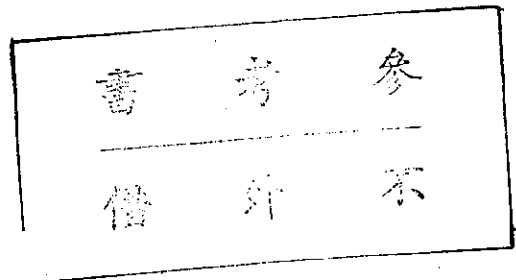


TR86-005

Computational Approaches in the  
Research of Chinese Computer-  
Topological and Geometric Descrip-  
tions for Chinese Characters

By Keh-Jiann Chen  
Associate Research Fellow  
Institute of Information Science  
Academia Sinica, Taipei,  
Taiwan, Rep. of China  
Tel: (02) 782-2002



中研院資訊所圖書室



3 0330 03 000057 9

0057

## ABSTRACT

This paper suggests computational approaches in the research of Chinese computer via utilization of the topological and geometric structure information of Chinese characters. A description language is designed for constructing the characters from radicals and the radicals from strokes. The parsing of the descriptions of a character provides the structure of this character which represented as the control points of composed strokes. Many statistical results about Chinese characters can be easily derived and hence support the research of Chinese computer. Furthermore, it is also capable of generating various fonts stroke by stroke in the writing sequence. The supported research areas include 1. character font generation 2. Chinese input and coding 3. character recognition 4. CAI on character spelling, and 5. data compression.

Key words: Chinese computer, font generation, Chinese input, CAI, computational approach.

## 1. Introduction

The study of computational approaches in the research of Chinese computer was motivated by the vast amount of effort had been paid in the design of Chinese keyboard and the dot matrices for character fonts manually. With the large capacity of character set, many simple routine jobs become tedious and time consuming. For instances, (1) the dot matrices design, for each Chinese character, it needs a lots of work to layout the dots on a  $n \times m$  rectangular grids by peoples to mimic the image of this character. (2) Researchers pay much of their time to design the Chinese keyboard in order to optimize the encoding. (3) On the character recognition, people tend to fail in the selection of better features due to lack of statistical data. The above tasks have common characteristics, i.e. they need to know the geometric and topological characteristic of all Chinese characters and then from which some common properties of Chinese characters can be found. The difficulties of finding the important properties of the characters is not caused by the complicate searching procedure but by the large set of search domain. Such kind of task is typically suit for the computational approach. Therefore we designed a description language which is used for describing the geometric and topological structure of Chinese characters in terms of radicals and strokes composition. Although, there were many description languages had been designed for composing radicals into character [1,2]. However, our language is the only one to describe radical structures in terms of strokes.

## 2. The Description Language

Basically, the Chinese characters are composed from radicals and the radicals are composed from strokes. First, we describe the way of constructing characters from radicals. For instance, the character 柳 is constructed from two radicals 木 and 卯 from left to right. If we device an operator  $\square$  which means compose two operands from left to right, then we are able to express 柳 as  $\square$  (木) (卯). However radical 卯, itself is also an composed character as  $\text{卯} = \square$  (夕) (卩). So,  $\text{柳} = \square$  (木) ( $\square$  (夕) (卩)). The syntax structure of this description language is as simple as follows

Tree  $\rightarrow$  operator (Tree) (Tree).  
Tree  $\rightarrow$  radical

There are 3 basic types of operators

1.  $\square$  : compose two radicals from left to right.
2.  $\square$  : compose two radicals from top to bottom.
3.  $\square$  : compose two radicals in overlapping.

Every Chinese character can be described from above three operators and its radicals.

However, for the convenience and the precision, the actual operators, we adopt, are in Figure 1.

Operators with parameters	Pictorial representation	Examples
$\emptyset 1$		杜 = $\emptyset 1$ (木)(土)
$\emptyset 1\{4\}$		唱 = $\emptyset 1\{4\}$ (口)(昌)
$\emptyset 1\{5,6\}^*$		林 = $\emptyset 1\{5,5.5\}$ (木)(木)
$\emptyset 2$		早 = $\emptyset 2$ (日)(十)
$\emptyset 2\{4\}$		号 = $\emptyset 2\{4\}$ (口)(方)
$\emptyset 2\{4,5\}$		昌 = $\emptyset 2\{4,5\}$ (日)(日)
$\emptyset 3\{0,2,4,0\}^{**}$		起 = $\emptyset 3\{0,2,4,0\}$ (走)(巴)

Note :\* Every character are supposed to fit into an unit square. The parameters of each operator denote the portion of the space occupied by the operand. The value of parameters is ranged between 0 and 1.

\*\* The parameters of  $\emptyset 3$  denote the boundary of second operand from top down, bottom up, left to right and right to left respectively.

Figure 1. The operators for composing radicals into characters

For instance, the character 榴 has the following tree structure.

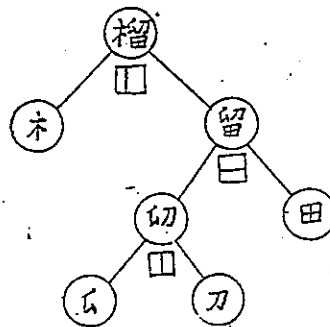


Figure 2. The tree structure of character 榴

The linear description of 榴 is  $\emptyset 1[4](\text{木})(\emptyset 3(\emptyset 1(\text{夕})(\text{刀}))(\text{田}))$ . Actually, each radical is coded as its radical code. Therefore, 榴 is  $\emptyset 1[4](0423)(\emptyset 3(\emptyset 1(0359)(0217)))(0514)$ . If the area is empty, we have ( ) denotes null radical.

The description of each character will be parsed to find the windows for the radicals in this character. If the window for a radical is smaller than the size of this radical, then this radical will be scaled to the proper size. Otherwise this radical will be shift to the center of the window. Usually the standard size of a radical is a unit square. For the purpose of character fonts generation, the representations of radicals are the geometric location of their respective containing strokes.

Each stroke in a radical is governed by three control points, i.e. head middle and tail point of the stroke. The geometric value of three control points for the strokes in a radical were derived from parsing of radical descriptions.

#### Compose strokes into radicals

No matters how the description language is designed, it should have the following functions (1) describing the length (x-direction), the height (y-direction), and the rotation angle of strokes, (2) describing the connection relations between strokes.

There are 19 different basic strokes as shown in Fig.3.

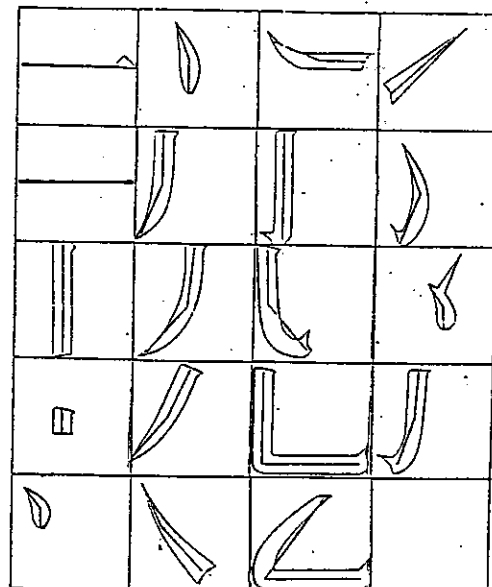
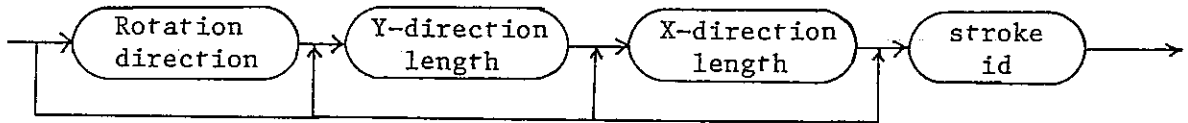


Figure 3. The 19 basic strokes of Shong style characters

However, the variations of each stroke are numerous. Therefore, the variations of a stroke should be able to be described from a standard one by applying length, height, and rotation operators. Every basic stroke has a standard size. For the convenience of character generation, we define the largest normal shape as the standard. The other variations are derived from the scaling and rotating operations.

The syntax diagram of the stroke representation is:



The detail notations for rotation and length as well as stroke identifiers see [3].

The relations between strokes may be connected or nonconnected. For the nonconnected relations, such as '𠄎', we use the operators as in the composition of radicals into characters. Therefore, "𠄎" is  $\square (\square (i) (\backslash)) (\square (\backslash) (\backslash))$ . For the connection relations, we devise an binary operator to determine the connected point. The reference points for a stroke are its control points namely H,M,T. For example, H denotes head point. MH denotes the center point between middle and head. We use the bisection operator to determine the exact connected location. Partial examples are shown in Figure 4.

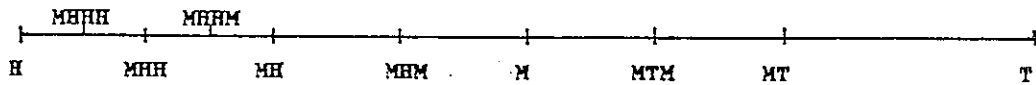


Figure 4. Bisection operators and their relative locations

Now, we are able to describe a complicate radical, for instance '鳥'.

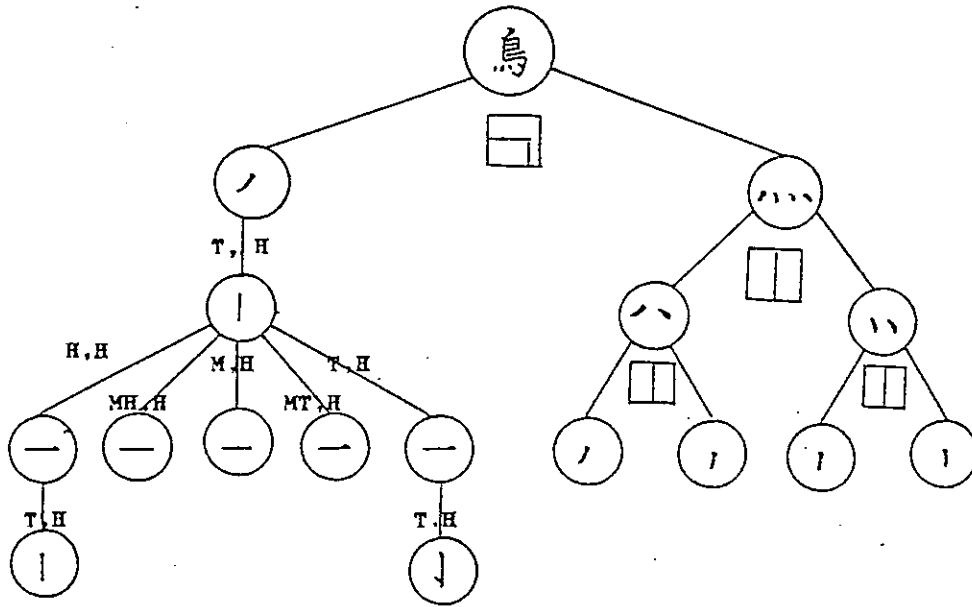


Figure 5. The tree representation of the radical '鳥'

#### Parsing of the description language

As we state before, the description of a character after parsed will derive the geometric information of each strokes in this character. The geometric information of the strokes is represented by the x,y coordinates of control points. For example, the control points of the character "鳥" is as shown in Figure 6.

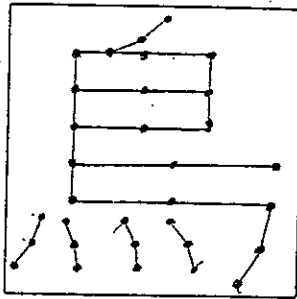


Figure 6. The control points of the character '鳥'

Since, characters are described as a composition of radicals by applying operators in the Figure 1 to form a tree structure. The parser looks the operators from top level down. Recurvely divides the windows into subwindows according to the operator. The top level window will be assigned by the user, usually a unit square. Then, the control points of the radicals are scaled and translated into the appropriated windows.

The parsing of radical description is a bit more complicate. The basic steps are as follows.

1. For individual stroke, get the control points of the stroke by scaling, or rotating the control points of standard strokes according to the operators.
2. If the strokes are connected by the operators, then find the connected points.
3. Each stroke is translated to the connection point according to the coordinates of the first stroke.
4. Find the minimax window of the radical, and shift the radical to the lower left corner of the unit square. If the radical size is greater than the unit; then scale to the unit size. Otherwise, keep it as it is.

After the completion of the parsing steps, each radical is represented by the coordinates of the control points of the contained strokes and the size of this radical. Most of radicals have size equal to unit square; some others may have size less than unit square. For example, the radical "一" has the size nearly zero at Y-direction and about 0.9 at X-direction. The size of radical is for the case when fitting the radicals into windows of a character. If the size of the window is less than the size of the corresponding radical than this radical will be scaled down to the appropriate size; otherwise, the radical will only be shifted to the center of the window. Here, 'the size' we mean the size in both X-direction and Y-direction.

### 3. Information provided by character descriptions

The descriptions of the Chinese characters provide the topological and geometric structures of each character from strokes to radicals then radicals to characters. Furthermore, we arrange the sequence of radicals and strokes in the description according to the writing sequence. Therefore descriptions of characters contain the following information.

1. The skeleton structure of each character serves the purpose for the generation of character fonts and partially for the character recognition.
2. The radical components of the characters serves the purpose for the studying of Chinese input or coding methods. Also, the representation of characters by radical compositions saves the storage in fonts generation[4] and character generation[5]. Figure 7 shows some sample results of Shong style font.



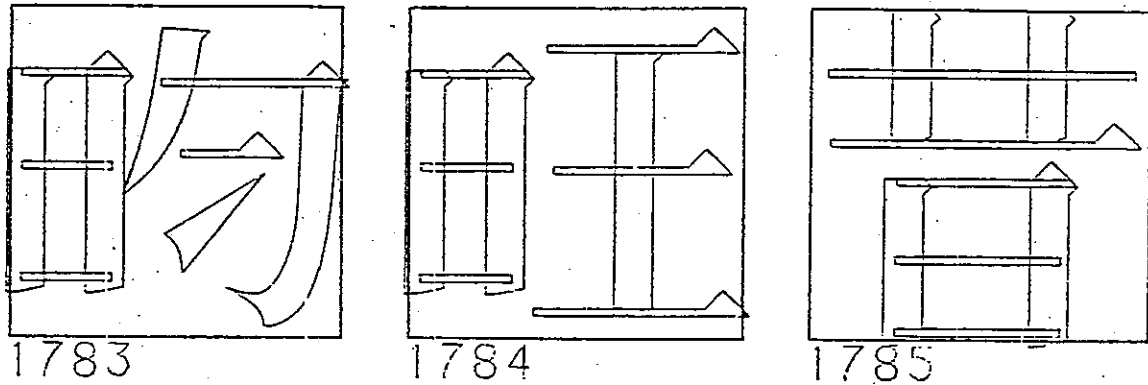


Figure 7. Some sample results of Shong style font.

3. Stroke sequence of each character provides a means of implementing low cost computer-aided instruction for Chinese character spelling [6]. Also, it provides the important sources of stroke coding method for Chinese which is useful in Chinese input and on-line character recognition [7].

The descriptions of character set can also provide many important statistical results about Chinese character such as, radical occurrence frequency, stroke distributions, stroke connection relations, corner shapes. In fact, any statistical result concerning about the structure of Chinese characters may be derived from these descriptions. Thus, they become very useful tool in the computational approaches for the research in Chinese computer.

#### 4. Conclusion

We had completed around 12,000 descriptions for Chinese characters. Such descriptions had been very useful in many researches conduct in our institute. Besides the fonts and dot matrices generation, many statistical results about the structure of Chinese characters had been completed from them. Here, we name a few.

1. The frequency distribution of radicals from 8000 characters, Table 1 shows the first set of 20 radicals which is quite agree with the result of Suen and Huang [8].
2. The character sets of the same stroke sequences, Table 2 shows the result. Here, we group the basic strokes into 7 different stroke types as shown in Table 3.
3. The character sets of the same strokes without sequence restriction, we found 363 groups. Table 4 shows the partial result.

5. References

1. Shi-Kuo Chang, "an interactive system for Chinese character generation and retrieval", IEEE Trans. on Systems, Man, and Cybernetics, Vol.3, 257-265, May 1973.
2. 謝清俊等 "中文字根的貯存和中文字的合成"，交大學刊，第六卷第一期，122-131, 1973.
3. 陳克健，鄭國揚，"ACCFONT - 中文字自動產生系統使用者手冊"，中央研究院資訊所技術報告 TR-83-005, 1983.
4. K.J. Chen, K.Y. Cheng & C.N. Chen, "On Construction and Generation of Chinese Characters", Intelligent System Imaging Technology and Software Engineering, Edit by S.P. Wang, P.99-114, 1984.
5. K.J. Chen & K.Y. Cheng, "A Model for Low Cost Chinese Character Generator", Proceedings of the International Conference of chinese Computing'85, San Francisco, U.S.A., 1985, P.F-1.1-P.F-1.10.
6. K.J. Chen & S.Y. Lin, "Computer Aided Instructions for Chinese Character Spelling with Compressed Coding", Proceedings of ICS'84, Taipei, Taiwan, R.O.C., P.306-310.
7. T. T. Hsieh, On-Line Recognition of Hand-written Chinese Characters, Master Thesis, Department of Electrical Engineering and Technology, National Taiwan Institute of Technology, Taipei, Taiwan, R.O.C.
8. C. Y. Suen and E.-M. Huang, "Computational Analysis of the Structural Compositions of Frequently Used Chinese Characters". Computer Processing of Chinese & Oriental Languages, Vol.1 Number 3 May, 1984, P.163-176.

ROOT NUMBER	ROOT1	ROOT2	ROOT3	ROOT4	ROOT5	ROOT6	ROOT7	ROOT8	ROOT9	TOTAL NUMBER	
口	301	314	303	574	336	86	8	2	1	0	1624
一	101	37	192	365	172	103	34	6	2	0	911
シ	328	557	47	0	0	0	0	0	0	0	604
日	419	120	191	166	75	19	5	0	0	0	576
木	423	395	36	8	0	0	0	0	0	0	439
ナ	410	295	107	10	0	1	0	0	0	0	413
月	421	153	54	147	43	10	3	1	0	0	411
イ	204	324	72	12	0	0	0	0	0	0	408
土	303	31	106	126	65	19	2	0	0	0	349
才	361	330	12	0	0	0	0	0	0	0	342
儿	205	1	27	127	114	14	26	13	1	0	323
去	605	223	19	50	22	7	0	1	0	0	322
工	202	82	193	24	10	3	0	0	0	0	312
田	514	44	92	82	63	22	1	0	0	0	304
口	302	36	86	91	55	8	17	1	0	0	296
言	702	210	35	15	6	15	1	0	0	0	282
金	801	273	5	1	2	0	0	0	0	0	281
又	229	7	61	102	68	36	3	2	0	0	279
又	343	185	58	23	7	1	1	0	0	0	275
又	323	97	167	9	2	0	0	0	0	0	275

Table 1. The first 20 radicals with the highest frequency distribution computed from 8000 characters

1	兀 九	19	叻 叻	37	牛 牛
2	刀 力	20	加 召	38	件 件
3	土 士 工	21	沼 加	39	捏 捏
4	千 广	22	管 筧	40	昼 睡
5	口 口	23	吊 引	41	茹 茹
6	己 巳 已	24	余 余	42	曼 晚
7	犬 犬	25	肝 旱	43	崑 崑
8	仄 木	26	杜 杜	44	緬 緬
9	天 六	27	味 景	45	巨 言
10	日 日	28	音 巨	46	丸 戈
11	且 且	29	涑 涑		
12	北 匕	30	肚 肚		
13	右 石	31	姐 姐		
14	叶 甲 申	32	胃 胃 昫		
15	失 矢	33	聶 聶		
16	田 由	34	毗 昆		
17	由 巳	35	彘 彘		
18	另 叻 叻	36	痲 痲		

Table 2. There are 46 groups of Chinese characters (found from 8000 characters) with same spelling sequences

	Strokes
1	一, 丿
2	丶, ㇇
3	丶, 丶, 丶
4	丨, 丨, 丨
5	丨, 丿
6	フ, フ, フ, フ
7	㇇, ㇇, ㇇

Table 3. 7 different stroke types for on-line character recognition

1	丈又太犬	26	仁年牛	51	倚勒架
2	上土士工	27	仇仇危	52	傲捷
3	不六户	28	今六本	53	歐靛
4	丑五正	29	介爪	54	兄四
5	丕丰平立	30	仍仿	55	元老
6	世占叶	31	仔竹行	56	亮祝
7	丙丙	32	佞床	57	入八
8	仞仞	33	侏侏	58	冂尸
9	俗囫	34	伊后	59	冒咄
10	凡凡尤	35	侔侔位	60	凡尤
11	主玉	36	位社社	61	由百
12	乃方	37	何向	62	刀力
13	欠欠欠	38	余余永	63	介旁
14	主朱	39	佟来采	64	切厄
15	九兀凡尤	40	侔侔空	65	刊市
16	了刁才	41	侔新	66	刊市
17	互巨	42	例初	67	列身
18	一十	43	悔帝	68	副柯固
19	亡口口	44	侶徊	69	加另叨召叻
20	交仗伏	45	促保味	70	助肯
21	亦仞	46	侏秋	71	劬周
22	亨伺	47	侏宫	72	勾材
23	京时	48	侔借	73	匸干
24	亭帝	49	侏姑	74	匸荒
25	亮柜	50	侔社	75	匸卓由音豆

Table 4. The first 75 groups of characters with same stroke combinations.

Appendix: The sets of Chinese characters with the same stroke combinations

1	丈又太犬	26	仁午牛	51	倚勅架
2	上土士工	27	仇仇庇	52	僻褲
3	不六宀	28	今未本	53	歐訖
4	丑五正	29	介爪	54	兄四
5	丕半平立	30	仍仿	55	先老
6	世占叶	31	仔竹行	56	兗祝
7	丙丙	32	伎床	57	入八
8	仞仞	33	件件	58	冂冂
9	俗囡	34	伊后	59	冒咍
10	丸凡尤	35	伴伴位	60	凡尤
11	主玉	36	住杜杠	61	凶古
12	乃方	37	何向	62	刀力
13	欠欠夫	38	余余永	63	夂夂
14	彳米	39	佟來采	64	切厄
15	九兀几尢	40	佯粹卑	65	刑帀
16	了勺才	41	併析	66	刑布
17	互巨	42	例初	67	列每
18	一十	43	侮帝	68	剗柯罔
19	亡口口	44	侶徊	69	加另叨召叻
20	交仗伏	45	促保咻	70	助肯
21	亦仞	46	俠秋	71	劓周
22	亨伺	47	倌宮	72	勾材
23	京耐	48	倅借	73	匸干
24	亭蒂	49	俯姊	74	匚荒
25	亮栖	50	倅社	75	匣卓卣昔茸

76	卞斗	101	啞豇	126	圩打
77	印芍	102	咧枷	127	丐拈
78	厝庫	103	咽莫	128	坐巫
79	廟哨	104	哀衷	129	坭屍
80	又大	105	哈桌	130	珂拈
81	叩引西	106	員唄	131	垣珥
82	叭合	107	咧涼	132	埃致
83	叮可弓	108	哺圃	133	培蚱
84	右石	109	哼師	134	堞塔
85	司吁	110	駁弔	135	錯銛
86	吉吐	111	唐砧匾	136	天天
87	吟含味東	112	售唯	137	失矢
88	叱囡	113	唾甜	138	夷奇
89	吃祀	114	啐栝粘	139	柔杵
90	呈告拒	115	唐啼	140	奶妨
91	吳茵	116	喃樽	141	如妄吐
92	吸度	117	善辜	142	姪姪
93	吻囡	118	詰罟	143	妓蔑
94	呆困杏東菜足	119	暗楫軛	144	姪妹枝
95	咒晒	120	嗝鞞	145	妾婢
96	世咕草草	121	嗇棺營	146	媪孀
97	叟姻	122	惶睡	147	姪痒
98	呷呻	123	喻翰	148	娶媪
99	命囡	124	圍害	149	婦婷
100	咎屎	125	在年	150	焚案

151	嫗嬋	176	局昞昂	201	惶惶
152	妞妞	177	居屈	202	愿甯
153	子手	178	屑翔	203	抃抖
154	孛孛	179	屏砢	204	捋拉拌
155	孿孿	180	峨嶷	205	拐招
156	冗杙	181	嶂崎	206	捏捏
157	守村	182	崐崑	207	摔圻
158	安坏	183	己巳己	208	援舡
159	宋林	184	帕帛	209	掉措
160	完罕	185	希府	210	掌第
161	宕柘	186	帙芬	211	支支
162	宙袖	187	幘景	212	救尅
163	客格	188	平立半	213	區單
164	宣桓軍	189	幸枉	214	升广
165	室桎釭	190	店茗	215	升斤
166	宰梓	191	庚奔	216	卑昕昇
167	宴晏焯蜡	192	疆彈	217	卒杯
168	宵梢	193	得曷	218	整槩
169	寂椒	194	從徠	219	斲藁
170	寄椅	195	志志	220	日日
171	富福	196	扭枉	221	旦耳
172	寐榎	197	帖柘	222	旭酉
173	寔櫻	198	悃棟	223	盱旱車
174	寵櫳	199	悼惜	224	昌盲
175	尖冰	200	憚惶	225	昔茸

226	昉昂	251	桀粧	276	评泣泮
227	昕昇	252	橙缸	277	泳涂
228	星皇重	253	桓軍	278	涑浞
229	昧是查	254	桡梳	279	渣湜
230	昱音	255	梲章	280	潼渥
231	晏焯蜡	256	梲程結	281	焯焯
232	晰皙	257	梲粘	282	焯蜡
233	暄暈暉	258	棄竦蛛	283	父爻
234	暝賂	259	棋粗	284	玖玖辰
235	朝腊開	260	棗棘	285	異豈
236	未本	261	棠棟	286	琇舒鈞
237	朱耒	262	棹喏	287	百白
238	杆辛	263	棹管	288	皇重
239	叔杖杖	264	楫軹	289	皙稗
240	杜杠	265	楫睥	290	相貞
241	杲杳	266	楸袂	291	省眇
242	枋杓	267	樟萸	292	眷眈
243	柄柄	268	槽糟	293	眸頂
244	楊桐	269	櫃輦	294	著皓
245	枸的	270	毗毘	295	睐貶
246	菓苓	271	汙汰	296	衲衲
247	柑莖草	272	沫沫	297	衲祐
248	柯罔	273	沾泄	298	禁策
249	校莢	274	沿洒	299	稗耜
250	株秣	275	汎泗	300	罨醉



301	站蚌	326	莛草	351	鄮隕
302	竝針	327	苧苻	352	鄰隣
303	竦蛛	328	菱蒞	353	醒酷
304	竹行	329	苻苻	354	鈐鉢
305	筓笳	330	草草	355	鈕鉦
306	鈞鉞	331	萹萑	356	鍾鎗
306	糊楨	332	菜菜	357	隍隍
308	禧谿	333	萱葦	358	階鄮
309	紊絞	334	蕃裸	359	馭馭
310	紉綱	335	薛藉	360	駕駟
311	紉素	336	衿袂	361	鳩亮
312	索絆	337	評証	362	麥俊
313	紉累紉	338	誅誅	363	甲由申田
314	網綱	339	誦飼	364	
315	缶舌	340	貴貼貫	365	
316	疔舩	341	踔踔	366	
317	肚肛	342	軌軌	367	
318	胃胃	343	迢迦	368	
319	胖胚	344	逞造	369	
320	脅脇	345	邢阱	370	
321	脈脣	346	邗防	371	
322	腊開	347	部陪	372	
323	腥腫	348	郵陸	373	
324	鈎鈎	349	都階	374	
325	芄芄	350	鄮隕	375	