

語言模式在中文語音辨識上的應用

王榮宗 王駿發

e-mail address: wangjf@server2.iie.ncku.edu.tw

成功大學資訊工程研究所

摘要

目前在中文語音辨識後處理方面，偏向研究一個辨認音只允許給一個候選音且得輸入完整句子的情況。但這通常會面臨需投入人力很多且無法實用化的問題。

因此，在本篇論文中我們提出一些方法以解決語言後處理系統在實用化時所會遇到的一些問題。(1)使用者只輸入片段句子的問題。(2)在辨認系統中若一個辨認音給許多個候選音時，系統處理速度及複雜度的問題。(3)當辨認系統的準確度不夠時，如何讓系統具有容錯能力的問題。(4)如何避免分析複雜文法，但又能處理複雜的句子的問題。

在運用語言模式的策略方面，我們利用一個橫跨辨認與後處理系統的計分機構來配合統計式的語言知識，來讓系統具有容錯的能力。並利用複合詞文法和非詞之詞[註]的觀念以解決詞庫過度膨脹與使用者只輸入片段句子的問題。且我們運用句型與複合詞文法來增加系統處理文句的深度和廣度，而在處理的過程中，我們是先將句子轉成短語形式並配合物件導向的方法把詞庫資料與複合詞文法封裝一起，以減低系統處理的複雜度。

[註]非詞之詞：在文法上並不屬於詞，但在實用化時，語者常會輸入系統的「詞」。例如：「的現代化」、「近半年來」、「並於昨天」。但這種詞相當的多，故得先用人工篩選出較有可能的「詞」。而這些詞又可分為兩類，一類是較有規則的，如「近半年來」，另一是較無規則的，如「個小時了」。

一、簡介

目前在中文語音辨認研究上有正確率不符需求的問題。又由於中文的特性，同音異字相當多，即使辨認系統能得到完全正確的音串，仍無法唯一決定那些是其對應的正確字串。況且在辨認系統的正確率尚未達到相當準確的階段時，如何充分利用語言知識來提高正確率，是一項重要課題。

對於中文而言，當辨認的音串要轉換成對應的句子時，通常須要經過斷詞或構詞的程序，因為詞才是句法上及語意上的最小元素，而決定如何構詞或斷詞通常有3種方法：

<1>法則式法[1]---利用語言學的知識，歸納出一些規則，藉以判定斷詞或配詞是否合法，甚至整句句子的合法性。好處是用少量的規則，就能處理大量的資料。缺點是完整而周詳的規則不易建立，且無法處理太多的例外狀況。

<2>機率統計法[2][3]--利用電腦分析大量語料庫，獲得字詞結合的數學模式，來計算其正確性有多大。好處是電腦統計方便，統計完後無須再使用人工預建的詞典。缺點是當語料庫的規模不夠廣或大時統計的結果易失之偏頗，且正確率有瓶頸。

<3>樣本比對法--利用人工收集的詞典用比對搜尋的方法，找出所有可能的組合，再利用給分規則或轉成Dynamic Programming [4]的最佳化決定。好處是規則簡易。缺點是因詞庫的容量限制無法將所有詞彙收集完整，而影響其正確性，且有組合爆炸的情形。

由於人類理解的方式是靠「物件的記憶」、「聯想」與「推論」，物件的記憶類似樣本比對、聯想類似機率統計法、而推論則類似法則式法。

所以在本論文中結合了樣本比對（詞庫收尋）、法則式法（複合詞文法，句型）及機率統計法(n-gram)來將正確率提昇。

而語言模式後處理系統在實用化時通常會遇到下列問題：

(1)辨認系統造成的問題：

由於辨認系統的準確度不夠，而造成候選音過多的狀況，甚至沒有正確答案在其中的問題。

(2)語者輸入所產生的問題：

由於句子太長或其他因素，語者只輸入片段句子，甚至不合文法的句子。

(3)詞庫完整性和新詞產生的問題：

中文詞任意組合都可能成為另一新詞，其語意、詞屬性，可能與原先的詞彙不同。且無一完整的詞庫，又無制訂標準，以致無法把所有詞彙蒐集完全。因此產生如何利用既有的詞庫，來解決所有狀況的問題。

(4)文法、句型規則制訂的問題：

中文語文資料一向欠缺較系統化的整理，而中文句型又千變萬化，再加上中文詞性變化多、屬性分類複雜，以致造成文法、句型在制訂時無法兼顧所有狀況的情形。

(5)語意分析及前後文的問題：

在中文特性中，具有相同音串，但有不同詞屬性的詞很多，這造成一個音串會有相當多的候選詞產生，導致了有些句子甚至得考慮前後文的關係才能決定出正確答案。

(6)處理句子深度及廣度的問題：

在實用化的考量下，系統應定位在能處理「任意中文句」的架構下，且需考慮系統處理句子的深度及廣度問題。

二、語音辨認與語言模式之整合架構

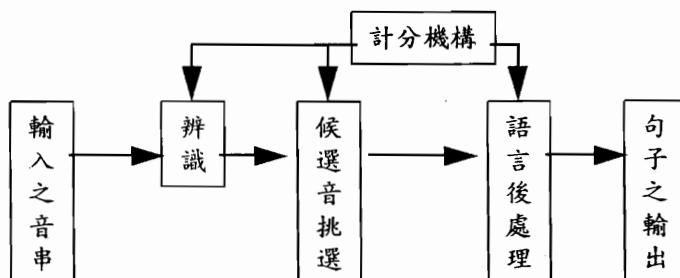
我們是針對語音辨識後處理，故應使系統具有下列特性。

可能因為句子太長，語者只輸入片段句子。因此語者輸入的音串，不論句子是否合乎文法，或只有部分句子輸入都得視為正確，但若是有完整的句子在裡面，須優先出來。

特性 1：系統能處正確的句子以外，還能處理短句或非詞之詞。

語音辨認時，就須與語言模式結合一起，共同挑選出最可能的候選音，以避免造成語音辨認結果主導後處理系統的現象，兩者需相輔相成。

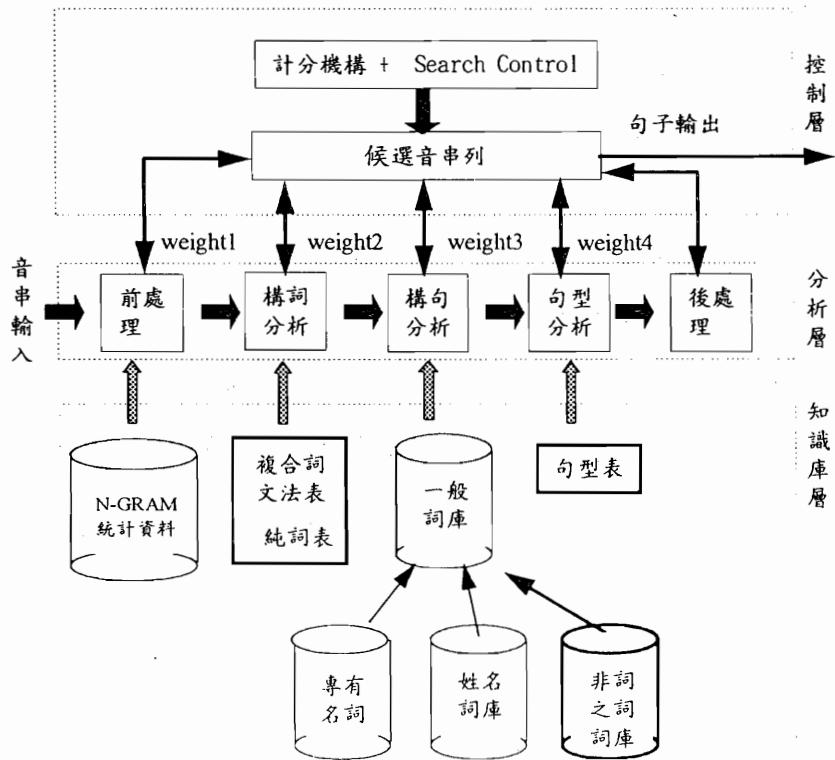
特性 2：具有一個計分機構，橫跨辨識與後處理部分，以結合低層的辨認與高層的語言知識。



不同的辨認方法須有不同的計分機構，使得辨認系統能與後處理系統緊密結合一起。

特性3：計分機構給分的判定公式與辨認給分公式相結合，使系統能整合一起且及時處理(realtime)。

其整合的架構如下圖，在語音串輸入後經由辨認器與計分機構共同挑選出一些候選音與可能的音韻，先利用純詞表將所有只有一種字的音先配上字，再將其結果送至構詞分析器，以節省詞庫收尋時間。再用搜尋比對的方法從詞庫中找出所有相對此音串的詞，並由「常用單字詞表」與「優先順序表」[5][6]將詞排定其優先順序並給其特徵分數。接著經由複合詞文法的分析，將可以合併的詞，再做一次合併，並付予新的詞屬性及特徵分數，使每個合併的詞轉成短語形式而使整句句子的複雜度降低。最後由計分機構從所有短語型態的詞串中，挑選出所有的組合去構句，並送至詞性分析器利用句型比對給分，而算出整句的特徵分數，最後由系統依據總體的特徵分數列出一些最可能的候選句。



在此為了克服中文句子組合複雜的特性，造成文法或句型難以制定，故先利用較簡易且短的複合詞文法把複雜度高的句子簡化成複雜度低的句

子，使整句句子變成短語的組合，並利用物件導向的方法將資料封裝，以減低系統處理的困難，如此一來句型或文法的規則只須用最簡的型式表示即可，且將來文法或句型要新增或修改也較容易。

<例>我是一個快快樂樂的中國人

從組合的詞組中有一如下的組合：

我 是 一 個 快 快 樂 樂 的 中 國 人

相對應的詞性如下：

我→Nhaa(常用的人稱代名詞)

是→Dbb(表說話者的評斷的副詞)

一→Nea,Dd(定詞，時間副詞)

個→Nfa(個體量詞)

快→VH13,Dd(能夠後接比較對象及兩者差額的述詞，時間副詞)

快樂→VH21(非作格動詞)

樂→VH21(非作格動詞)

的→De,Ta(副詞)

中國→Nca(專有地方名詞)

人→Nab(個體名詞)

我們希望它經過複合詞文法合併後能變成如下的組合。

我 是 一 個 快 快 樂 樂 的 中 國 人

且詞性更改為系統的屬性：

我→nn(名詞)

是→V2(動詞)

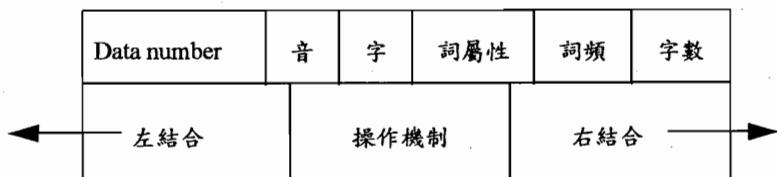
一個→cl(量詞)

快快樂樂→av(形容性述詞)

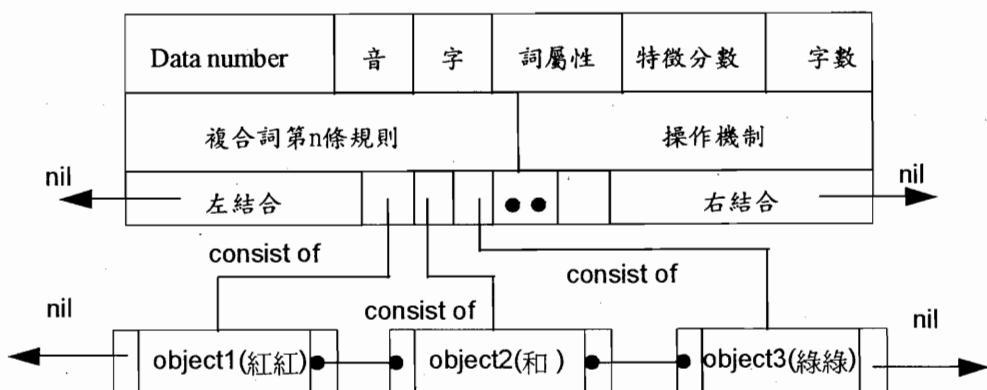
的→de

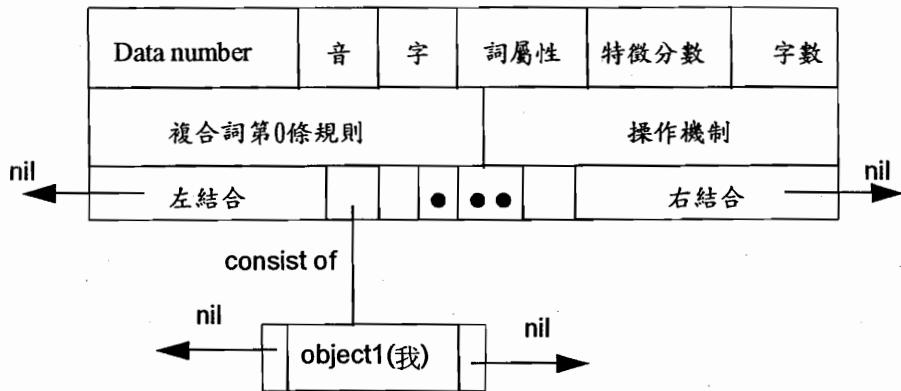
中國人→nn(名詞)

但由於中文詞彙的組成現象相當活躍，任何組成都可能成為新詞，若不用一些合成的規則先加以分析，會造成不斷遭遇新詞的困擾（像此詞沒存在詞庫，但可以用合成的方法產生）。且由於新詞的產生往往會造成詞性變形而導致文法或句型分析上的困難，故與其在文法及句型上使用繁瑣的規則，倒不如將這種問題利用物件導向封裝資料的方法並配合複合詞文法來解決。因此在資料單元的結構上，除了資料本身，尚須包含詞的合成規則。為了將模組的外部行為（句型詞性的比對）與模組內部的建構（詞的合併與詞屬性的更新）分隔開來，故採取物件導向的知識表示法。因此我們將每個詞(object)的資料結構表示為



在此我們把每一條複合詞的文法建成類別(class)（為了一些無法與別人結合的詞也能包含進去，故特增加一條複合詞文法但其內容為empty）。如此一來，我們就可以利用複合詞文法來彼此傳遞訊息，而把詞(object)再合併。而不能再與別人合併的最小單位我們定義為一個block或稱為過渡之詞。而每個Class的結構表示為





對於每個block來說，可以將其視為具有某項特定詞屬性的黑盒子而不必考慮block裏頭的資料及實際內部的運作情形，且文法或句型的分析主要就是詞屬性間關係的分析，並不須要知道內部實際的資料，故此種資料結構會減低系統處理的複雜度。且可將過渡之詞對外所呈現的行為與觀念和內部的建構細節分離開來，讓外界可以用一種較高的抽象式觀點來看待此物件。

三、語言模式技術

在以往的系統上，語言處理模式系統總是毫無抗拒的接受辨認器所給的候選音，這常會造成辨認器主導整個系統的好壞。若辨認器不好的話，就算語言處理模式再好，通常也回天乏術。所以當辨認器在選擇候選音時就須把語言知識的資訊給考慮進去，因此辨認候選音時的給分函數應為 $G(x)=weight_1 * \text{辨認給分} + weight_2 * \text{語言統計給分}$ 。

辨認器要挑選出一個候選音時，必須考慮它與前幾個候選音的關係，而算此候選音出現的機率，即為n-gram，但為了實用性只採用tri-gram。

設W為詞串列 $W = W_1 W_2 W_3 \dots W_n$

W的機率為 $P(W) = \prod_{i=1}^n P(W_i | W_1 W_2 \dots W_{i-1})$

化為trigram的形式 $P(W) = \prod_{i=1}^n P(W_i|W_{i-2}W_{i-1})$

而其中 $P(W_i|W_{i-2}W_{i-1})$ 的機率估算由下式來計算

$$P(W_i|W_{i-2}W_{i-1}) = f(W_i|W_{i-2}W_{i-1}) = \frac{C(W_{i-2}W_{i-1}W_i)}{C(W_{i-2}W_{i-1})}$$

上式中的函數 C 代表其參數字串曾被訓練過的次數。

但由於統計的關係會造成某些音串是正確的，但卻不會出現在訓練資料中而造成字串的機率 $P(W)$ 變為零，故再將trigram的公式作平滑化處理，取trigram、bigram、unigram的內插而得到相對機率。

$$P(W_i|W_{i-2}W_{i-1}) = q_3 * f(W_i|W_{i-2}W_{i-1}) + q_2 * f(W_i|W_{i-1}) + q_1 * f(W_i)$$

$$\text{其中 } q_3 + q_2 + q_1 = 1$$

但辨認器每個音會給許多候選音，故須算與所有候選音的關係

$$P(W_i|\sum_{n=1}^N W_{i-2n}W_{i-1n}) = q_3 * f(W_i|\sum_{n=1}^N W_{i-2n}W_{i-1n}) + q_2 * f(W_i|\sum_{n=1}^N W_{i-1n}) + q_1 * f(W_i)$$

N為候選音的個數

若是要求更精確的利用語言學知識，著實應該在辨認公式上加以修改，但這又隨著不同的辨認方法而有不同的作法。而在辨認時就考慮語言知識還有二項重要的目的，由於辨認器的關係可能候選音中沒有正確的答案，這會造成後處理相當大的麻煩，除非後處理系統有容錯的能力，否則通常很難解決此類的問題。但我們可以利用加入統計的語言學知識(如 trigram)盡可能的將這種問題避免掉。另一目的是當使用者輸入非詞之詞時，由於非詞之詞並不合乎句型，甚至複合詞文法，所以沒有任何的規則可遵循，而非詞之詞本身就是統計的資料，因此假如在辨認時就加入統計的語言知識這對非詞之詞的辨認率會有明顯的幫助。

而為了簡化自然語言處理系統的運算步驟及能利用更多的語言知識，音串必須轉成詞串，而詞串必須轉成短語型式。為了將詞合併變成短語型式，我們將複合詞的定義予以延伸擴大使之包含了複合名詞、地方副詞、時間副詞、定-量式複合詞、名-方式複合詞、數量詞、方位詞...等。

經由觀察的結果，發現大部份的非詞之詞是由介詞、助詞及少量副詞所組成，故利用複合詞文法可解決大部分片段句子輸入所造成的問題。而在複合詞合併的過程當中，計分機構也會依其特徵減其特徵分數，而所有詞合併完畢後就可以進入構句的步驟了。構句的主要動作是從音串的第一個音開始到最後一個音按照各音開頭之詞群已排定的順序，一一予以組合。由於每個過渡之詞都有其特徵分數（特徵分數愈小愈好），因此我們把構句問題轉成從sequence找出一條最小score的路徑，且字數符合需求。經過了構句的步驟後，其正確的句子幾乎在所列的候選句中前200名內了，但由於系統只列出前10名，故須靠句型比對將其正確的句子排名提昇至前10名內。

而目前的剖析器常採取Top-Down(由上而下)或Down-Up(自下而上)的方式進行，其方法通常是在某一範圍內使用某一條規則來判斷。但這些規則都是區域性的(Local)，因此很難顧及其整體的正確性，因此常發生就區域而言可能是正確的，但當這些區域接合之後就會有不太正確的情形產生。且剖析處理常耗費不少時間，又有時加入新的文法會發生與原有規則衝突的狀況。而複合詞文法本身就是一些短句的文法，況且句子到此已經變成短句型式的組合了，其彼此之間也只是詞性和詞性之間的關係，故我們捨棄文法的剖析(parser)，改用句型比對來配合計分機構給分，使其正確句子排名提前，且若將來有新的句型結構出現時，或使用者自定句型結構，通通只要加入句型表即可，如此可避免修改文法的困難。

每個詞(object)初始時系統都會給予其相同的基本分數，而計分機構主要的工作就是在每個步驟中判斷「詞」是否符合某些特性，若是則計分機構會減其特徵分數，故物體的特徵分數愈低，其排名就愈前面，最後將各個步驟的特徵分數乘上不同的比重值，再從所有組合中挑選出前幾名的句子當做輸出。將各步驟所得的分數乘上比重值的主要原因是為了避免當片段句子輸入時會因通不過句型的比對，而導致沒有結果輸出的情形。況且各部份的處理對於正確句子的形成影響不同，故須給不同的比重。

我們現就對各步驟的減分標準逐一說明：

(1)前處理：依其tri-gram的統計機率減分，機率越大減分越多。

(2)構詞分析：減分標準是依詞數，若此詞的字數愈多其減分就愈多，然後再比較是否合乎優先順序表[5][6]，若符合則再減分。

(3)構句分析：減分標準是看組成的句子中是否包含有下列的詞，若有則每項各再減分。

1. 「是」
2. 「有」
3. 數詞，如「三」、「十」
4. 定詞，如「這」、「那」
5. 代名詞，如「我」、「他們」
6. 專有名詞，如「地名」、「人名」

這目的是要避免包含不正確長詞的句子太多時，會造成正確的句子無法進入排名的情形。其依據主要是參考[5][6]。

(4)句型分析：

若此句所有的過渡之詞詞性都符合句型中的詞性則特徵分數減分。

而後處理這部分只做2個動作。當候選句中有相同的句子時，刪掉其排名較後面的句子，因為從不同的組合方式或由於具有不同詞性的關係，但其最後結果有可能相同，為了使別人還有機會進入前幾名，因此刪掉相同的句子。若無任何候選句輸出時，則將每個音用最常用字代替而輸出。

四、實驗結果與討論

實驗設定環境

1. 辨認系統：採用台康公司的語音辨認卡，單音輸入，每個音給9個候選音。

2. 實驗資料庫：

(a) 測試語料庫：十二冊國小國語課本課文，從中選取第7冊當測試樣本。

國立編譯館[Elementary Chinese]兩冊中所有的例句，約8000句測試句。

(b) n-gram的統計資料：統計1000個非詞之詞的trigram資料，而從隨機選取50句當測試資料。

因為後處理部分與文章的內容相當有關，又沒有一標準的測試檔案，所以我們乃隨機選取國小課本，與"Elementary Chinese"中的例句當做測試樣本。選國小課本當測試樣本的原因是裡面有包含新詩，而"Elementary Chinese"中的例句包羅萬象涵蓋了大部份的用語。而系統中用的句型也是"Elementary Chinese"中的句型，共有206條。因在中文裡有很多句子其混淆性若不用語意分析或利用前後文的關係根本無法決定其正確性，如「向警方開槍」、「像警方能幹」，「甚麼叫做容易」、「甚麼叫做溶液」，但語意分析耗費人力頗大，而整體預期的效果也不可知，因此系統尚未加入語意分析，而改取前10名的句子當做輸出。

實驗一：辨認時未加入語言知識(tri-gram)，完整句子輸入(此正確率是指整句句子全對)，分別統計國小課本與國立編譯館例句的正確率：

正確率	累計正確率 (國立編譯館)
Top 1	71.3%
Top 2	76.2%
Top 3	81.6%
Top 4	83.7%
Top 5	84.3%
Top 6	84.9%
Top 7	86.1%
Top 8	88.6%
Top 9	89.8%
Top 10	90.1%

正確率	累計正確率 (國小課本)
Top 1	71.3%
Top 2	79.2%
Top 3	84.6%
Top 4	89.7%
Top 5	91.3%
Top 6	93.6%
Top 7	95.8%
Top 8	96.1%
Top 9	96.9%
Top 10	97.9%

雖然句型規則是根據國立編譯館所制定的，但在實驗中卻發現國立編譯館的正確率反較國小課本低！其原因是字與詞的關係，因國立編譯館的句型包羅萬象，故其例子所用的詞有很多在我們的詞庫中並沒有儲存的。
(例)你喜歡看中國片還是外國片？

由於“中國片”詞庫沒有，且複合詞文法也無法組合出，“片”也沒存在詞庫)故造成錯誤。而從實驗中發現大部分的錯誤都是由於系統未儲存此詞彙，或是因複合詞無法組合此詞彙造成的，這通常把此詞彙加入後即可解決。

實驗二：隨機選取50句非詞之詞當測試樣本，辨認給分時未加入trigram統計資料，這50句非詞之詞中，2字詞佔20個，3字詞佔20個，4字詞佔10個，以下是部份例句的辨識情形([]表示正確答案，()表示候選句的名次)

[有減少]

- (0) 有減少
- (1) 有剪草
- (2) 有鮮少
- (3) 眼臉少
- (4) 眼臉掃
- (5) 眼臉草
- (6) 眼臉少
- (7) 眼臉早
- (8) 眼臉吵
- (9) 眼臉找

[查一下]

- (0) 長一下
- (1) 長於嚇
- (2) 長於夏
- (3) 長於架
- (4) 長於嫁
- (5) 長於駕
- (6) 長於向
- (7) 長於象
- (8) 長於謝
- (9) 長於穴

[九月底]

- (0) 九玉里
- (1) 九月底
- (2) 九樂里
- (3) 九錄野
- (4) 九玉體
- (5) 小月比
- (6) 小月筆
- (7) 小月起
- (8) 小月體
- (9) 小月體

〔五 輛〕

- (0) 我 另
- (1) 傷 羨
- (2) 火 烈
- (3) 五 列
- (4) 五 輛
- (5) 五 令
- (6) 武 令
- (7) 我 烈
- (8) 我 滅
- (9) 舞 步

〔直 到〕

- (0) 時 報
- (1) 碩 大
- (2) 食 道
- (3) 遲 到
- (4) 直 到
- (5) 直 道
- (6) 執 傲
- (7) 雌 豹
- (8) 衡 道
- (9) 十 道

〔你 也〕

- (0) 你 野
- (1) 你 演
- (2) 你 瘋
- (3) 你 影
- (4) 你 引
- (5) 你 舐
- (6) 你 臉
- (7) 米 野
- (8) 米 演
- (9) 米 瘋

正確率	累計正確率
Top 1	60.1%
Top 2	63.8%
Top 3	67.2%
Top 4	71.7%
Top 5	73.3%
Top 6	78.9%
Top 7	80.1%
Top 8	82.5%
Top 9	83.6%
Top 10	83.9%

實驗三：與實驗二相同的50測試句，但辨認給分時有加入tri-gram統計資料

國小課本	累計正確率
Top 1	63.8%
Top 2	68.4%
Top 3	69.4%
Top 4	72.5%
Top 5	76.9%
Top 6	81.2%
Top 7	84.7%
Top 8	86.9%
Top 9	87.1%
Top 10	88.9%

由實驗觀察中發現，非詞之詞的正確率與辨認時單音的排名有相當大的關係在，若辨認的單音在候選音串中排名越前面，則其正確率越高，原因是非詞之詞幾乎是把候選音所有的組合配出，然後到詞庫搜尋出來。所

以候選音排名越前面越可能對，故在辨認時加入tri-gram或bi-gram對其正確率有相當大的幫助，且落入前五名的機會也大大增加。

實驗四：我們隨機選取一些較特別的句子與國音輸入法做個比較，由於沒有測試大量資料，故此部份僅供參考（[]表正確句子）

[山崩地裂我不怕 頂天立地站的牢...（國小課本—新詩）]

（國音）山崩地裂我不怕 頂天立地佔的勞

（系統）山崩地裂我不怕 頂天立地站的牢（1）

[籃下 咱一球莫展...（報紙）]

（國音）藍下 咱一球莫展

（系統）攏下 咱一球莫展（1）

[全自動滾筒電鍍工廠誠徵...（廣告）]

（國音）全自動滾筒電鍍工廠成爭

（系統）全自動滾筒電鍍工廠誠徵（1）

[晴天的映照把你染藍了...（新詩）]

（國音）晴天的應召把你染藍了

（系統）晴天的映照把你染藍了（1）

[一上來即...（非詞之詞）]

（國音）一上來及

（系統）一上來急（1）

[你喜歡中國片還是外國片...（日常用語）]

（國音）你喜歡中國片還是外國片

（系統）你喜歡中國騙還是外國片（1）

[子音與母音間不穩定區段的資訊...(論文部份擷取)]

(國音)子因與母因肩部穩定區斷的資訊

(系統)子音與母音兼不穩定區段的資訊(6)

[計算其週期變異數...(論文部份擷取)]

(國音)計算其週期的便藝術

(系統)計算其週期的辯異數(4)

[接近口語化...(短語)]

(國音)接進口語化

(系統)接近口語化(1)

[的定義是...(非詞之詞)]

(國音)的定亦是

(系統)的定義是(1)

[在上式中我們用到...(論文擷取)]

(國音)在上市中我們用到

(系統)戴上是鐘我們用到(1)

ps. 系統掛號中的數字是表示此句是第n名出現

實驗五：各處理步驟中weight比重的實驗，及減分多寡的實驗：

由於此部分應該實驗大量的測試資料才能決定出最好的比重及減分多寡，但限於人力及無一標準的測試庫，所以我們用句型中的例句當測試資料。而wieght 的測試值也只分佈一小範圍，但由於效果不錯故就採用weight1=0.3，weight2=0.25，weight3=0.3，weight4=0.15，而減分是每符合條件減一分。

但構詞給分會不會與構句給分相衝突呢？因構詞給分是依長詞優先，所以詞中包含的詞屬性一定比較少，但構句給分卻是屬性越多越有利。

(例)我 是 中國人 ... (a)

我 是 中國人 ... (b)

在構詞分析時(a)由於長詞優先的法則故較占優勢，而在構句分析時(b)由於屬性較多所以較占優勢。但因為句型規則是以最簡型式表示，所以不會有 $s+v+n$ 與 $s+v+n+n$ 的狀況，且經過過渡之詞的合併 $n+n$ 會變成 n ，故仍是(a)優先順序較高。那會不會有長詞中可以分成兩種以上的不同詞性的詞呢？在我們實驗時並未遇到此情形，但為了解決這問題我們對於構詞分析與構句分析給不同的比重值，也就是說當此種狀況發生時，詞屬性較多者較佔優勢。

第五章 總結

在本論文中，我們嘗試性的提出一些方法，以解決語言後處理在實用化時所會遭遇到的問題，而目前本系統只粗具規模，仍有些瓶頸尚未突破，未來的系統仍可做以下的改進：

(1) 對於詞彙加以更進一步的研究，包括非詞之詞的分析，並結合複合詞，以克服片段句子輸入及詞庫膨脹的問題。

(2) 進一步分析複合詞文法與句型規則，使其更有規則，並擴大其應用範圍。

(3) 充分結合短句模式與句型，以減少系統的複雜度，並使系統隨時具有擴充性。

(4)利用計分機構以避免語音評分主導或是語言評分主導，並利用統計的語言知識以克服候選音中無正確音的問題，使系統具有容錯能力。

(5)研究特徵分數的給分方式及各步驟的比重加權，並用大量語料庫做測試。

參考文獻

[1]Hsieh , M.L , T.T.Lo and C.H Lin "Grammatical Approach to Converting phonetic Symbols into Characters" Proceeding of National Computer Symposium , Taipei , 1989 , 453-461 。

[2]C.K.Fan and W.H.Tsai , "Automatic word identification in chinese sentences by the relaxation technique" , Proc. of National Computer Symposium , 1987 , pp.423-431

[3]張俊盛、陳志遠、陳舜德，"限制式滿足及機率最佳化的中文斷詞方法"中華民國第四屆計算語言學會研討會論文集,147-165

[4]Sproat.R. "An Application of Statistical Optimization with Dynamic Programming to Phonemic-Input-to-Character Conversion for Chinese" , Proceedings of ROCLING 3 , 1990 , 379-390

[5]A language Processing Model for Mandarin Speech Recognition Using Score Unit and Grammar-Based Markov Language Model—Chung-Hsien Wu,Jhing-Fa Wang,ROCLING ,1994

[6] 謝子陵 - "A Linguistic Decoder for Mandarin Speech Recognition Using a Scoring Parser", 碩士論文, 國立成功大學, 1993

系統優先順序表：

1. 以常用單字詞為開頭的多字詞。
2. 常用單字詞。
3. 非以常用單字詞為開頭的多字詞。
4. 普通單字詞。

系統句形規則

由於句型太多，故只舉例說明之：

其主要根據是從國立編譯館"Elementary Chinese"上下冊(正中書局)，中所列的句型再修改。

句型23：

名詞 複詞 動詞 了 嗎？

你們	已經	懂	了	嗎
學生	已經	走	了	嗎
他們	已經	說	了	嗎
老師		講	了	嗎
你朋友		問	了	嗎

句型19：

主詞 動詞 數量詞 名詞

我 要買 三塊 錢的茶葉
外國學生 一萬塊 錢的紅蛋
他父親 三千三百元 的古董

系統複合詞規則

複合詞中用到的詞類，其代號及意義說明如下：

ado：前置副詞

adv：一般副詞

av：形容性數詞

b：後綴詞

cl：量詞

cj：連接詞

de1：的

de2：得

dt：定詞

loc：位置詞

n：名詞

q：數詞

v1：不及物動詞

v2：可及物或不及物動詞

v3：及物動詞

$N_1(n)^*(loc/b)$

$N_2(av)^* de1(n)^*$

$N_3 N_* (cj/de1 \quad N^*)^*$

$N_4 dt(cl_*) N_*$

$AV: ((av)^* cj(av)^*)^*$

$CL_1: (q)^*((av)cl(cl))$

$CL_2: ((q)^*(cl))^*$

$CL_3: ((q)(av)cl(cl))$

$Adj: (adj)^*(cj)(adj)^*$

$ADV_1: adv(adv)(del)$

$ADV_2: av \ del$

$V_1: (ado)^*(ADV_1)_{v1}$

$V_2: (ado)^*(ADV_1)_{v2}$

$V_3: (ado)^*(ADV_1)_{v3}$

$ASP: (asp / de2)(asp)$

符號說明：

1. 括弧(): 有括弧者表示「可以不存在」，無括弧者表示「必須存在」。
2. 星號*：有星號可以出現多個。
3. 十字+：某條規則裡有任一含十字者存在就能成立，如果有多個含十字者存在，則應按順序出現。
4. 斜線/：以斜線串連的各項只能有其中一個存在。
5. *：代上式中任一同詞性的複合詞文法。