

新中文字根集合的設計

The Design of a New Chinese Character Component Set

杜敏文 Min-Wen Du

Department of Computer Science, N. C. T. U.

(Received August 9, 1976)

ABSTRACT — A new character component set for chinese input/output of a computer has been designed for the following two ends. First, the maximum number of components necessary to identify a Chinese character can be reduced to only 4. Second, characters composed from components may have better appearance. The first end is achieved through introducing new character components. The second purpose is approached by introducing character component variation concept.

摘要——本文敘述如何設計用以作電算機中文輸入輸出的一組新字根。該組字根設計目標有二。即一，最多只要用四個字根即可單一區分每個中文字。二，使由字根組合出來的字能較為美觀。我們增加了一些新字根，同時使用字根變形的的方法以達到這兩個目標。

一、簡 介

在交大巴完成的CHIPS中文計算機輸入輸出系統，解決了字碼輸入、字形貯存以及內用字碼等基本問題〔1〕。但對使用者而言，却有兩點需作改進。第一，雖然平均每個中文字只需打1.98個字根鍵，但一個字最多却須打至九個鍵，例如：麤：鹿七七鹿七七鹿七七。第二，由於有些字根在不同的字中有不同的形狀大小，一個字根字形設計不易使所有由該字根組合的字都很美觀。本文將從設計新字根集合着手來解決這兩個問題。不同字根數從原來的460個增至532個使得在新字根系統中最多打字根鍵4下即可尋得系統內8532字中任何一字。同一個字根給予數種字根的構想也使得由字根組成的字外形較為美觀。

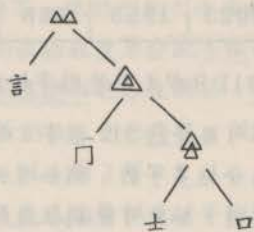
二、CHIPS系統的文字表示式及字形說明檔結構

在CHIPS系統中，每個中文字是以一個二元樹來表示〔1〕。每個樹的外節（External node）為字根，內節（internal node）為關係符號。每個部份樹的根（root）都代表一個字的部份字體。關係符號則表示部份字體的相對關係。在系統中共使用三種關係符號，即橫連△△，直連△和包含△。字根數則共460個（不包括同字根異形體，如才，木，圖一），可表48,713個中文字〔1〕。

為了節省貯存空間，使字形說明檔內每個單元定長，以及尋字方便，CHIPS中將不同數目字根的字分別建表貯存。同時，每個字都用抽出式波式法（Extracted Polish Form）表〔1〕。每個文字表示式分成三個部份，一個字根序列，一個關係符號序列及一個字根關係符號向量。根

據這三個部份就可把文字從字根組合起來。

例一：調的二元結構樹如圖二。

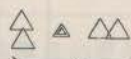


圖二 “調”字的二元結構樹

其抽出式表示式為

言 冂 士 口

字根序列



關係符號序列

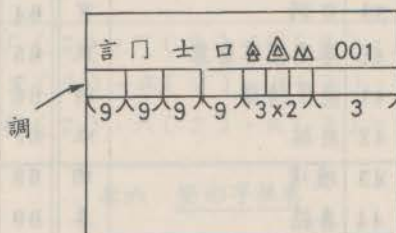
0000111

字根關係符號向量

字根序列及關係符號序列排列順序各為二元結構樹以波氏法 (Polish Form) 表時的相對順序，字根關係符號向量中 0 表字根，1 表關係符號，此向量表二元結構樹以波氏法表時之字根與關係符號排列位置。

由於以抽出式波氏法表示，字根序列與關係符號是分離的，將此種文字表示式列表建構後字根序列即可當做尋找鍵尋找文字表示式。若將字形說明檔依照字根大小順序排列，則可使用二元搜尋法 (Binary search) 有效率地找到為某一字根序列的字。

例二、調有 4 個字根，所以應於 4 個字根的字形說明檔中尋找。圖三為 4 個字根字的字形說明檔。CHIPS 系統中共用 499 個字根，三種關係符號，所以每個字根用 9 筆 (bit) 表，每個關係符號用 2 筆表。4 個字根的字其字根關係符號向量長度可省為 3 筆 [1]。



圖三 4 個字根字的字形說明檔

三、減少尋字所需字根鍵的數目

表四為在 CHIPS 系統中，所有 8532 個字使用字根數目的統計。由這個統計我們可看出超出 5 個字根的字共 1208 個。這對使用者是一大負擔。是否可能不用全部的字根序列即可單一地找到一個字呢？我們嘗試的方法是用字根序列前面 (左起) 的幾個字根來單一區分所有的 8532 個字。經過程式測試的結果，發現若是只用每個字最前三個字根來尋字，將有 933 個字不能區分。若是用最前 4 個字根來尋，則僅有 213 個字不能區分。所以我們選擇用最前 4 個字根為尋找鍵的方法。不能分的 213 個字則設計新字根使能單一分出。

字根數	1	2	3	4	5	6	7	8	9
字數	271	2178	2923	1953	888	221	79	18	1

表四 CHIPS 中文字使用字根數目統計

表五為CHIPS 中用 4 個字根不可區分的 213 個字，以及新增用以區分此 213 個字的字根。增加字根的原則有二，一為選擇可區分較多字的，例如增加口則許多字皆可由不可分而成可分。其二為由使用者看來該字體成為一整體，如此可便利往後的由字分字根工作。

		增根		增根		增根		
1	番鷗	岡	27	償儻	尚	53	傲覽整整整整整整整	
2	雁罷	能	28	睿整寂	尙		督幣弊	甬
3	叢棘	棘	29	參耗	爻	54	尚峒	峒
4	蟬蟻	古	30	能態熊	能	55	從聳聳	峒
5	筦筦	元	31	兔剝	兔	56	愆衍	峒
6	簽籤	僉	32	芻雜敝	芻	57	蕭蕭	峒
7	讀讀	買	33	淳濼	古	58	痲癒	峒
8	麟麋	鹿	34	潮幹漸	卓	59	歲劇翻	峒
9	鑿擊	々	35	滄漱	僉	60	衛變	峒
10	郵靈	靈	36	澄滂澄滂滂	炊	61	僉斂	峒
11	鐵錫	古	37	湯盪燙	易	62	倉歛穌	峒
12	鐘鐃錦	尚	38	澆澆	鹿	63	鄉嚮響響	峒
13	晉鄙戢	晉	39	寅戢	寅	64	恐響響響響	峒
14	殄餐	歹	40	寒寒寒寒寒寒	寒	65	杼柳	峒
15	卒頓	卒	41	墜墜墜墜	隋	66	檀攬	峒
16	熟塾	古	42	捺振	双	67	檻攬	峒
17	裏瓢	裏	43	壇壇	回	68	系繡	峒
18	就驚	古	44	喜話	喜	69	景瞭影顯	峒
19	豈鷓鸞鸞鸞	回	45	尋郭	尋	70	早旰	峒
20	鼓馨整整整	鼓	46	曉噤	古	71	只叭	峒
21	區歐鷗歐歐	區	47	啜啜	双	72	顯歛	峒
22	匿愚	若	48	叻另	另	73	薛藥	峒
23	厭歷壓壓壓壓壓	厭	49	号頸鴉郭	号	74	敬榮驚警警	峒
24	旭咎	旭	50	單戰鄴戰戰戰	留	75	觀歡顛勸勸	峒
25	佼傲	交	51	吧邑	吧	76	莖莖	峒
26	倥僂	古	52	頃員	員			峒

表五 用 4 個字根不可區分的字

四、字根變形及組合字形之美化

一個字根在不同的字裡有可能需要不一樣的形狀大小才能使組合出來的字美觀。在〔1,2,3〕中使用定比重的辦法，根據每個字根的複雜程度來分配字根佔據的空間。如此解決了字根在字中的大小問題。同一字根不同形狀可用字根變形的方法解決。字根的變形可分為3類：

1. 一般變形：所組成字多而本身筆劃數又少的字根，須較多的變形方能適合各種不同組合的字。這類字根如、（點），ノ（撇），八（捌）等。

2. 偏旁變形：一些字根單獨成字及成為偏旁部首時字形不一樣，例如木，才；金，鈗；竹，𥵓，等等。

3. 包含變形：包含根包含位置應隨着被包含的部份複雜程度擴大縮小。為此一些包含根都給予兩個字根以示不同的包含位置。

表六所列為根據以上三種考慮所設的變形字根。

表七為原系統字根與新增字根總表，共 616 個。

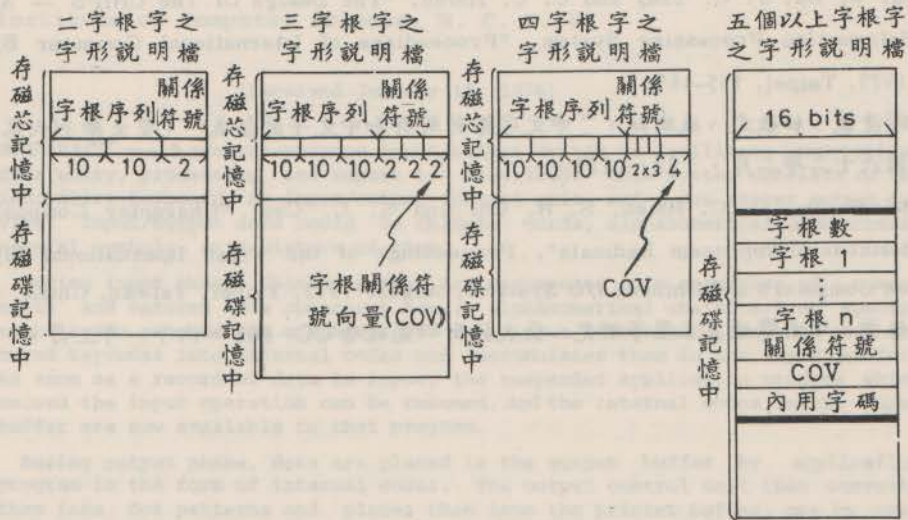
1. 一般變形：、（4個），ノ（4），一（4），丨（2），八（6），丿（2），厶（2）， ，乂（2），十（2），丁（2），卅（3），口（2），心（2）。
2. 偏旁變形：九（2），儿（2），大（2），土（2），乃（2），弋（2），弓（2）， 女（2），小（2），木（2），王（2），戈（2），尹（2），尺（2）， 气（2），毛（2），支（2），元（2），六（2），尤（2），夫（2）， 牛（2），攴（2），衣（2），瓦（2），四（2），𠂇（2），禾（2）， 疋（2），爪（2），耳（2），瓜（2），走（2），更（2），兔（2）， 克（2），金（2），尚（2），麥（2），糸（2），竹（2），食（2）。
3. 包含變形：乙（2），尸（2），冫（2），冂（2），山（2），口（2），宀（2）， 尸（2），又（2），夂（2），戶（2），夂（2），戔（2），夂（2）， 鳥（2），齊（2），夫（2），夫（2）。

表六 變形字根表

五、字根碼之訂定，字形說明檔安排及字形說明之搜尋

由於一個字的字形說明是由字根序列作搜尋鑰來尋找，而在鍵盤上為了便利使用者，變形的字根只用一個鍵。也就是說，鍵盤上的一個鍵有可能對應系統內的幾個字根。所以字根碼的訂定，字形說明檔的安排，都必需依此設計以加速字形說明搜尋所需時間。

圖八為字形說明檔安排方法。



圖八 字形說明檔設計

一個字根的字不需字形說明檔，因各字根本身就是字。2個，3個，4個字根的字各歸併於一檔內。屬於較常用的4千字的字〔4〕存於磁蕊記憶中，否則存於磁碟中。5個字根以上的字存於磁碟中。常用字存於磁蕊記憶中可增進字形說明搜尋的速率。

2, 3, 4個字根的字形說明檔，存於磁蕊記憶部份與存於磁碟部份可各別按照字根序列大小排列。5個以上字根的字形說明檔則按照前4個字根鍵排列。如此若出現的字存在磁蕊記憶中，可用2元搜尋法迅速找到；若存在磁碟記憶中，可先建好一索引檔，只一次磁碟的讀入即可找到。

所謂的字根序列大小此處須加以說明。由於一個字鍵有可能對應幾個字根，字鍵輸入字碼實際上只能對應到一組變形字根的共同部份。所以變形字根的字根碼設計，以能彼此相鄰而 Hamming 距離愈短愈佳。例如點共有4個，字根碼各為 0000000000, 0000000010, 0000000001, 0000000011。其左方8個數字皆同。表七中由左至右，由上至下，字根的順序就是設計成使變形字根共同部份可最大者。

六、結 論

本文設計了一個新的中文字根集合，使〔4〕中的8532個字依照由左而右，由上而下，由

外而內的簡單規則分析，只用最前 4 個字根即可單一地定出任何一個字。本文提出的字根變形方法使中文字由字根組合時可有不同的選擇，因此使組合成的字較為美觀。

七、誌 謝

本文中尋找有相同 4 個字根的程式由李立明先生完成，謹此誌謝。

參考文獻

- [1] M. W. Du, J. C. Tsay and C. C. Hsieh, "The Design Of The CHIPS — A Chinese Information Processing System, "Proceedings of International Computer Symposium 1975, Taipei, 155-165.
- [2] 謝清俊、杜敏文、戚樹紅，"中文字根的貯存和中文字的合成"，交大學刊六卷一期，民國六十二年二月，122~131。
- [3] M. W. Du, C. C. Hsieh, S. H. Chi and S. C. Chu, "Character Composition and Methods to Represent Radicals", Proceedings of the First International Symposium on Computers and Chinese I/O Systems, August 1973, Taipei, Taiwan, China, 63-78.
- [4] 林樹，中文電腦基本用字研究，交大技術研究報告 CC-601，六十一年三月。