

# 《世載堂雜憶》舊版造字轉碼說明

中研院資訊所文獻處理實驗室  
中研院史語所漢籍電子文獻工作小組  
2008/5/15 陳建安 製作

一、《世載堂雜憶》一書(世載堂雜憶.xml)使用舊版造字 113 個，字頻 567 次，詳如附件一。這 113 個造字中，92 個可轉成 Windows XP 能顯示的字，字頻 451 次；另外 21 個字必須轉成構字式，字頻 116 次。

二、附件一的造字分析表說明如下：

甲、編號：Big5 造字空間為 6217 個，編號由 1 到 6217。

乙、造字：舊版造字

丙、字頻(txt)：造字在「.txt」文件的出現次數

丁、字頻(xml)：造字在「.xml」文件的出現次數


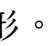
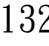
戊、Big5：造字的 Big5 碼

己、Unicode：造字所對應的 Unicode 碼

庚、WinXP：造字在 Windows XP 的對應字形

辛、構字式：Windows XP 若無對應字形，則改採用構字式

壬、備註凡例：

- 1、校對問題，舊版漢籍錯字，可用程式全部取代：在舊版漢籍電子文獻中即存在的錯字，因校對時的疏漏而未更正，持續留存在新版漢籍電子文獻中；若該造字的所有頻次，皆屬於錯誤使用的錯字情形，可以用程式全部取代為正確字形。如編號 923 的「火产兼」字，原字為「癩」。
- 2、標記問題：標記問題主要是漢籍電子文獻在舊轉新的過程中，採用了新的標記語言，而在修改標記時，對於某些標題句做了增刪的動作，因為這些標題句的內容包含了缺字，所以這些更改標記的動作造成 txt 檔與 xml 檔缺字頻次不符的情形。如編號 2130 的「憲」字，xml 檔在 55 頁新增了一個註解，因而造成頻次不同，處理方式為將所有頻次皆修改為「憲」。
- 3、異體字問題：新版漢籍考量到使用者檢索及使用時的便利性，將用字原則改為除專詞等特殊情形之外，一律改用標準字呈現。如編號 4132 的「血」係「眾」字之異體，故以「眾」字取代。又編號 4067 的「綉」字，係「繡」字之簡體字，亦以標準字的「繡」字取代。

- 4、Unicode 字型呈現差異：Unicode 字型與舊漢籍造字有些微差異，但只是字體風格差異，實際上仍為同一個字，因此仍取 Unicode 字型。如編號 4158 的「譌」字，Unicode 字型呈現為「譌」，實際上仍為同一字。
- 5、待造字：Unicode 及漢字構形資料庫皆未收錄的舊漢籍造字，正在等待補造字中，所以「造字」欄空白無法看到字形。如編號 5749 的「片𠄎庸」。

三、Unicode 目前收錄的漢字總數為 70194，分屬於三個不同區段，詳如表一。目前 Windows XP 只支援 CJK 認同表意文字區的 20902 個字，內碼為 4E00-9FFF。所以造字編號 1704 的「榑」字，Unicode 編碼為 3BAE，由於 Windows XP 並不支援，仍須使用構字式「木𠄎挈」。

表一、Unicode 的字數及編碼區段

Unicode	字集子集合	新增字數	新增編碼區段	總字數	WinXP
1.1 版	CJK 認同表意文字區	20902	4E00-9FFF	20902	支援
3.0 版	CJK 認同表意文字擴充 A 區	6582	3400-4DFF	27484	不支援
3.1 版	CJK 認同表意文字擴充 B 區	42710	20000-2A6D6	70194	不支援

四、txt 檔與 xml 檔頻次不同且造字分析表「備註」欄有備註問題情況者，詳細取代內容可參考附件二的手動取代表。附件二的手動取代表的欄位說明如下：

- 甲、編號：Big5 造字空間為 6217 個，編號由 1 到 6217。
- 乙、造字：舊版造字字形。
- 丙、檔案原文摘錄：摘錄檔案(世載堂雜憶.xml)中包含該造字的文句段落，其中紅字底線的部分，為需要手動取代的原文。文句後方標示(全部取代)者，表示執行全部取代的動作，不一一列舉所有原文。
- 丁、手動取代結果：變更過後的檔案內容，其中藍字底線的部分，為修改後的結果。文句後方標示(全部取代)者，表示執行全部取代的動作，不一一列舉所有原文。

附件一、《世載堂雜憶》造字分析表

編號	造字	字頻 (txt)	字頻 (xml)	Big5	Unicode	WinXP	構字式	備註
296	𠄎	11	11	FBEC			𠄎	
297	𠄏	6	6	FBED			𠄏	
298	𠄐	2	2	FBEE			𠄐	
299	𠄑	3	3	FBEF			𠄑	
300	𠄒	3	3	FBF0			𠄒	
308	∞	1	1	FBF8			∞	
449	斤	1	1	FCE8			斤	
628	𠄓	5	5	FDFE			𠄓	
665	𠄔	1	1	FE64			𠄔	
717	𠄕	1	1	FEBA			𠄕	
752	𠄖	2	2	FEDD			𠄖	
758	𠄗	1	1	FEE3			𠄗	
923		1	1	8EEB			𠄘火产兼𠄙	𠄘，校對問題，舊版漢籍錯字，可用程式全部取代。
1167	網	1	1	90A5	7D97	網		
1354	禩	1	1	91C3	79A9	禩		
1565	絲	1	1	92F9	7DDC	絲		

編號	造字	字頻 (txt)	字頻 (xml)	Big5	Unicode	WinXP	構字式	備註
1704	榑	1	1	93E7	3BAE		木𠂇挈	
1871	蟲	2	2	94F1	87C1	蟲		
1900	桂	1	1	954F	88BF	桂		
2046	踪	1	1	9644	8E2A	踪		
2050	蹠	1	1	9648	8E4F	蹠		
2130	窻	1	2	96BA	6119	窻		窻，標記問題，人工取代。
2142	冢	1	1	96C6	7758	冢		
2151	冫	1	1	96CF	961D	冫		
2274	鎡	1	1	97AD	9348	鎡		
2374	冢	2	2	9852	51A1	冢		塚，異體字問題。
2570	𠂇	1	1	9979	21C08		𠂇𠂇兀	
2806	勒	1	1	9AEA	976D	勒		
3789	冲	2	2	8154	51B2	冲		沖，異體字問題。
3799	刂	1	1	815E	5226	刂		劫，異體字問題。
3804	効	6	6	8163	52B9	効		
3805	勅	1	1	8164	52C5	勅		
3814	却	9	9	816D	5374	却		卻，異體字問題。
3816	厓	1	1	816F	5393	厓		

編號	造字	字頻 (txt)	字頻 (xml)	Big5	Unicode	WinXP	構字式	備註
3829	咏	1	1	817C	548F	咏		
3835	噍	1	1	81A4	564D	噍		
3840	坂	2	2	81A9	5742	坂		
3842	堃	1	1	81AB	5803	堃		
3846	壻	8	8	81AF	58FB	壻		Unicode 呈現差異。
3855	孽	4	4	81B8	5B7C	孽		孽，異體字問題。
3858	寃	5	5	81BB	5BC3	寃		寃，異體字問題。
3860	尠	2	2	81BD	5C20	尠		
3864	峯	16	16	81C1	5CEF	峯		峰，異體字問題。
3883	徧	2	2	81D4	5FA7	徧		
3886	悞	2	2	81D7	609E	悞		
3891	慙	1	1	81DC	6159	慙		慙，異體字問題。
3909	擡	2	2	81EE	64E1	擡		
3921	皙	7	7	81FA	6673	皙		皙，異體字問題。
3929	黎	4	4	8243	68C3	黎		
3930	憵	7	7	8244	6901	憵		
3966	烟	30	30	8268	70DF	烟		煙，異體字問題。
3969	煊	19	19	826B	714A	煊		

編號	造字	字頻 (txt)	字頻 (xml)	Big5	Unicode	WinXP	構字式	備註
3974	牀	8	8	8270	7240	牀		
3990	甄	2	2	82A2	750E	甄		
3996	疎	1	1	82A8	758E	疎		
4000	疴	2	2	82AC	75B4	疴		
4013	盃	1	1	82B9	76CC	盃		
4016	着	15	15	82BC	7740	着		著，異體字問題。
4034	穉	2	2	82CE	7A49	穉		
4036	窰	2	2	82D0	7AB0	窰		
4043	笄	1	1	82D7	7B53	笄		
4067	綉	2	2	82EF	7D89	綉		繡，異體字問題。
4077	罇	4	4	82F9	7F47	罇		
4080	羣	37	37	82FC	7FA3	羣		群，異體字問題。
4081	翱	34	34	82FD	7FFA	翱		
4090	脚	4	4	8347	811A	脚		腳，異體字問題。
4094	館	1	1	834B	8218	館		
4098	芪	1	1	834F	82AA	芪		
4105	菓	1	1	8356	83D3	菓		
4111	莼	25	25	835C	8493	莼		

編號	造字	字頻 (txt)	字頻 (xml)	Big5	Unicode	WinXP	構字式	備註
4129	虬	4	4	836E	866C	虬		
4132		63	63	8371			血𠄎𠄎	眾，異體字問題。
4134	衛	38	38	8373	885E	衛		衛，異體字問題。
4139	袴	1	1	8378	88B4	袴		
4142	褻	2	2	837B	8943	褻		
4146	覩	8	8	83A1	89A9	覩		
4158	譌	2	2	83AD	8B4C	譌		Unicode 呈現差異。
4173	躄	1	1	83BC	8EAD	躄		
4184	迹	11	11	83C7	8FF9	迹		
4193	鈎	3	3	83D0	920E	鈎		
4213	鞞	1	1	83E4	97BE	鞞		
4220	穎	1	1	83EB	9834	穎		穎，異體字問題。
4234	羸	1	1	83F9	9A58	羸		
4237	髡	1	1	83FC	9AE0	髡		髡，異體字問題。
4242	鬪	3	3	8442	9B2D	鬪		Unicode 呈現差異。
4247	鮎	2	2	8447	9B8E	鮎		
4256	鷄	5	5	8450	9DC4	鷄		
4258	麇	7	7	8452	9E90	麇		

編號	造字	字頻 (txt)	字頻 (xml)	Big5	Unicode	WinXP	構字式	備註
4259	麇	1	1	8453	9E95	麇		
4267	鼈	1	1	845B	9F08	鼈		
4278	亘	4	4	8466	4E98	亘		
4308	刼	1	1	84A6	5227	刼		劫，異體字問題。
4409	回	1	1	854C	518B	回		
4431	逵	2	2	8562			辵△彖	
4515	呪	4	4	85D8	546A	呪		
4702	場	1	1	86F6	5872	場		場，異體字問題。
4768	妬	1	1	8779	59AC	妬		妒，異體字問題。
4940	韵	2	2	88AA	97F5	韵		
4996	斲	2	2	88E2	65B5	斲		
5176		5	5	89F9			冫△冗	沉，異體字問題。
5220	毌	2	2	8A66	6BCC	毌		
5290	携	1	1	8ACE	643A	携		攜，異體字問題。
5298	昨	2	2	8AD6	65FF	昨		
5410	枌	1	1	8BA9	67AC	枌		
5437	并	1	1	8BC4	5E77	并		并，異體字問題。
5472	瓦	3	3	8BE7			瓦	瓦，異體字問題。



編號	造字	字頻 (txt)	字頻 (xml)	Big5	Unicode	WinXP	構字式	備註
5514	𧯛	2	2	8C52	25095		執𧯛皿	𧯛，異體字問題。
5551	賚	3	3	8C77	8CF7	賚		
5556	隸	7	7	8C7C	96B7	隸		隸，異體字問題。
5559	欸	2	2	8CA1	6B35	欸		欸，異體字問題。
5666	汚	7	7	8D4D	6C5A	汚		汚，異體字問題。
5706	敍	23	23	8D75	654D	敍		敍，異體字問題。
5749		1	1	8DC2			片𧯛庸	待造字。

附件二、《世載堂雜憶》手動取代表

編號	造字	檔案原文摘錄	手動取代結果
2130	憲	<u>闡</u> (全部取代)	<u>憲</u> (全部取代)