://olp.dfki.de/OntoLex07/>

November 11th , 2007, Busan, South-Korea

Endorsed by ACL SIGSEM (<http://www.aclweb.org/sigsem>) and
The NEON project (<http://www.neon-project.org/>)

The OntoLex workshop series is concerned with the interface between knowledge representation in ontologies and the representation and use of linguistic knowledge as encoded in (multilingual) lexicons. Previous OntoLex workshops were held in Sozopol (2000), Las Palmas (2002), Lisboa (2004), Jeju (2005), Genoa (2006). The extraction of ontology elements from textual data is a pre-requisite for many Semantic Web applications. Hence the interface between ontologies (describing objects in a domain) and the lexicon (describing the linguistic features of terms that refer to such objects) is becoming increasingly important. We invite submission of papers on the following topics in this interdisciplinary research area:

- * The lexicon/ontology interface
- * The lexicon in Semantic Web applications
- * Ontologies in knowledge-based Natural Language Processing

Important Dates

- August 1, 2007: Submission Deadline
- September 1, 2007: Notification
- October 1, 2007: Camera-ready Version
- November 11, 2007: Workshop

Submission

Papers must be in English, not exceed 10 pages and should be formatted in LNCS style (Springer format for Lecture Notes in Computer Science series). For complete details, see Springer's Author Instructions (<http://www.springer.com>). Papers must be submitted electronically through the conference submission site (<http://iswc2007.semanticweb.org>)

Organizing Committee

Paul Buitelaar, Competence Center Semantic Web - DFKI, Germany
Key-Sun Choi, Semantic Web Research Center - KAIST, South-Korea
Aldo Gangemi, Laboratory for Applied Ontology - CNR, Italy
Chu-Ren Huang, Institute of Linguistics - Academia Sinica, Taiwan
Alessandro Oltramari, Cogito S.r.l., Italy

PACLIC 21 Call for Papers

Seoul, Korea, November 1-3 2007

<http://paclic21.snu.ac.kr>

Annual meetings of the Pacific Asia Conference on Language, Information and Computation (PACLIC) have provided places where scholars in theoretical and computational linguistics mainly from the Pacific Asia region get together to share new ideas about language, information and computation.

We invite prospective authors to submit abridged versions of original and unpublished papers on all aspects of theoretical and computational linguistics. Possible topics include but are not limited to: (1) language, (2) information, and (3) computation. Priority will be given to papers that are interdisciplinary and cover a wide range of languages in Asia. Outstanding papers will be considered as candidates for the best paper prize.

Abridged papers should not exceed 3,000 words including the bibliography. Author information should appear on a separate page with (1) the topic area of the paper: (language and/or information and/or computation), (2) the title of the paper, (3) name(s) of the author(s), (4) affiliation(s), (5) postal address(es), (6) e-mail address for correspondence, and (7) the preference for oral or poster presentation. Accepted papers will be published in the Conference Proceedings. Proceedings from the last five PACLIC's (PACLIC 16-20, 2002-2006) are indexed in ISIProceedings (formerly ISTP). All future PACLIC Proceedings will be submitted to ISI for indexing. The deadline for paper submission is June 30, 2007. Papers should be submitted online at <http://paclic21.snu.ac.kr>.

Important dates of PACLIC21:

- Due date for abridged paper submission: June 30, 2007
- Acceptance notification to authors: August 15, 2007
- Due date for the camera-ready full paper: September 15, 2007
- Due date for early registration: October 1, 2007
- PACLIC21 meeting: November 1-3, 2007

Inquiries: program committee chair Hee-Rahk Chae at paclic21@hufs.ac.kr

Invited Speakers:

Harry Bunt (Tilburg University)

Jerry Hobbs (University of Southern California)

Byung-Soo Park (Kyunghye University)

語音韻律簡介

江振宇、蕭希群、*余秀敏、**廖元甫

國立交通大學電信工程研究所

*中華大學外文學系

**國立台北科技大學電子工程研究所

gene.cm91g@nctu.edu.tw balade29@yahoo.com.tw

*kuo@chu.edu.tw **yfiao@ntut.edu.tw

1. 前言

當我們透過聲音與他人相互溝通時，我們聽到的不只是對方所要表達的內容，同時也接收到對方所發出的高低起伏的情緒、抑揚頓挫的口氣、以及輕重緩急的語流，因此，往往同一句話、隨著說話人的「語音表情」，會讓我們感到對方說得興高采烈、或義憤填膺、或語氣堅定、或感傷悲戚，這裡所謂的「語音表情」，指的就是語音的韻律變化（prosodic variations）。一般所謂的韻律變化是指人們利用(1)音量的大小（energy levels）、(2)音高的高低（pitch contours）、(3)說話速度的快慢（speaking rates）、以及(4)停頓的長短（pause durations）等四個韻律參數，來建構語音的表情，藉以強化語者本人所要表現的語意、情感、重點、立場等等。一個人說話若缺乏了語音的韻律變化，將會說出聽起來單調、貧乏的聲音，而無法引起他人的注意，因此也就無法達到溝通的目標。因此，從語音合成的角度來看，若要系統發出讓人耳聽起來流暢、自然、引人注意的話語，勢必要加入韻律訊息，才能發展優質的語音合成系統；反過來說，站在語音辨認的立場，韻律訊息的運用，能增加語音內容的瞭解、提高系統辨識的能力，因此，對一特定語言韻律訊息的掌握，有助於各項語音處理技術的發展，而本文的目的即在介紹語音韻律的產生和變化，我們將在第 2 節介紹語音產生的過程，以來說明影響韻律的因素以及產生的機制；第 3 節介紹國語語音韻律的特性，提供讀者簡單易懂的一些概念；接著我們在第 4 節介紹國語文字轉語音技術的核心「國語韻律的產生」，第 5 節則介紹目前建構國語韻律模型的方法，最後筆者將討論一些國語韻律模型未來的研究方向。

2. 語音與韻律產生的過程

圖一為一個韻律模型與韻律訊息產生的流程 [1]，我們將由左往右說明這個流程。從圖一的最左欄開始，我們可以看到輸入部分的資訊（input information），表示一個韻律訊息的產生，一開始人們必須具備三大類的資訊：(A)語言的資訊（linguistic information）(B)副語言的資訊（para-linguistic information）以及(C)非語言的資訊（non-linguistic information），其中：

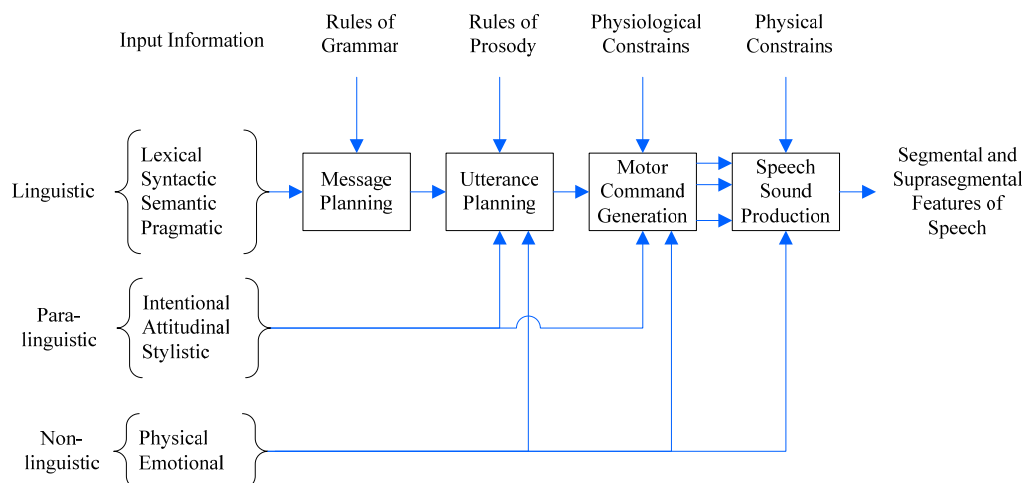
A. 語言的資訊（linguistic information）包含下列四個層次的資訊：

1. 詞彙層次的資訊（lexical information），例如「詞」、「字」以及「詞/字所對應的詞性、意義與發音。
2. 句法層次的資訊（syntactic information），指詞或詞組在句中所扮演的句法功能。
3. 語意資訊（semantic information），指詞或詞組內含的語意特徵。
4. 語用資訊（pragmatic information），指說話情境對語言使用的影響。

B. 副語言學的資訊（para-linguistic information）：是指語者心理內在的意圖、態度、或是說話的風格為語言表現所帶來的影響，例如若要強調「這個東西不重要」中的「不」字，我們會將此字說得較大聲、或是較長一點。儘管這類資訊不屬於語言層次的資訊，

但卻會影響語言最後的表現，因此，稱此類的資訊為副語言學的資訊。

- C. 非語言的資訊 (non-linguistic information)，表現出與年齡、性別、個人的習慣、情緒狀態以及生理結構等對語言的影響，由於這些因素跟語言學的資訊沒有直接的關係，因此，稱此類的資訊為非語言的資訊。



圖一、一個完整的韻律模型與韻律訊息產生流程

接下來，往右四個矩型代表這三類輸入資訊所歷經的四重資訊處理階段，第一階段主要是透過語言的語法規則 (rules of grammar) 綜合各層的语言資訊、規劃出要表達的訊息 (message planning)；第二階段則是在處理完第一階段的語言訊息後，加進 (1)副語言的資訊，如表達各式的語氣、態度、風格，以及 (2)非語言的訊息 — 如性別、情感、生理資訊，這兩類的資訊需要運用語言的韻律規則 (rules of prosody)，也就是調整各種韻律參數的數值，來將此二類的資訊帶進話語表現的籌劃當中 (utterance planning)；接下來第三階段，則是歷經生理的制約 (physiological constraints)、由大腦基於語言表現的籌畫、來下達牽動發音器官的運動神經的指令 (motor command generation)；第四個階段，則是人們透過運動神經與肌肉組織的生理活動 (physical constraints)，發出語音訊號 (speech sound production)，而從語音的產生中，我們最終得到音段與韻律的訊息 (segmental and prosodic/supra-segmental features of speech)。

若從國語韻律訊息的合成的角度來走一次圖一的流程，一開始我們為系統輸入的語言資訊是以文字所寫成的文章形式來呈現，每一篇文章都內含階層性的語言架構，從字 (character) 所構成的最小語意單位詞 (lexical words) 開始、到由幾個詞所構成的詞組 (syntactic phrases)、到幾個詞組構成的句子 (sentences)、到最後以數個句子組成一篇文章。輸入的文章透過語法規則的分析、韻律規則的調整，最後由語音合成器產生包含國語音段與韻律訊息的聲音。以國語一字、一音、一調的特色來看，其語音韻律是以音節為最基本的發音單位，繼而由數個音節構成韻律詞 (prosodic words)、幾個韻律詞再構成韻律片語 (prosodic phrases)，最後由幾個韻律片語構成整個篇章 (discourse)。

以下，我們從最小的韻律單位開始、分層介紹國語的韻律特性，由於較高階層的韻律單位 (韻律詞，韻律片語) 較為複雜，我們將放在第 5 節介紹此韻律架構及其模型，而第 4 節將會介紹目前已成熟的語音韻律應用，也就是文字轉語音系統中的國語韻律產生。

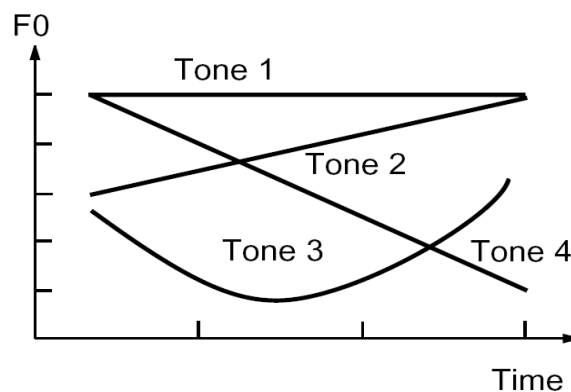
3. 語音韻律的特性

3.1 音節特性及基頻軌跡

先從國語語音的語音基本單元來說明，國語是一個音節搭配一個具有辨別語意功能的聲調，音節的結構也有層次，最上層的兩個音節要素可以分成「聲母」(initial)與「韻母」(final)，第二層中「韻母」又可再細分為「介音」(medials)與「韻腳」(rime)，最後第三層「韻腳」的部分則再次分為「主要元音」(nucleus)和「韻尾」(coda)，這些結構要素中只有「主要元音」是音節的必備成分，其他的要素都是由子音所構成，而且是可有可無的音節成分。國語音節中的子音和母音的組合，必須符合國語的音節的結構以及音節組合律(phonotactics)。

除了基本音節音段資訊以外，國語是一種聲調語言(tonal language)，因此聲調的表示很重要，也是十分重要的韻律特性之一。國語的聲調總共有五個，分別是一聲(tone 1)、二聲(tone 2)、三聲(tone 3)、四聲(tone 4)以及輕聲(tone 5)，聲調的特性尤其顯現在母音的基頻軌跡上，圖二表示的為一聲到四聲基頻軌跡的基本形狀及其相對位置，我們可以發現到我們注音符號上的聲調符號，和實際上的基頻形狀很類似。輕聲的基頻軌跡原則是受到前後語音聲調的影響，音節長度通常比其他四聲的聲調短，音量也較小，因此比較沒有固定的基頻形狀，最明顯且常見的例子就是單字詞「的」。

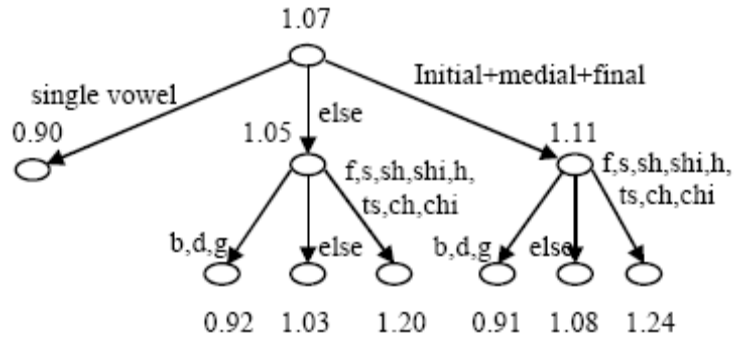
另一個值得注意的國語聲調特性為變調(tone sandhi)規則，最著名的例子就是當連續念兩個三聲的詞，如「雨傘」，第一個三聲要念成二聲，然而當多個三聲詞連在一起組合成一個句子時，變調的方式則涉及句法資訊的影響，例如「我有好幾把小雨傘」 [2]，一般人只會讀成「23223323」而不會說出「22223223」這種變調等等，國語變調與語法的關係，讀者可以詳見於 [2]。



圖二、國語四個聲調的基頻軌跡圖

3.2 國語音節長度的特性

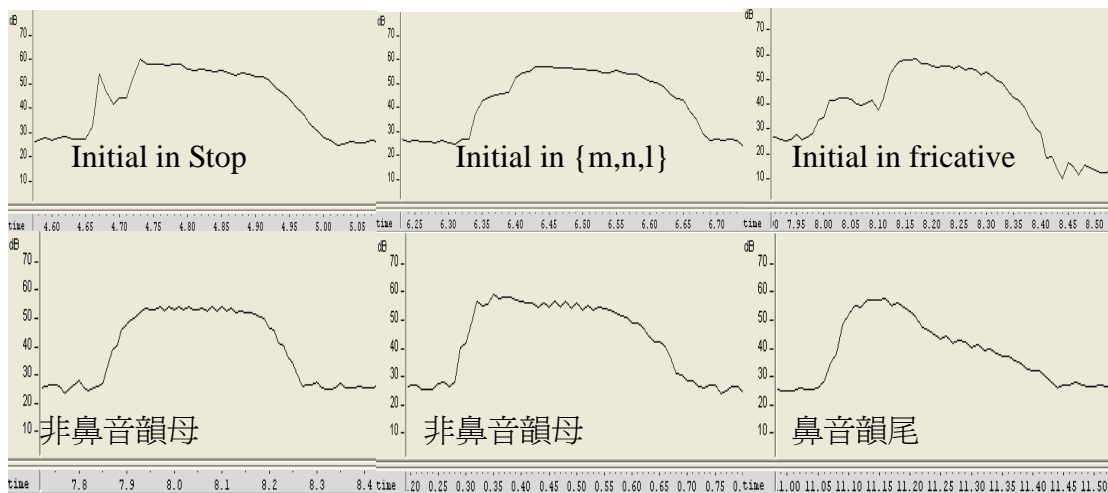
國語帶聲調的音節組合(tonal syllable)大約有 1300 個，而如果扣除聲調的部分，基本音節(base syllable)大約為 417 個，由於音節有各種語音結構(phonetic structures)，也就是說，可以由單一個主要元音構成一個音節(如「衣服」的「衣」字的音節)，也可以多到聲母+介音+主要元音+韻尾等四個成分都齊備來構成一個音節(如「香水」的「香」字的音節)，因此 417 個基本音節的音長並不相等，一般來說只有單一主要元音的音節比較短，四個音節成分都齊備的音節則相對的較長，聲母為摩擦音(fricative)或為塞擦音(affricate)的音節較長，若聲母為塞音(stop，如ㄅㄆㄇ)的音節，其音節長度較短，而且比平均音節的長度還要短。圖三 [3]中表示的為經過大量語音資料統計後，音節長度與音節成分的決策樹關係圖，決策樹上節點的數字代表相對於平均音節長度的比例，大於 1 便表示較長，相反的，小於 1 便表示較短。



圖三、依據音節的語音結構的分類所得的音節平均長度

3.3 國語音節之能量與停頓

國語音節的能量變化是與發音的方法很有關係，也就是說不同的聲母、韻母組合會有不同的音節能量輪廓(energy contour)，最明顯的例子是含有母音「a」的音節，通常比含有母音「i」的音節平均能量還要大，圖四顯示的為幾個主要語音結構組合下音節能量的大致輪廓。



圖四、幾個主要語音結構組合下音節能量的大致輪廓

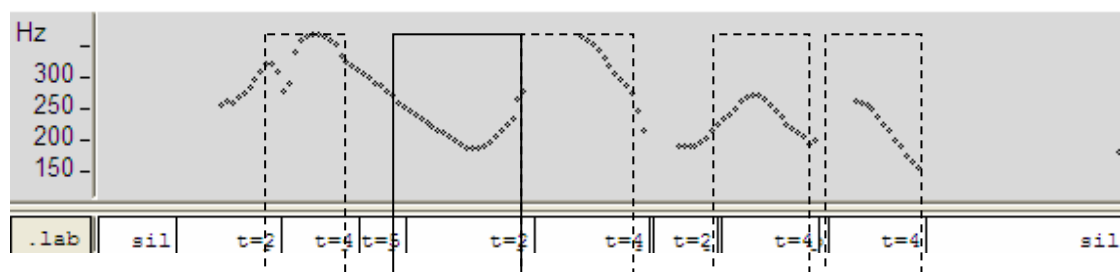
至於停頓方面，由於音節已經是國語語音發音的基本單位了，所以在音節中聲母和韻母或是韻母和韻母的交界處，以及整個音節裡面都不會有停頓發生，而在某些音節發音之前，會有短暫的停頓，這些音節通常是塞音(stop)或塞擦音(affricate)為聲母之音節，這種音節前的短停頓，是不一定完全都會發生，在快速的流利語音之中，相較於較慢速的語音或獨立詞，發生的機率就會少很多。

3.4 大單元之國語韻律

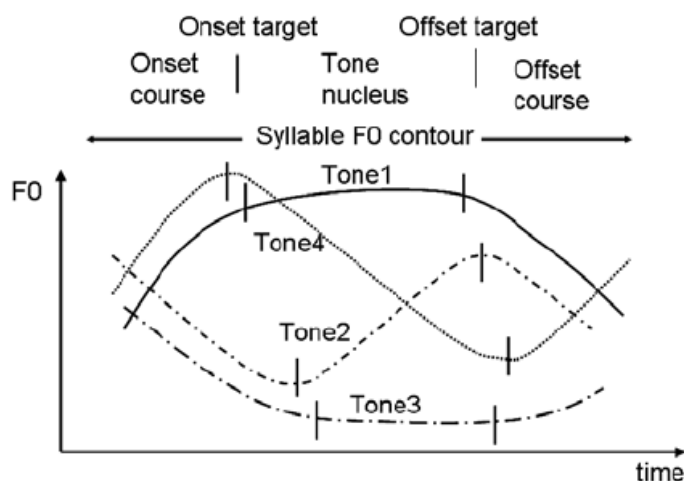
3.1 到 3.3 節之間，介紹了音節層次的國語韻律基本特性，觀察的焦點放在較小的單元，然而連續語音(continuous speech)或自然流利語音(spontaneous speech)的韻律，是由這些小單元所組成的大單元，這個大單元的韻律變化雖然是由音節單元所構成，但由於前後音節和音節間的協同發音(coarticulation)的關係，以及高於音節層次的參數影響，連續語音或自然流利語音的韻律變化顯得十分複雜。

3.4.1 前後音節和語調的影響

如圖五所示，在連續語音之下的一句語音，基頻軌跡的變化已經和我們所認知的形狀有所變化，左邊第一個實線框為一個二聲的音節，然而此音節的基頻軌跡已經和標準的有所差別，看起來已經像是一個三聲，這是因為受到前後語音(context)影響，雖然基本上二聲在頻率是上是由低走向高，然而由於前面的四聲和輕聲頻率較高，聲帶震動無法瞬間由快到慢，因此就產生了二聲前半段與前面音節平順相接且逐漸下降的基頻軌跡，因此就有學者提出核心聲調（tone nucleus） [4]的理論，如圖六所示，這個理論是說我們可以把一個音節的基頻軌跡描述成起始調（onset course）、核心聲調（tone nucleus）以及結尾調（offset course）三個段落，起始調（onset course）指的是音節基頻軌跡受到前面音節影響的不穩定部份，核心聲調(tone nucleus)為一個音節表現聲調的穩定核心部分，而結尾調(offset course)為相接於後面音節的不穩定基頻軌跡的部分，此理論認為聽者對聲調的辨認，是只把核心聲調（tone nucleus）的部份當作辨認的目標，相關的聲調辨認實驗結果某種程度上驗證了這個理論。

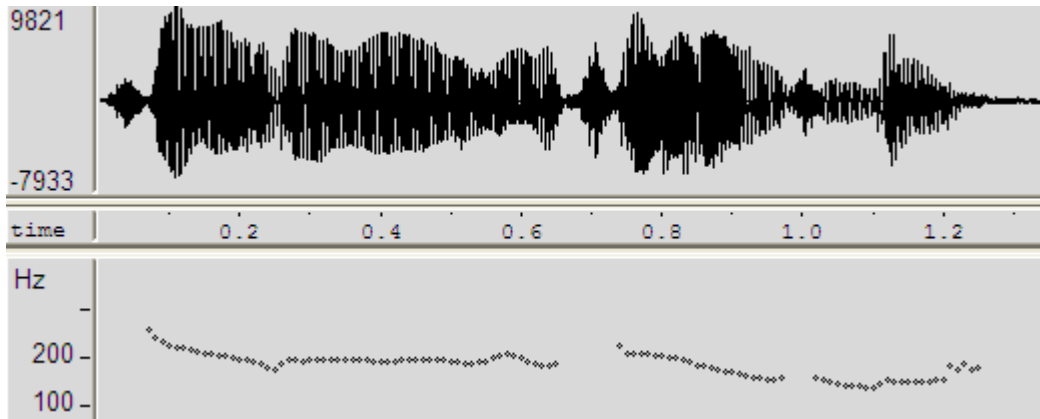


圖五、受前後文及語調影響之 pitch contour

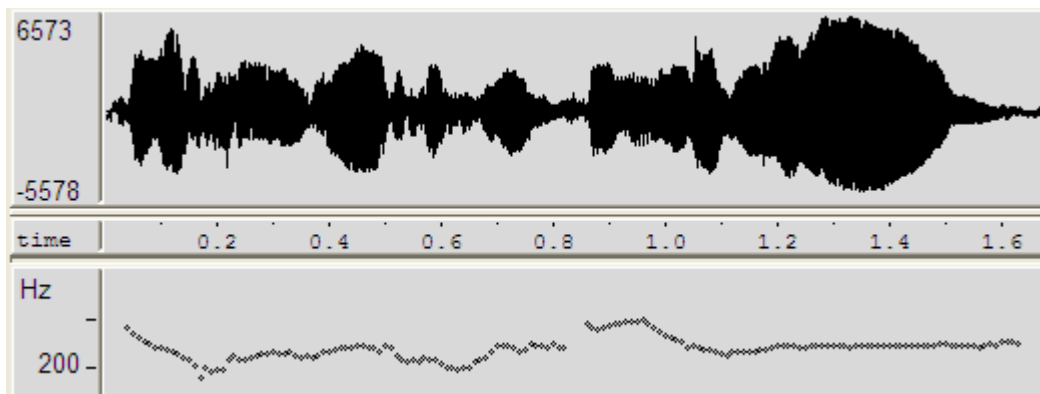


圖六、四個聲調的 tone nucleus

再來我們來觀察圖五中四個虛線框起來的四聲，我們可以發現到越到句尾基頻軌跡越低，這是由於語者講話時的自然現象，我們稱之為語調的變化(intonation)，最明顯的語調例子如圖七及八所示，一般來說敘述句(declarative utterance)基頻軌跡，會隨時間由高慢慢的變化至低處，而疑問句(interrogative utterance)基頻軌跡，可能會在句尾處上揚。



圖七: 敘述句「這個東西是這樣子的」

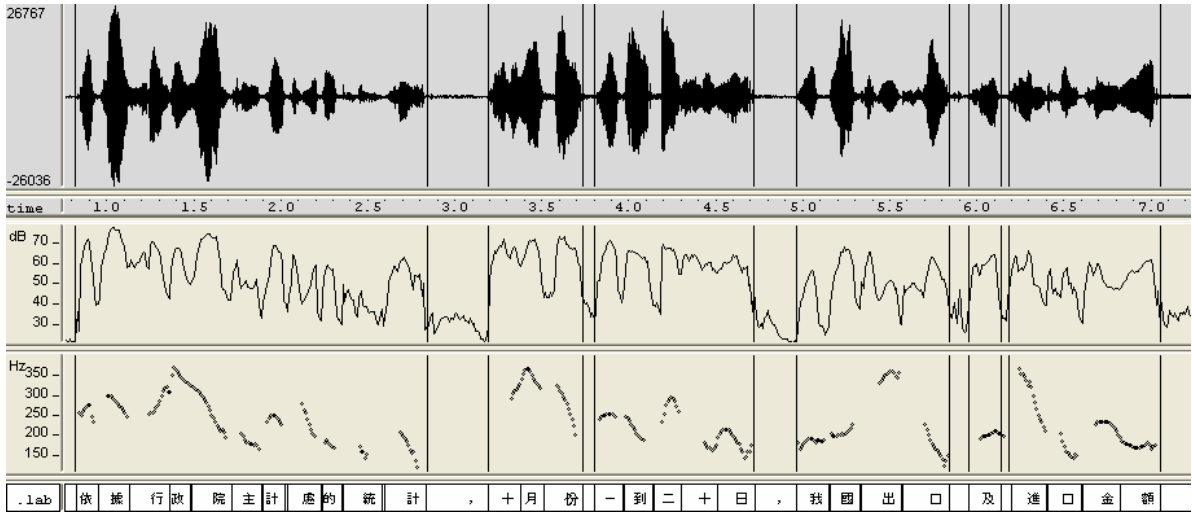


圖八: 疑問句「這個東西是這樣子的嗎?」

前面我們先說明了流利語音和單音節語音韻律的差別，接下來，我們將在下一節的圖九中，顯示一個連續語音的波形、能量軌跡、基頻軌跡、音節切割位置以及對應之文字，我們藉此圖來對較大的國語語音韻律單元特性，以語音韻律的四個面向做較深入的介紹，並和語言參數做對應的比較。

3.4.2 停頓

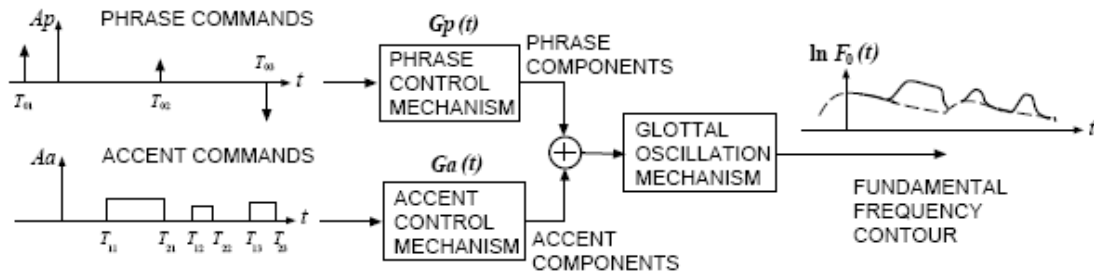
圖九的文句是「依據行政院主計處的統計，十月份一到二十日，我國出口及進口金額…」，由這個例子我們可以發現到，通常人們在朗讀文章時，在標點符號出現的地方，會產生明顯的停頓，來表示另一語意段落的開始；然而在由多音節所構成的詞中，其內部的音節之間 (intra-word boundary) 幾乎不會產生停頓，如「行政院」中的「行/政」和「政/院」交界處。停頓的產生往往在詞和詞邊界的交界處 (inter-word boundary)，才有機會產生停頓，如「十月份一到二十日」中的「份/一」之間，以及「出口及進口」中的「及/進口」的交界處，這是因為 (1)「詞」是基本的語意訊息單位、(2)朗讀者為了能成功的傳達訊息，因此會選擇在適當的詞邊界處做停頓。但同時我們也可以發現到，雖然「詞」是基本的語意單位，但語音上幾個詞可能連著念，組成一個較大的語音單元，經過分析 [5]，除了標點符號對應的停頓以外，在語法片語 (syntactic phrase) 或一些功能詞 (function word) 的邊界，也較容易發生明顯的停頓，在這個例子中，「一到二十日」為一個名詞片語，「及」為連接詞 (功能詞的其中一種)。



圖九: 連續語音之聲音波形、能量軌跡、基頻軌跡、音節切割位置以及文字內容

3.4.3 基頻軌跡

我們先觀察以標點符號分開的語句，大致說來，可以發現到一般語句中的基頻軌跡存在著類似浪潮（tide）般的整體變化，也就是音調的高低在一個句子中並不是一個維持平穩的趨勢，而是有逐漸下降的趨勢，並且有變化幅度越到句尾越趨減緩的趨勢，通常這種趨勢發生在以長停頓分開的各語句之中（通常伴隨著標點符號），而在這個逐漸下降的趨勢之中，伴隨著局部高低起伏浪潮般起伏，而這些局部的變化由音節的聲調以及前後音節互相影響所造成，在這裡我們可以再回顧應證一下 3.4.1 所提到的前後文音節和語調的影響，概念上，我們可以將逐漸下降的趨勢視為語調（intonation）的變化，而聲調的變化影響著局部浪潮的變化，這兩種成分相加起來，就變成我們觀察到的基頻軌跡，在這裡要介紹一個著名基頻軌跡的疊加式(superposition)韻律模型「Fujisaki pitch model」 [6]，如圖十所示，這個模型將基頻軌跡分為兩個部份，一個為一開始稍微向上然後緩慢下降的 phrase component，另一個為局部性駝峰的 accent component，這個模型原本是用於日文語音韻律，後來也同時應用到不管是聲調語言（tonal language）或語調語言（intonation language）的各種語言 [1]，對於國語語音韻律，一般來說 phrase component 是對應到語調的變化，而 accent component 是對應到聲調的變化。



圖十: Fujisaki pitch model

另外，在這些以長停頓分開的各語句之中，常有部分明顯的 pitch reset 現象發生，即下一個音節的 pitch level 突然增高的現象，例如在圖九中的「十月份 一 到 二十日」，「到」和接下來的「二」雖然都是四聲，然而「二」的基頻軌跡高度明顯比「到」還要高了許多，這樣的 pitch reset 通常是因為語者想要加重某個字或詞而造成的，類似於 3.4.1 所提到停頓

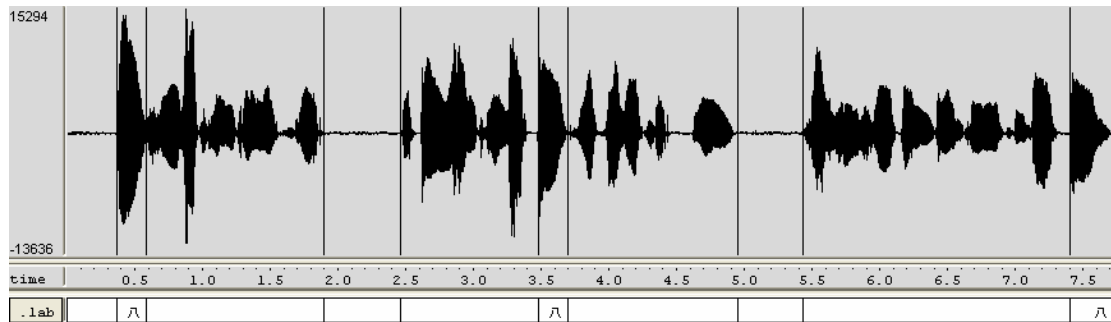
發生的地方，這些 pitch reset 的信號，大部分是對應到語法上的片語邊界、功能詞，而在詞內(intra-word)位置則鮮少發生。

3.4.4 音節長度

在連續語音中，通常語句的最後一個字或是詞的最後一個字，會比一般句中或詞中的音節長度拉長一些，如在圖九中的「行政院」、「統計」、「十月份」、「出口」，這就是所謂「lengthening effect」，通常這些拉長的詞尾代表的是類似停頓的功能，以拉長音節的方式代表訊息的單位的結束，但我們也可以發現到，並不是所有的詞尾都會發生lengthening effect，一般認為，相同於停頓以及基頻軌跡，lengthening effect較容易發生在語法片語、句尾或有標點符號的地方。

在3.2中提到不同音節依據本身的 phonetic structure 而有不同的音節長度特性，受到此一因素的影響，在圖九中的例子可能比較沒辦法很明顯的看到句尾 lengthening effect，因此我們舉另一個著名的句尾 lengthening effect 例子於圖十一，在這個例子為一個設計過的文章段落如下 [7]：

「八是一個常見的字，一般人常把八字掛在嘴邊，講話時動不動就提到八」



圖十一: Lengthening effect

三個「八」分別在句首、句中和句尾，我們可以明顯的看到句尾的八特別拉長，顯示出 lengthening effect。

3.4.5 能量軌跡

一般來說，在句首的音節通常有較高的能量，當語音發展到句尾以後，能量下降的趨勢變得很明顯，如圖九中的所示的能量軌跡，這是一個很自然的現象，因為語者在講一句話的時候，一開始有較足的氣，但講到句尾的時候漸漸地沒有足夠的氣，所以才會有能量軌跡逐漸下降的趨勢。然而上述的是一般的通性，語者有時候也會在他想要加重的音節上，以較大的能量來表達而造成局部的能量增強，如同其他上述之韻律參數，這些加重的位置通常發生在語法片語的開始音節。

4. 國語韻律的產生 — 文字轉語音技術

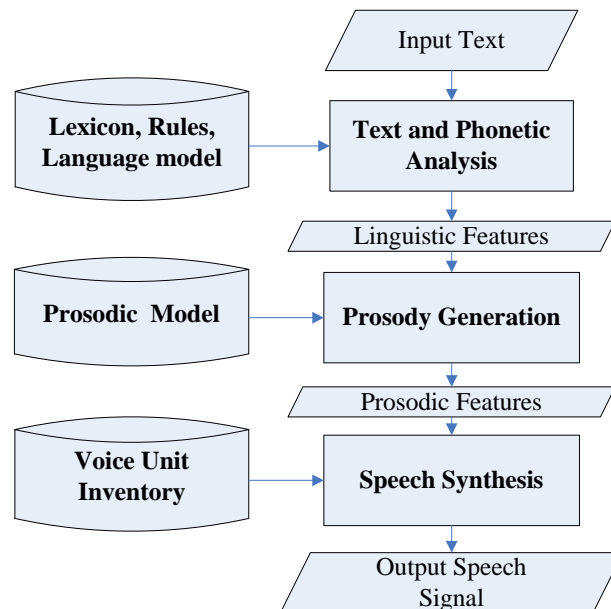
由前一個小節的討論，我們介紹了由音節單位以及較大以語句為單位的國語韻律特性，這些特性分別表現在語音韻律的四個面向上，而這些韻律的變化又和語言參數是有所關聯的，因此如果我們能夠建立語言參數和國語韻律的關係，也就能達成由文字轉成語音(或稱語音合成)的工作，在這一節我們來介紹一些目前為止成熟的語音合成應用技術之核心 — 國語韻律的產生。

4.1 國語文字轉語音系統簡介

圖十二是一張表示文字轉語音系統的工作流程圖，輸入語音合成系統的為文字，輸出為聲音信號，系統可以再分為（1）文字分析（text and phonetic analysis）（2）韻律產生（prosody generation）（3）語音合成（speech synthesis）三大部分。不同於英文，國語的書寫在詞和詞之間不會以空白表示，因此文字分析的部份負責將輸入的國語字進行斷詞的動作，一旦斷完詞便可得到一些語言參數（linguistic features），這些語言參數由低層次至高層次可以是：

- 音節層次：如音節讀音、聲調、音節於一個詞的位置
- 詞層次：如詞長、詞類
- 語法片語層次：如動詞片語、名詞片語
- 語法樹層次
- 語意層次

接下來將這些語言參數輸入至韻律產生器（prosody generation）以來產生韻律參數（prosodic features），也就是我們所介紹的語音韻律的四個面向，最後語音合成器依據韻律參數合成語音。



圖十二: 文字轉語音系統的系統方塊圖

一般來說，韻律產生的方式分為兩種，一種為規則法（rule-based approach） [2, 8-10]，另一種為資料驅動法（data-driven approach） [11-16]。

4.2 規則法之國語韻律產生

規則法（rule-based approach）利用一些前人觀察到的韻律現象，再加上一些簡單的統計結果而設計出一套韻律規則。這類的規則通常和語言層面的訊息有很大的關連性，必須先針對語句做詳盡的分析，才有可能套用適當的規則。一般常見的規則方法有：在各音節 pitch contour 的決定上套用 tone-sandhi rules，以及在考慮不同的 syllable、initial、final 及 tone 的影響之下，syllable duration 的各種分佈情形 [2]。另外，也可以分析句子中標點符號的形式（‘，’ ’。’ ’？’等）以決定在套用整體抑揚頓挫（intonation）在句子中該如何分配，

或是分析詞串的詞類(POS)、句法結構 (syntactic structure)、語意等，在特定的語法結構組合下套用特定的韻律規則。此類韻律模型的建構動輒產生數十個甚至上百個韻律規則，並且需要大量的語言學知識和複雜的人工處理過程始能產生，所以通常以這種方法所分析的語料對象沒有辦法非常龐大，僅能針對有限資料作分析，所得到的韻律模型有可能非常適應於某些語料，但是在模型的調適性上有很大的限制，無法依據不同的語料特性做適度的調整。

4.3 資料驅動法之國語韻律產生

資料驅動法 (data-driven approach) 則利用資料本身的參數，使用一些數學模型自動產生韻律模型，由資料本身的特性自動學習出韻律規則，建立語音和文字之間韻律參數的關連性。此種方法最大的好處為，我們不需要嚴格定義太多語言行為，也不必瞭解太深入的語言學知識，只要藉由適當的參數選取以及套用適合的數學模型，就可以由大量資料自動產生出效果不錯的韻律模型。

資料驅動法所使用到的數學模型，最主要為一些圖像識別 (pattern recognition) 用途的工具如決策樹 (decision tree) 以及類神經網路 (artificial neural network, ANN)，另外還有用於統計學上的線性迴歸法 (linear regression)，在此段落，我們舉一個已經成熟且用於商業軟體的韻律產生器來介紹 [14]。

如圖十三所示，基本上是用遞迴式類神經網路 (Recurrent Neural Network, RNN) 當作架構 (類神經網路的型態之一)，此 RNN 包括兩層隱藏層、輸入層及輸出層用以模擬人類韻律訊息產生的方式，輸出為音節之基頻軌跡、能量準位、子音長、母音長以及音節間的停頓等，輸入分為兩個部份，一為音節層次語言參數 (syllable level linguistic features)，另一個為詞層次語言參數 (word level linguistic features)，音節層次語言參數如下：

- $I(S_j)$ 、 $I(S_{j+1})$ 為現在音節 S_j 和下一個音節 S_{j+1} 之子音 (initial) 類型
- $F(S_j)$ 為現在音節 S_j 之母音 (final) 類型
- $T(S_j)$ 、 $T(S_{j+1})$ 為現在音節 S_j 和下一個音節 S_{j+1} 之聲調類型
- $L(S_j|W_i)$ 為現在音節在詞裡面的位置 (詞的開始、詞的中間、詞的結尾)

詞層次語言參數如下：

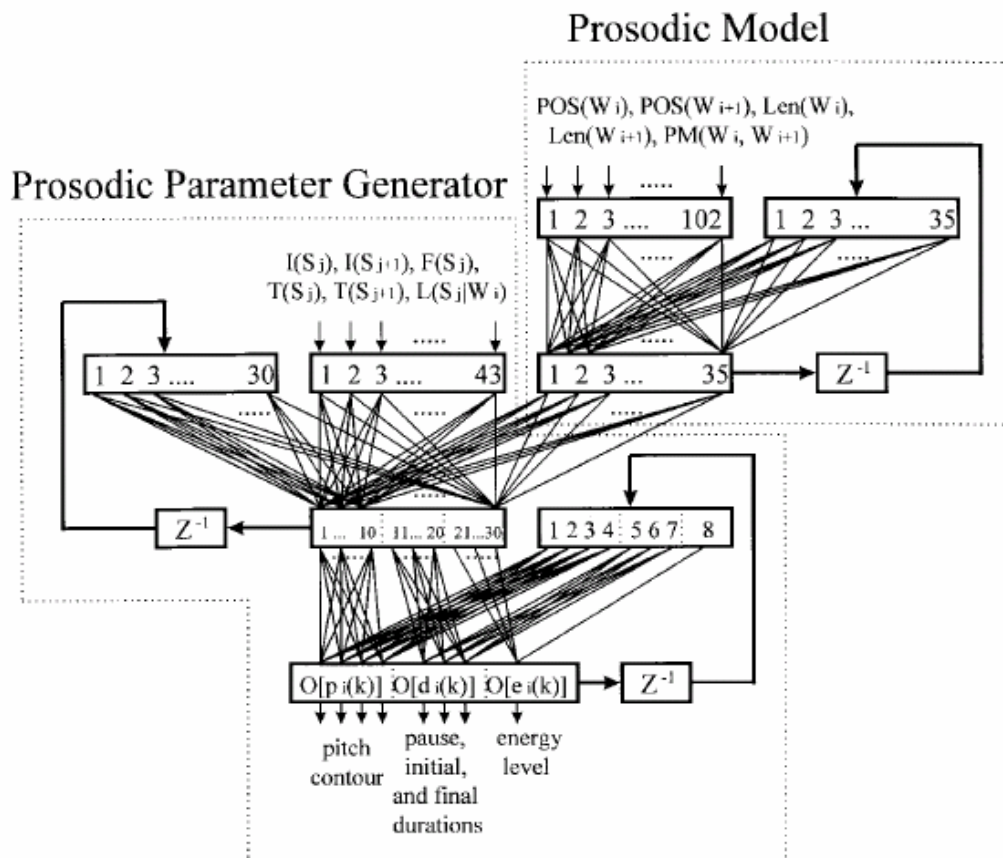
- $POS(W_i)$ 、 $POS(W_{i+1})$ 為現在這個詞 W_i 和下一個詞 W_{i+1} 的詞類
- $Len(W_i)$ 、 $Len(W_{i+1})$ 為現在這個詞 W_i 和下一個詞 W_{i+1} 的詞長
- $PM(W_i, W_{i+1})$ 為詞 W_i 和 W_{i+1} 中間的標點符號

此遞迴類神經網路 RNN 依其功能可區分為兩部份，一為韻律模式 (prosodic model)，第二部份為韻律產生器 (prosodic parameter generator)。

韻律模式主要由一個詞層次語言參數輸入層及第一隱藏層所組成，它藉由詞層次語言參數輸入進而分析與探索隱藏於文句中以及對應語音之韻律片語結構，此韻律模式以詞為單位輸出每個詞之韻律狀態 (prosodic state) 參數。

第二部份的韻律產生器，它是由一個音節層次參數輸入層、第二隱藏層以及輸出層所構成，它接受第一部份韻律模式所產生之韻律狀態與音節層次參數為輸入，並以音節為單位產生每個音節之韻律訊息。經過大量語料訓練，此遞迴類神經網路 RNN 可以自動學習並記憶人類說話之發音規則，如三聲接三聲之變調規則，進而產生適當之韻律，合成自然

流利語音，經由實驗結果數據中顯示，大部分合成之韻律訊息與原始之韻律訊息十分相近，由此可證明以類神經網路 RNN 以及韻律模式為基礎之韻律訊息產生器，效果很不錯。



圖十三：RNN 韻律產生器之架構

5. 國語韻律模型的建立

上述以 data driven 方式所得到的韻律產生模型雖然能夠得到不錯的韻律效果，但是經過這些韻律模型所預估的韻律參數和語言層次的關聯性非常模糊，無法很明確地分析出哪些語言參數對於韻律的產生造成何種影響。如果用比較直覺的概念說明，我們可以把此類韻律產生模型視為一個黑盒子，我們只能在輸入語言層次參數之後得到一組韻律參數輸出，但是黑盒子內部到底發生了什麼事情卻沒辦法由輸入的語言參數中得知。增加或減少一項語言參數，對於輸出的影響程度如何、影響哪一個韻律行為等，完全沒有辦法得到一個明顯的對應關係。這樣的模型也許適用於較簡單的語音合成系統，但是若要進一步考慮更高層的語法、語意與韻律之間的關係時，就變得非常不直覺且難以分析。

為了解決上述問題，有兩種主要方法開始被提出。第一種為「階層式國語語音韻律模型」，這個模型是以人工方式，以聽者的聽覺為準，將聽覺上感覺音節間有停頓的語音位置，予以停頓標記 (break labeling)，然後藉由這些停頓標記，來建構階層式的國語韻律模型。第二種為「Latent prosody modeling method」，此種方法假設韻律的形成受到某些潛在的 (latent) 影響因子影響，而這些影響因子可以藉由統計方式分離出來，並且可將各因子和語言參數做關聯性分析 [5, 17, 18]。

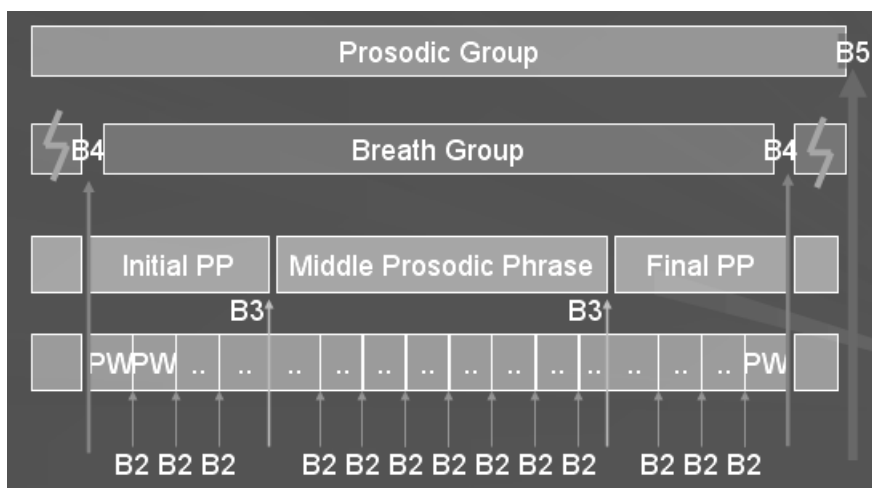
5.1 階層式國語語音韻律模型 [7, 19]

此種模型的建構，以人工標記的方式，對語料標記出聽者聽覺上的停頓位置，此標記方式是以 ToBI system [20]為依據，針對國語語料做 break labeling。當標示完語料後，依據階層性韻律結構 [7]建立模型，將階層性的語言影響因素扣除後，得到各種 typical patterns，包含 typical tone pattern、typical duration patterns 等。此種方法定義了六種 break types，其定義如表一所示。

表一：階層式國語語音韻律模型中的 break type 說明

Break type	Definition	Characteristics
B0	Reduced syllabic boundary	Syllable truncation often occurred in fast or informal fluent speech.
B1	Normal syllabic boundary	Usually with no identifiable pauses, but more of a psycholinguistic unit for native speakers.
B2	Minor phrase boundary	Perceived as a boundary where a slight tone of voice change usually follows.
B3	Major phrase boundary	A clearly perceived pause.
B4	Breath group boundary	Perceived end of exhale cycle followed by inhaling to begin another breathing cycle. It could be where a speech paragraph ends where trailing occurs with final lengthening coupled with weakening of speech sounds. But the speaker may still go on by breathing but not ending the speech paragraph.
B5	Prosodic group boundary	A complete speech paragraph ends by final lengthening coupled with weakening of speech sounds. The speaker makes a complete stop, take a new breath, and begin a new speech paragraph.

在標記完之後，我們由這些標記建構結層式的韻律結構如圖十四所示，其中 PW (prosodic word) 為以 B2 為分界的韻律單位，通常是由數個詞所構成，PP (prosodic phrase) 是以 B3 當作分界由數個 PW 構成，而 Breath group 是由幾個 PP (Initial PP、middle PP 以及 final PP) 組成，Prosodic group 是以 B5 當作分界由數個 breath group 構成的韻律單位。

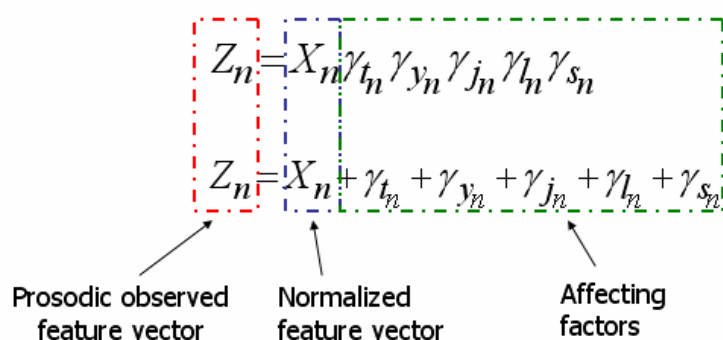


圖十四：階層式國語語音韻律模型

爲了找到這些韻律單位的 patterns，我們可以利用簡單的 linear model 來找到這些韻律單元的所對應的韻律參數，詳細的內容請參考 [7,19]。

5.2 Latent prosody modeling method [5, 17, 18, 21, 22]

階層式國語語音韻律模型的建構，是以聽者聽覺爲主的人工標記，然而此方法要求大量的人力，而且不同的標記者對於語料有不同的聽覺感受，進而造成語料標記的不一致，考量到節省大量人力以及希望能夠自動處理大量語音資料，「Latent prosody modeling method」孕育而生。此方法爲使用統計方式（Statistical method），考慮各種影響因素，這些影響因素可分爲「低層次的語言參數」以及「韻律狀態（prosodic state）」兩大類，低層次的語言參數包含聲調以及音節類型，而韻律狀態這項影響因素我們視爲一個 latent variable，代表著除了低層次語言參數以外的影響，如詞、語法片語、語意等等，而最後我們觀察到的韻律參數（prosodic observed feature vector）以這些影響因素（affecting factors）的相加（additive）或相乘（multiplicative）來表示，如圖十五中所示。



圖十五：Multiplicative and additive latent prosody models

我們可以發現到此方法中的低層次語言參數，直接影響的是 3.1~3.3 討論的音節固有的特性，再來，回想在 3.4 中討論到的大單元韻律的特性，除了低層次語言參數以外，還受到高層次語言參數或其他未知因素的影響，這三者造成了韻律的變化，然而到底這些因素本身各貢獻了多少影響給韻律的變化，是不知道的，而且同時我們也發現到高層語言參數對於韻律變化的貢獻也是較爲複雜，例如並不是所有的詞或是語法片語邊界都造成停頓，比較確定的是這些低層次語言參數對於大單元韻律的影響，這樣的方法試著用「韻律狀態」來吸收所有不確定的影響因素，接下來配合一些演算法（EM algorithm、CART 等）來得到低層次語言參數以及韻律狀態的影響參數，我們再進一步分析到底哪些影響因素造成了「韻律狀態」的變化。目前此種「Latent prosody modeling method」已應用在 prosody break labeling 上了，詳見 [5]。

6. 國語韻律模型未來研究方向

截至目前爲止，國語韻律的研究仍然以朗讀式的語音（read speech）爲主要研究對象，對於一般性的語音，如自發性語音（spontaneous speech）、情緒語音（emotional speech）等，則著墨甚少。主要原因爲國語語法規則並不非常明確，尤其是在越自然的語音下，規則性被破壞得越嚴重。譬如在流利語音中，Repair 和 Lexicon Cues 頻頻出現的現象 [23, 24]，常使得語句中的語法結構被破壞殆盡，使得語句分析變得異常困難。另外，在語流速度方面，自然語言節奏（tempo）的變化性又比 read speech 要來得大 [24]，造成韻律模型建構上的

不易。即使如此，國語韻律研究仍舊必須往自發性語音及情緒語音推進，畢竟韻律研究的最終目的終究是要處理真實世界裡人類的自然語言現象。

在國語自發性語音部份，未來研究方向有漸趨使用韻律訊息作為語音辨識的訊息來源之一，主要是利用停頓點、pitch jump、energy contour 等語音特性，搭配語言層次的分析，抽取隱藏在韻律背後的語法語意訊息，搭配現有的語音辨識技術以達到提昇辨識率的目的。

在國語情緒語音方面則大至可分為兩個主要方向，一個是情緒語音的辨識 [25-28]，另一個則是合成出情緒語言。目前所能達到的情緒語音辨識程度仍然不盡理想，僅能針對四到六種情緒做出偵測。主要原因為情緒語音語料庫收集不易，大部分的語料都來自語者刻意使用某種情緒讀出特定的講稿，因此很難收集到真正自然的語料。再加上情緒分類的模糊特性，我們通常很難界定一個句子到底含有多少百分比的高興、憤怒、哀傷等成份。這也造成情緒語音辨識的研究很難有很明確的評估方式，而只能以簡單的 MOS 等方式做評估。在國語情緒語音合成方面，目前的研究也只能做到以整句合成的觀點，將情緒成份完全套用進去。然而我們知道，情緒的表達並不見得是以句子為單位，有時候一些關鍵點的停頓、音量增強、音調提高等，就能將情緒完整表達出來，所以其實這方面的研究仍舊值得我們繼續努力。

另外，韻律研究也可以往認知的方向繼續下去，以瞭解人類為何可以輕易地偵測出各種韻律現象背後所隱含的情緒、動機、暗示等成份。這需要更多領域的專家們提供所學合力完成的。

7. 結語

韻律訊息的偵測與韻律模型的建構所牽涉到的學門非常廣泛，舉凡認知心理學、生理學、語言學、語音學到各種工程科學（控制、通訊、計算機科學）等，在各個領域都有深入的探討。以語音合成的觀點而言，以往對於韻律的掌握性非常不足，常造成合成出來的語句不甚自然等問題。另外在語音辨識方面，一直以來希望將韻律訊息加入語音辨識單元的需求越來越強烈，這使得韻律模型的建構在近幾年來越加重要。另外，對於各種應用層面諸如對話系統、語言學習系統等需要解析更上層語意的應用系統裡，韻律訊息的偵測與否將有機會成為理解語意的關鍵因素。

參考文獻

- [1] Fujisaki, H.: "Prosody, Information, and Modeling- with Emphasis on Tonal Features of Speech," *Keynote Paper at the Workshop on Spoken Language Processing*, January 10-11, 2003, Mumbai, India.
- [2] Lin-Shan Lee, Chiu-Yu Tseng and Ming Ouh-Young, "The Synthesis Rules in a Chinese Text-to-Speech System," *IEEE Trans. in Acoustic, Speech and Signal Processing*, Vol.37, NO.9, September 1989.
- [3] S.H. Chen, W.H. Lai, and Y.-R. Wang, "A New Duration Modeling Approach for Mandarin Speech," *IEEE Trans. On Speech and Audio Processing*, vol. 11, no. 4, 2003.
- [4] J.-S. Zhang and K. Hirose, "Tone Nucleus Modeling for Chinese Lexical Tone Recognition," *Speech Communication*, vol. 42, no. 4, pp. 447-466, 2004.
- [5] Chen-Yu Chiang, Hsiu-Min Yu, Yih-Ru Wang and Sin-Horng Chen, "An Automatic Prosody Labeling Method for Mandarin Speech," *Interspeech 2007* (accepted)
- [6] Fujisaki, H. and H. Sudo, "A Generative Model of Prosody of Connected Speech in Japanese," *Annual Report of Eng. Research Institute*, 1971, 30, pp. 75-80.
- [7] C. Tseng, S. Pin, and Y. Lee, "Speech Prosody: Issues, Approaches and Implications," in Fant, G., H. Fujisaki, Cao and Y. Xu Eds. *From Traditional Phonology for Mandarin Speech Processing*,

- Foreign Language Teaching and Research Process*, Beijing, China, 2004, pp. 417-438.
- [8] Ren-Hua Wang, Qinfeng Liu and Difei Tang, "A New Chinese Text-to-Speech System with High Naturalness," *Proc. of ICSLP96*, pp.1441-1444, USA 1996
- [9] Fu-chiang Chou, Chiu-yu Tseng and Lin-shan Lee, "Automatic Generation of Prosodic Structure for High Quality Mandarin Speech Synthesis," *4th Int. Conf. on Spoken Language Processing*, Philadelphia, PA, USA, Oct. 1996, pp. 1624-1627.
- [10] Ren-Hua Wang, Qinfeng Liu, Yu Hu, Bo Yin, and Xiaoru Wu, "KD2000 Chinese Text-to-Speech System," *ICMI2000*, LNCS 1948, pp.300-307, 2000.
- [11] Chung-Hsien Wu and Jau-Hung Chen, "Template-Driven Generation of Prosodic Information for Chinese Concatenative Synthesis," *ICASSP'99*, vol.1, pp.65-68, Phoenix, USA, 1999.
- [12] Ming-Shing Yu, Neng-Huang Pan, and Ming-Jer Wu, "A Intonation Prediction Model that can Outputs Real Pitch Pattern," *The Seventh Conference on Artificial Intelligence and Applications (TAAI 2002)*, pp. 784-788, 2002.
- [13] ShaoHuang Pin, Yehlin Lee, Yong-cheng Chen, Hsin-min Wang, and Chiu-yu Tseng, "A Mandarin TTS system with an Integrated Prosodic Model," *ISCSLP*, pp.169-172, 2004.
- [14] S. H. Chen, S. H. Hwang, and Y. R. Wang, "An RNN-based Prosodic Information Synthesizer for Mandarin Text-to-Speech," *IEEE Trans. Speech and Audio Processing*, Vol.6, No.3, pp.226-239, May 1998.
- [15] http://sovideo.iis.sinica.edu.tw/NeGSST/Index_PTT.htm
- [16] Sun Lu, Yu Hu, and Ren-Hua Wang, "Polynomial Regression Model for Duration Prediction in Mandarin," *Journal of Chinese Information Processing*, 2005.19(1).
- [17] Sin-Horng Chen, Wen-hsing Lai and Yih-Ru Wang, "A new duration modeling approach for Mandarin speech," *IEEE transaction on speech and audio processing*, vol. 11, no.4, Jul 2003, pp. 308-320.
- [18] Sin-Horng Chen, Wen-hsing Lai and Yih-Ru Wang, "A statistics-based pitch contour model for Mandarin speech," *J. Acoust. Soc. Am.* 117(2), Feb. 2005, pp. 908 – 925
- [19] Tseng, Chiu-yu, Pin, Shao-huang, Lee, Yeh-lin, Wang, Hsin-min and Chen, Yong-cheng, "Fluent speech prosody: framework and modeling," *Speech Communication*, Vol.46,issues 3-4,(July 2005), Special Issue on Quantitative Prosody Modelling for Natural Speech Description and Generation, 284-309.
- [20] Beckman, M.E. and G.M. Ayers, "Guidelines for ToBI Labeling," 1994, <http://www.ling.ohio-state.edu/phonetics/ToBI/main.html>.
- [21] Chen-Yu Chiang, Xiao-Dong Wang, Yuan-Fu Liao, Yih-Ru Wang, Sin-Horng Chen, Keikichi Hirose, "Latent prosody model of continuous Mandarin speech," *ICASSP 2007*
- [22] Chen-Yu Chiang, Yih-Ru Wang and Sin-Horng Chen, "On the Inter-syllable Coarticulation Effect of Pitch Modeling for Mandarin Speech," *Proc. of Interspeech 2005*, Lisboa, Portugal, pp. 3269-3272
- [23] Tseng, S.-C., "Repairs in Mandarin conversation," *Journal of Chinese Linguistics*, 34(1), (2006) 80-120.
- [24] Tseng, S.-C., "Linguistic Markings of Units in Spontaneous Mandarin," *ISCSLP 2006*, LNAI 4247, pp. 43-54, 2006.
- [25] Sheng Zhang, P.C. Ching, and Fanrang Kong, "Acoustic Analysis of Emotional Speech in Mandarin Chinese," *ISCSLP2006*, pp.57-66, 2006.
- [26] Tsand-Long Pao, etc, "Detecting Emotions in Mandarin Speech," *Computational Linguistics and Chinese Language Processing*, Vol. 10, No. 3, Sep. 2005, pp346-362.
- [27] Yuan J., Shen L., Chen F., "The acoustic realization of anger, fear, joy and sadness in Chinese," *Speech Prosody 2002*, France.
- [28] Chung-Hsien Wu and Ze-Jing Chuang, "Emotion Recognition Using IG-based Feature Compensation and Continuous Support Vector Machines," in *Proceedings of 2006 Speech Prosody conference*, Dresden, Germany, 2006.

