

# 中文字根之分析

謝 清 俊 (交通大學工學院)

黃 永 文 (中原理工學院電子工程系)

林 樹 (交通大學工學院)

(1972年12月16日收到)

**摘要**——本文討論之主要內容，為分析中文字形結構之方法，其目的在減省計算機輸出中文字形所需之字形儲存空間。此法係將中文視為由字根拼成之二度空間字形，進而提取其共同點，減少重複，並加系統之整理，以求得一類似語文模式 (Linguistic Model) 之結構。經統計結果，此模式可產生 48,713 個已有之中國字，若以九千通用字計，可將現用之字形儲存空間減少至二十分之一以下，且對於中文輸入計算機，提供一種新方法。

## I. 基本用字研究

我國文字，自殷商時代，已有甲骨文約三千字，至漢說文解字書，收九千餘字，清康熙字典，集字四萬二千餘，民國五十七年，印中文大辭典，得四萬九千餘字。然於茲四萬數千字中，今日能靈活運用者，祇占極少數，餘者可謂死字或垂死之字。古今重要典籍，其使用單字，罕有逾三千者，如易經，有單字 1,595 字，<sup>(1)</sup>三民主義為十六萬餘字鉅著，僅有 2,134 單字。<sup>(2)</sup>至於時下常用字數，亦三、四千字足矣。聯合報中文自動鑄排機，採常用字彙 2,376 字，羊汝德先生編訂新聞常用字彙三千字。然為應社會上各方面之需要，計算機所需字彙，應由現代社會通用文字中，盡量擴大蒐集範圍。因此，本文所述研究工作，採用中文電腦基本用字研究一書，列 9,129 字，略歸為四類，計收「最常用字」1,857 字，「次常用字」2,068 字，「間用字」2,182 字，「罕用字」2,425 字，「異體字」597 字，已較一般常用字彙超出甚多，足供應用。

再者，此書中字彙及其出現頻率係集蔡樂生先生著常用字選等十一種資料統計而成，統計中，共取字樣 (Samples) 數為 2,022,604，其中最常用字占全部字樣 97.34%，次常用字占 2.27%，間用字占 0.27%，罕用字占 0.12%。各字均附有其在全部分字樣中之出現次數。可供中文字根分析之參考。<sup>(3)</sup>

## II. 以往之字根研究

二十世紀以還，由於計算機科學之飛躍進展，自第一部計算機 ENIAC 問世，迄今雖祇二十餘年之歷史，而其應用範圍日益廣大，應用項目日漸增多，儼然成為歐美各國促進科學、經濟、文化發展之重要工具；然而我國由於文字之字數太多，字形複雜，以致難於利用計算機處理中文資料。近年來，國內外不少專家學者，曾致力於中文資料處理自動化之研究，嘗試中文輸入、輸出之各種方法：或以一字一碼，或取字之部份筆畫為碼，或以字音為碼……等，然都難以推廣應用，究其原因，均礙於所需儲存空間太大，機器製作困難，操作費時費事。為謀此等問題之解決，遂有中文字根之研究。此法係將每一中國字用其所含之字根 (Radicals) 以及指示字根關係位置之定位符號 (Operators) 表示，作為中文資料輸入輸出之藍本。茲將其較重要者，略述如下：

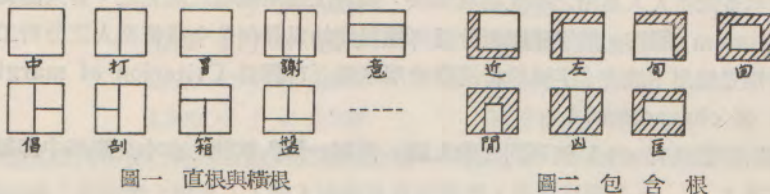
*Fujimura* 及 *Kagaya*:<sup>(4)</sup> 彼等計用字根二十一個，每個字根上有二至三個連接點；以不同之定位符號表示字根相連接之關係。此法固然有其優點，但其表示法相當冗長，而且為遷就所定義之定位符號，而犧牲傳統筆順，尤其是兩個不相連接之筆畫，諸如“=”、“彳”等，無法加以描述，更是嚴重缺憾。

*Toshiyuki* 等<sup>(5,6)</sup> 或 *Haruhisa* 及 *Shigeo*<sup>(7)</sup>，其所用方法大致相似。利用許多定位符號，以表示字根間相關位置。依其觀點，用於數量較多之中國字，惟有增加定位符號數目，或定位符號下不能表示之字視為一字根。因此，難免淪於定位符號數目繁多，規則瑣碎，以及字根增加，應用困難之弊。例如 *Toshiyuki* 使用五種不同之定位符號 S. K. L. J. G. 分別表示困、字、進、司、原五類結構字形之字根相關位置<sup>(5)</sup>；其實，此種相關位置即“木”放在“口”中，“子”放在“宀”中，“隹”，放在“辵”中，“亠”放在“冫”中，“京”放在“厂”中的位置，在中文字根本身性質中已具備，今用五種不同之定位符號加以區別，豈非徒增不必要之記憶與乎操作之煩雜。

又以字合成之美觀言之，兩個字根所用定位符號相同，但其被包含字根位置之大小却不一定相若，如“問”與“岡”二字，*Toshiyuki* 同用定位符號K，但“口”在“門”內與“山”在“門”中所佔位置之大小，却未加注意；其實，“山”所占位置應較“口”為大，方能符合字形美化原則。

綜上所述，以往之字根研究，泰半由於未能深切了解中國文字之結構特性，更未能考慮字形之美化問題，以致成效欠佳，難於推廣應用；因此，吾人從事字根分析，勢非着重於國字結構特性之探究不為功。

### III. 中文結構之形式



圖一 直根與橫根

圖二 包含根

圖一為楷書文字可區分為數部份之常見形式，當然尚有若干較複雜之形式（倪耿先生共列出68種）<sup>(8)</sup>，但在應用計算機數學模式時均無問題。須特別注意者，端為圖二所示之包含根。所謂包含根係指一字根之位置有包含其他字根之勢，此種字根，依其包含之形狀與位置，可概分為七大類，（如圖二所示），茲各舉一字為例，以說明之，如「近、左、旬、回、閒、凶、匡」等字中之「斤、工、日、口、月、×、王」均為被包含之字根，其餘部份即分別依序為圖二所示包含根之例。包含根與被包含根常有一定之相關位置，不容輕易更動，此特性對於定位符號之省略，有極大助益，被包含根之地位及大小與字形之美觀甚有關聯，因此字形組合之設計，必須注意及此。

圖一中，字根分為橫向或直向排列，可稱之為橫根或直根，甚多字根只作橫根或直根用，且可分為左橫根、右橫根、上直根及下直根，其縱橫比例亦多有一定者，可減省字之復合作。

### IV. 重複原理與字根

各國文字均有其重複（Redundancy）之處，中國文字之重複處在於字之構形。多數字可以析分為數部份。若一字可分為數部份，每部份各有其他字之部份與之相同者，則各部份稱為字根。若取字不加析分時，則此字亦可稱為字根。字根實即重複之源由。

在諸字中共含之字根為數甚衆，設若一概漠視之，不加利用，而視字彙為整個之個體，則無異不解中文之外籍人士，視中文如圖畫，而以處理圖畫之方式處理中文，用之於傳真，固無不可，用之於計算機，則頗有商榷餘地。

設採一字所占方格劃分為32×32小格，小格中有墨者（或墨多者）以黑點或「1」代表，無墨

者(或墨少者)以空白或「0」代表,如此一字所需之記憶容量為 $32 \times 32 = 1024 \text{ bits} = 128 \text{ bytes}$ ,以貯一萬字計,共需記憶容量 1280 Kbytes! 此項數字,頗為驚人,非中、小型計算機所足以勝任者。此外,鍵盤上將有極多之字數,則打字工作又非專業不辦矣。

究其原因,實因視國字如圖畫所致,造字原則如形聲、會意等,導致字形重複之處,絲毫未加利用之故。其實國字雖多衍變,大部仍井然有序,字根之數甚為有限,並非任何形狀之圖案均可能出現者,若然,則前述巨額之記憶容量固無法減省矣。

復有一種構想,謂中國文字不過點橫直撇捺等而已,歸併為八類,良可應付,賦以編號,設八鍵之鍵盤足矣,窮所有國字,未有出此範疇者。法似甚善,然欲計算機能辨認分別一字,則一字一碼,一碼一字(One-to-one Correspondence)為必要條件,在此分筆法中,不易作到,此點與注音符號分類法之困難相同。倘加入人為編號(一如電碼),則有強記背誦與專業訓練之苦事,難以推廣應用。此外,尚有筆順次序,簡、異字體及輸入費時等主要困難。

以上僅就輸入而言,若輸出則分筆法幾不能適用,蓋上一筆與下一筆之相對位置及長度變化,過於繁複,處理極端不易,更遑論字之整齊與否矣。

### V. 字根分析之邊際效用原則

綜上所述,關於中文字根之重複性,如利用不足,則記憶容量與鍵盤規模龐大,成本過昂,影響中文電腦化之推廣;如利用過度,則輸入速率過緩,且軟體程式過繁,成本亦增,倘有字形過劣之問題不易解決。前舉兩例恰為二極端情形,以一字為一圖案,絲毫不加析分者,利用不足(毫無利用),以一筆為一單位,將字全部分解至極限者,則利用過度。

欲求一成本至低而人人易用之中文處理辦法,須於上述兩極端之間覓之。吾人認為兩極端間存在一最佳(Optimum)解答。國字究應析分至何種程度始為最佳?今為使吾人之分析工作有一理論根據可循,特探究提出下述之「字根分析邊際效用原則」(可譯為 Criterion of marginal utility for analysis of characters):

假設一輸入裝置具有  $n+1$  個不同之輸入鍵,而每一輸入訊號(在本文應用中即每一中文字)平均需按鍵  $s$  次。

今考慮其中一鍵應否取消,設此鍵所相應之輸入訊號,即中文字,可分為二部份(即分為二字根),若其中有一部份與其餘  $n$  鍵所代表者均不相同,則此訊號字如分解為二部份為無意義,因必須新設一字根鍵,故雖原鍵因分解而取消,並未減少總鍵數,而徒然增加平均按鍵次數。

若一鍵所對應之字可分解為二字根,而此二字根均已包含於其餘之  $n$  鍵中,則如將此鍵取消,可收減少鍵數之利,其所代表之字仍可經其他鍵而輸入(如何將字根再組合成字見<sup>(9)</sup>),鍵數則自  $n+1$  減少為  $n$ 。

但此時按鍵次數即增加,原按一次即可,今則需按二次(另尚需按組合鍵一次,以使兩部份復合為一,惟組合鍵極少且常用,故其費時可略而不計),至於其影響平均按鍵次數如何,須視此一字(或字根)之常用程度而定。設  $f_i$  為字之出現頻率,此字在全部字樣  $\sum_{i=1}^N f_i$  中出現  $f$  次( $N$  為全部字數),則此字出現之機會為  $f/\sum_i f_i$ ,當此字(或字根)出現時,按鍵次數則增 1,故平均按鍵次數應自  $s$  次增加至  $s+f/\sum_i f_i$  次。

在應用輸入裝置時,鍵數愈少愈方便,又每一訊號平均按鍵次數愈少亦愈方便,故使兩者之乘積為極小值,頗具意義(其實際意義,在後文中將再討論)。

第  $n+1$  鍵如不取消,則此乘積為  $(n+1)s$ ;如取消,則積為  $n(s+f/\sum_i f_i)$ ,此鍵取消與否,當視兩乘積之大小而定,若

$$(n+1)s > n\left(s + \frac{f}{\sum_i f_i}\right),$$

則當取消之。上式即

$$s > nf / \sum_i f_i^2,$$

或

$$f < \frac{\sum_i f_i^2}{n} \tag{1}$$

此處  $\sum_i f_i^2$  為取全部字樣數，按基本用字研究一書，即 2,022,604。s 本隨 n 而變，理論上每一 n 確定後，可統計一新平均按鍵次數 s，但此統計工作甚為繁巨，故可以嘗試誤差法 (cut and try) 決定之。事實上，s 為一變化甚緩慢之函數，當 n=8,532 時，s=1，n=496 時，經統計 s 之值為 1.9。計算得 f 諸值如下（通常 n 甚大於 1，故 n+1 與 n 值近似）：

鍵盤鍵數	分解為二部份之字 其出現次數應小於(約值)
50	76,000
100	38,000
200	19,000
300	12,000
400	9,500
500	7,600

在一字分解為三部份或以上，則 (1) 式應修改推廣。設一字分爲 r 部份，則此字出現時增加按鍵 r-1 次，故平均按鍵次數爲  $s + (r-1)f / \sum_i f_i^2$ ，得分解條件爲

$$f < \frac{\sum_i f_i^2}{n(r-1)} \tag{2}$$

在 n=500 時，計算得下表：

$f > 7,600$	不分解
$3,800 < f < 7,600$	可分解為二部份
$2,500 < f < 3,800$	可分解為三部份
$1,900 < f < 2,500$	可分解為四部份

前文中述及使用方便之程度，此詞似甚為抽象，不易確定，惟細推使用時吾人之主要考慮條件，無非學習速率（即學習之難易），輸入速率及成本等項，今分別論之。

學習速率與鍵數 n 有關，因如欲操作純熟，須記憶各鍵之位置，如記憶之事物在人腦一般之容量以內，吾人認為記憶所需之學習時間大致與欲記事物之量（此處即鍵數）成正比；如此則可假定學習速率與 n 大致有線性關係。

輸入速率以尋鍵時間及平均按鍵次數兩者為主要決定因素。尋鍵時間有二情況，一為熟諳者所需，包括記憶之讀出及移動手指時間，一為生手所需，包括鍵面文字之逐行尋找，兩者均不難理解為大約與 n 成正比。至於平均按鍵次數即 s。

成本復可分為機件費用、打卡員薪資及計算機租金。機件費用包括按鍵開關數及線路材料以及裝置人工，易知其均與 n 成正比。打卡員薪資視輸入速率成正比，與平均按鍵次數最有關聯。至計算機租金亦視輸入速率而定，由於計算機運用之快速，吾人可預見輸入過程將成為瓶頸，故租金一項，主要亦與 s 有關。

綜上所述，使用之不方便程度，其中各項因素均大致與 n 及 s 成正比。若成正比，則(1)(2)二式為必然結果。但二式之成立不以正比為必要條件，事實上，若不方便程度為 ns 之任意遞增函數，如  $\phi(ns)$ ，均可得相同結論。

為探討起見，設此項不方便程度  $\phi(n^p s^q)$  為與 n 之 p 次方及 s 之 q 次方有關，則一鍵應予取消分解之線性近似條件為

$$\frac{\partial \phi}{\partial n} \Delta n > \frac{\partial \phi}{\partial s} \Delta s, \tag{3}$$

式中  $\Delta n=1$ ， $\Delta s=(r-1)f / \sum_i f_i^2$  故得

$$pn^{p-1}s^q\phi^1 > qn^p s^{q-1}\phi^1 \frac{(r-1)f}{\sum_i f_i}$$

或

$$f < \frac{p}{q} \frac{\sum_i f_i}{n(r-1)} \tag{4}$$

即分解為  $r$  部份時各  $f$  限值之比例不變，僅相差一  $(p/q)$  因數。

### VI. 字根之分析及頻率統計

由上理論，理想之字根選擇，為每分解一字重新計算頻率分配一次，逐字進行，求其邊際效用，至無再應分解者為止，但事實上由於繁複程度，難以作到，實際選字根之過程如下：

中文以左右及上下可分為二部份之字數最多，尤以左右可分者居多。第一步先將左右可明顯分為二部份之字分解，相同之字根歸併，並將其所屬原字頻率相加，如此可將 9,129 字減少為 3,256 字及字根。第二步復將上下可明顯分為二部份與包含結構之字分解，結果字及字根數減少至 621，經簡單分解，已將原有字數減少約十四倍。

須注意者，如前論，應保留出現次數大於 (2) 式條件之字，不予分解，如此共保留 25 字。此外，一字分解後，其字根之一，無在他字中共用者，為孤獨之字根，亦不予分解。

再將前述原則應用於字根，作第三步之分解，並檢查第一、二步中不合原則之分解，予以恢復，如此共得字根 448，原保留之 25 字以及罕用字根 23，合計字根共為 496。（見表一，表二）此數雖仍較吾人所樂意之數為多，然如再加分解，則不免支離破碎，蕪雜而輸入費時，已如前述。

496 字根中，含有常用熟字數 305（為整體之字），其出現次數超過所有字數之半。當出現時，僅需按鍵一次。

字根之使用率累計百分比如下：按常用次數之序排列，最前五字根占 11.3%，最前 25 字根占 30%，最前 50 字根占 49%，最前 100 字根占 66.7%，前 200 字根占 84.9%，前 300 字根占 95%，前 368 字

表一：中文字根表甲（依頻率高低排列）

口	丨	日	白	儿	之	門	木	一	言	三	女	月	冫	人	文	才	也	口	小	冂	疋	才	丕	彳	人	立
走	小	广	寸	文	糸	我	乂	夕	貝	了	目	十	田	禾	又	才	心	大	力	广	疋	八	丕	彳	人	立
上	方	王	巾	在	巾	西	竹	尸	回	來	目	頁	尔	古	示	里	立	二	佳	疋	土	至	工	彳	人	立
止	車	生	兰	虫	去	巾	几	戈	金	艮	門	者	山	耳	日	尤	中	勺	佳	疋	干	乍	天	欠	厂	
心	弓	用	犬	于	馬	巳	几	戈	丁	八	者	山	耳	日	尤	中	勺	佳	疋	干	乍	天	欠	厂		
母	彳	少	牛	五	又	章	正	四	雨	看	重	水	艸	其	彳	弟	月	戊	石	豆	衣	為	米	足	下	
手	直	長	木	丁	更	皮	勿	之	韋	七	丹	雨	木	艸	羊	而	刀	永	山	衣	亥	非	一	告	亡	
产	东	夕	吕	卜	夂	纒	巾	纒	舟	牙	木	艸	艸	艸	羊	而	刀	永	山	衣	亥	非	一	告	亡	
臣	古	气	マ	身	产	夂	九	高	舟	牙	木	艸	艸	艸	羊	而	刀	永	山	衣	亥	非	一	告	亡	
X	太	求	収	巳	乃	予	魚	求	十	廿	木	艸	艸	艸	羊	而	刀	永	山	衣	亥	非	一	告	亡	
黑	斗	甫	羽	申	麻	巾	兆	飛	戈	办	无	央	辰	共	彳	半	予	刃	制	坐	子	鬼	丙	片	帶	
非	丘	拉	毛	東	夂	州	兆	飛	戈	办	无	央	辰	共	彳	半	予	刃	制	坐	子	鬼	丙	片	帶	
卜	土	凡	束	末	牛	乎	母	川	久	戈	承	瓜	甲	易	東	凡	羊	史	門	史	骨	具	乙	段	瓜	
耒	反	乘	巨	产	由	小	川	戈	承	瓜	甲	易	東	凡	羊	史	門	史	骨	具	乙	段	瓜	瓜	瓜	
有	麗	夕	束	兼	互	丞	衰	甚	与	屯	丙	艸	艸	艸	羊	史	門	史	骨	具	乙	段	瓜	瓜	瓜	
少	广	喪	夂	兼	互	丞	衰	甚	与	屯	丙	艸	艸	艸	羊	史	門	史	骨	具	乙	段	瓜	瓜	瓜	
身	包	窟	夂	兼	互	丞	衰	甚	与	屯	丙	艸	艸	艸	羊	史	門	史	骨	具	乙	段	瓜	瓜	瓜	
尸	尤	冂	夂	兼	互	丞	衰	甚	与	屯	丙	艸	艸	艸	羊	史	門	史	骨	具	乙	段	瓜	瓜	瓜	
☆	的	是	有	他	這	國	們	說	個	就	要	全	到	以	你	時	那	裡	和	道	得	家	麼	後	樣	
井	盟	可	了	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂	夂

- 說明：1. 本表依字根出現頻率之高低由左而右，由上而下順序排列。
- 2. ☆為酌留“常”用字“井”為罕用字根。
- 3. 本表計收字根 448 個，酌留常用字 25 個，罕用字根 23 個，總計 496 個。

表二：中文字根表乙（依筆畫多少排列）

(一畫)	(二畫)	(三畫)	一	丁	乙	フ	丨	レ	く	ノ	厶	(二畫)	上	ノ	ナ	又	了	十	力	匕	二	厂	
丁	乃	尸	匕	九	マ	乙	厶	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹
入	乚	夕	乚	力			(三畫)	止	之	之	之	之	之	之	之	之	之	之	之	之	之	之	之
三	士	下	子	己	九	才	又	乃	卅	巳	乚	勹	子	也	厶	才	寸	大	上	口	小	止	山
少	女	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕	夕
方	斗	为	木	丕	王	白	主	尤	巴	犬	戈	五	开	云	产	尹	夫	夂	艸	土	冫	冫	冫
牛	代	反	乚	丕	井	尸	弓	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹
宀	内	母	文	心	介	夂	戶	欠	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹
艸	花	头	事	半	疒	古	去	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹
兵	羽	共	心	血	手	王	介	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹	勹
(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)	(七畫)
的	留	常用	字	他	以	有	全	你	那	的	到	和	是	导	彼	們	個	時	家	這	國	得	就
罕	用	字	根	之	也	口	子	可	凡	又	心	囟	囟	囟	囟	囟	囟	囟	囟	囟	囟	囟	囟
罕	用	字	根	之	也	口	子	可	凡	又	心	囟	囟	囟	囟	囟	囟	囟	囟	囟	囟	囟	囟



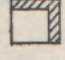
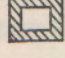
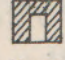
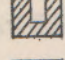
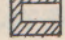
說明：本表按字根筆畫多少排列；筆畫相同者，按具起筆點，橫、直、撇、捺順序；起筆相同者，按頻率高低排列。

根占99%，其次50字根占0.8%。

VII. 字根之分類

由前述中文結構之形式，可知每一中文字根與另一字根間，常有一定的相關位置，且其縱橫比例亦多有一定，此一特性，對於字根組合所用定位符號之省略，有極大助益，因此若能將字根按其特性，妥為歸類，則表示字根間相關位置之定位符號，似可省略，如此必可使操作方便快速。茲將前節分析所得字根，略加歸類如下，以供進一步字根分類之參考。

(一)包含根：係指一字根之位置包含有其他字根之勢。如“間”字中，“門”為包含根。可概分為七類：(下圖所示，設一單字之形狀為口，斜線部份為包含根位置，空白部份為被包含根所佔位置。)七類之中，尚有程度之不同，茲舉例如下：

-  : 之、之、九、走、毛、尤、久……
-  : 宀、宀、宀、尸、戶、尸、戶、夕、尹、……
-  : 勹、勹、气、戈、戈、戈、戈、幾、鳥……
-  : 口
-  : 門、門、羸、戊
-  : 丨
-  : 乙

(二)左橫根：係指一字根隨其後出現之其他字根，必接於其右。如“位”字中之“亻”為左橫根，此種字根為數甚多，諸如讠、扌、彳、艹、礻、示、礽、斗、艹、彳、子、彳、冫、彳、彳、彳、彳……等是。

(三)右橫根：係指一字根，在其前出現之字根，必位於其左，如“到”字中之“刂”為右橫根，又如：欠、頁、欠、東、艹、斗、爻……等是。

(四)上直根：係指一字根，隨其後出現之其他字根，必接於其下一區位、如“箱”字中之“竹”為上直根，又如艹、夕、宀、彳、冫、艹、卜……等是。

(五)下直根：係指一字根，在其前出現之字根，必位於其上，如“第”字中之“弟”為下直根，又如氏、乚、儿、手、爻、止……等是。

(六)中根：凡不屬於上述任一情形之字根、統歸為中根，諸如口、日、田、十、王、山、女……等是。

### VIII. 中文字之語文模式

各字根可以點矩陣或線段之形式儲存資料於計算機，輸出時則以點矩陣為之。點矩陣即前述將字根所在之方格分為許多小方格，按黑白依序為「0」、「1」之訊號。此法尚可分為二類，一為單線字，即一橫或一直僅佔一行（例如藝術字體），優點可用較小矩陣，節省記憶容量，另一為粗體字，每筆粗細可變（例如毛筆字），優點為字形逼真美觀。線段之形式則視筆劃為直線或折線，以其起點終點位置資料儲存，可更節省記憶容量，此法用單線字。

字根儲存之形式，所佔空間之大小，及其分類等之資料結構 (Data Structure) 在<sup>(9)</sup>中詳細討論。

輸入時另制定四種基本定位符號 (Relational Operator)，即 ( ) △, ⊕, 及 ⊗, 其優先次序即如所列。△表包含，隨於△後之字根為△前字根所包含。⊕及⊗表上下及左右組合，括號之作用一如數式中者，為最優先。因包含根以局部居多（即僅包含一根最為常見），故令居優先。又圖一第二行前二形較後二形為常見，故令⊕較⊗優先，在後二形中使用時則加括號。

舉例以明之，「休」「類」「盟」「葡」諸字之按鍵法如下：

休	亻△木
類	米△犬△頁
盟	(日△月)△皿
葡	艹△匍△艹

當括號中僅有二字根時，前後括號之一尚可能省略。

二字根作縱向或橫向復合時，其相對大小比例即同於長或寬之指標數之比例<sup>(9)</sup>。例如「休」字二部份作橫向復合時，若「亻」根寬指標為1，「木」根寬指標為3，則「亻」占寬度 $\frac{1}{4}$ 「休」占寬度 $\frac{3}{4}$ 。

根據以上之分析，可將任意一中文字分解成爲一個字根與定位符號之表示式，此結構即爲一個有優先序的文法結構 (Precedence Grammar)，可稱之爲中文字之語文模式或中文模式，用 Backus Normal Form 可將之表示如下：

<中文用字> ::= <文字> / <符號>

<符號> ::= + / - / × / ÷ / ? / . / , / : / ; / ( / ) / { / } / [ / ] / \* / \$ / # / / / \ / ……等等  
並加上大小寫英文字母以及希臘字母數學符號等，多少不拘。

<文字> ::= <字根形> / ( <字根形> ) / <文字> <定位符號> <字根形> / <文字> <定位符號> ( <文字> )

<字根形> ::= <字根> ⊕ / <字根> ⊗ / <字根> 8 / <字根> ∞ / <字根>

<字 根>::=496個(如表一)

<定位符號>::=△/△/△

說明

1. 品符號表示成「品」形重疊，例如車品即成轟。
2. ∞符號表示橫向重複二個，例如木∞即成林。
3. ∞符號表示橫向重複三個，例如口∞即成□□□。
4. 8符號表示縱向重複二個，例如火8即成炎。

以上四個符號，只是爲了使用方便而設，故稱之方便符號。

以上中文模式，是一種再生形式(Recursive form)，寫成此形式，純爲表示簡潔方便，事實上，並非再生形式之結構，因爲表示式長度，實爲有限長度，在倪耿論文中，曾統計此形式約有六十八種之多<sup>(8)</sup>。

在以上中文模式中，定位符號有優先次序，若將括號及品8∞符號考慮在內，可排列爲

品, ∞, 8, ∞ > (> △ > △ > △ >)

此結構提供了必備之資料，可將中文表示式組合成中文單字，其中方法之一便是用結構樹分析法(Structure Tree)詳如本期中文字根的貯存和中文字的合成一文。<sup>(9)</sup>

## IX. 中文模式與中國字

本文研究之中文模式，雖僅有496個字根，然若用以組合中國字，則除本研究所採用之9,129個單字毫無困難外，尙可組合成歷來所有中國字之97.61%!!

我國文字之總字數，歷經數千年之累積，迄今已達五萬字之譜。張其昀先生等編纂之中文大辭典，係歷來所有字書中單字數之最多者，該書網羅經史子集，旁採類書、叢書、輯佚書、字書、辭典……等資料，全書計收單字49,905字，堪稱集古今中國字之大成<sup>(10)</sup>。以本研究之中文模式，逐字加以核試，結果於該49,905個單字中，計可組合48,713字，僅有1,192字難以組合。然於此吾人必須特加留意者，端爲此類難以組合之字，泰半已成死字或爲異體字<sup>(8)</sup>。諸如：“禮、尗”等字，其義未詳；“屮、用、禱”等字，各爲“鬲、用、禱”之本字；“𠂔、膏、夙、昔”等字，各爲“龜、壹、官、齒”之古字；“饒、屯”等字，各爲“饒、中”之籀文；“冎”係“𠂔”之篆字；“圖”爲古器；“飛”爲飛之簡體；“𠂔”與“滂”同；“楸”爲“椒”之俗字……等是。

再者，誠如前述，中文模式具有再生性質，因此除可組合已有之48,713個中國字外，尙可產生新字，對於日後爲應時代、社會需要而增加之新字，未嘗不無裨益！

## X. 鍵盤設計

前文中所得496字根各以一鍵代表，作成一鍵盤，以爲輸入之始端機器。鍵盤中字根之排列法與使用之方便與否有極重要之關係，其主要之考慮包括記憶之方便，尋找之方便及手指移動平均距離之減少等，以增加速率並降低輸入者之勞累，庶幾爲大眾所樂用而推廣流行。

設計時考慮之因素可列舉如下：

(一)根之類別：凡屬包含根者，極易區別，可分列之；有爲明顯之左橫根、右橫根、上直根及下直根者，亦可分列之；此外，有左右上下中俱可明顯應用，筆劃甚少而極常用者，可列中間；尙有獨字常單獨應用，極少組合爲他字之根者，可列於較偏僻處。

(二)字形：字形相似者並列，如「己、巳、巳、巳」等，對於記憶極有助益。



- (三)筆劃：筆劃極少以及極多之字，一目了然，為極易分別之事，惜筆劃適中者，不易分別。
- (四)首筆分類：取一字之首筆及第二筆等，視其為點、橫、直、撇等加以分類。
- (五)部首：因字典之通用，熟習部首者極衆，但有難查之字。對於字根，則不易適用。
- (六)注音：注音可一如英文字母排列，消除許多不便，但有同音字問題及不諳注音者。對於不成字之字根無法應用。
- (七)字義：同義字、反義字及相似詞等可並列，減少記憶困難，但僅能局部應用。
- (八)詞彙：考慮字與字間之關聯，即連詞之應用，如並列，則可減少移動距離及尋找時間，亦只能局部應用。
- (九)頻率：以熟字僻字分列，可增輸入效率，鉛字盤中已有應用，但本法在實用上僅能作粗疏分隔。
- (十)歌訣：可助記憶，例如與地命相之類，藉歌訣傳千百年之久。如將無可分類之處，編為歌訣，或亦可助益。

## XI. 結 語

- (一)本研究係根據中文電腦基本用字研究一書所列單字 8,532字及其異體字 597字，合計 9,129字，經分析統計而求得 496個字根，由此，字形之儲存空間將可減少約 20倍；再者，字根之筆畫數遠較單字為少，應用資料壓縮法 (Data Compression Schemes) 亦可使其儲存空間大幅減少，<sup>(9)</sup>合而觀之，其儲存空間總計至少可縮小 40倍以上。
- (二)該儲存空間，已縮小到適合使用 Mini-Computer 或 Special Wired Read Only Memory 實際製造，故有望利用 LSI 技術，大量生產，以利國人使用。
- (三)此組字根尚可於計算機中再分，使用時以軟體組合，更可節省儲存空間，此點亦頗有研究之價值。
- (四)由於文中之中文模式具有再生性質，故用此 496個字根，不僅可組成 48,713個已有單字，而且尚可利用語文模式，創造新字，此一特點乃為其他方法之所無。
- (五)加碼 (Encoding) 之方法可與其他中文加碼方法互相轉換，僅需一對照表而已，而此對照表亦遠較本文第三節中所述之儲存空間為小。
- (六)用此方法組合單字應不困難，訓練容易，使用方便；至於如何將字根組成單字以及字形之美化問題，本期另有專文論述。<sup>(9)</sup>
- (七)目前每字之按鍵次數，若符號不計，平均僅為 1.9；而△，△△，△△△等符號之省略，或可用字根分類法為之，乃為尚待研究之問題。
- (八)應用中文字根編排索引，對於中文資料之檢索以及中文字典之編排，提供一種新式方法。
- (九)以此字根設計之中文鍵盤，其鍵數較一字一鍵法減少約 20倍，如此不僅可使鍵盤之製造大幅減少，操作亦大為簡便，對於中文打字機及計算機中文輸入裝置之應用，頗具價值。

## 誌 謝

本研究之完成，承蒙杜敏文教授提供許多寶貴意見；關於字根之分析統計，復蒙劉冠軍、柳肇輝、林明達、黃玉美、林家榮、曾七雄、秦順諸先生協助，對本研究貢獻良多，謹此一併致謝。

## 參 考 文 獻

1. 見燕京大學索引。
2. 彭仁山：教育研究二十三期，中山大學，民國十九年。
3. 林樹：中文電腦基本用字研究，交通大學工學院，民國六十一年。

