

「閩客語典藏計畫」—《廈英辭典》之缺字處理說明

中研院資訊所文獻處理實驗室
中央研究院語言所文獻語料小組
2008/12/12 丁玫伶

壹、《廈英辭典》的缺字狀況說明：

「閩客語典藏計畫」為中央研究院語言學研究所「漢語典藏與典藏架構」的五個子計畫之一，擬收錄《廈英辭典》等辭典文獻建立閩南語、客家語語料庫。文獻首先以外包打字，之後經過助理們二次電腦螢幕逐字校對，未來將放上網路供檢索之用。在文獻數位化過程中，若遇到電腦無法呈現的缺字，即以構字式¹或符號來表示。由語言所助理蔡瑋芬提供了一校完成的《廈英辭典》word電子文獻檔，目前的缺字表達方式，共有以下兩種：

1. 缺字大部份使用構字式表示。例：「艇」字用「骨 Δ 廷」表示。
2. 《廈英辭典》檔案中以「●」符號表示無法辨認或無法用構字式組合之字。

因此此次工作是處理《廈英辭典》的缺字表達方式，使其可以無誤地呈現在二校電子文獻檔。

貳、《廈英辭典》的缺字問題與解決方式

缺字處理是使用程式分析電子文獻檔中的構字式，與漢字構形資料庫比對，整理出漢字構形資料庫未收之字，再與原書進行校對，確認缺字字形是否正確，最後進行造字。處理過程中發現的問題，說明如下：

一、構字式的問題

構字式的表達有一定的準則，而《廈英辭典》缺字的構字式存在

¹ 構字式為中研院資訊所文獻處理實驗室開發，有關構字式的表達與使用方式，請參考〈[構字式的處理技巧](#)〉一文。

以下幾點問題：

1. 構字式使用 Unicode 部件

由於漢字構形資料庫目前無法支援 Unicode，因此無法判讀構字式中的 Unicode 部件，構字式的 Unicode 部件應改為 Big5 部件。例：檔案中「剷」字的構字式「疒△兆」的「疒」為 Unicode 部件，編碼為 7592，修正為 Big5 編碼 8BCC 的「疒」。

2. 未正規化的構字式

構字式的部件與組成有一定的標準，因此有以下的問題時需進行正規化：

(1) 可成字的部件卻被拆分為更小單位的字根

構字式以使用成字的部件為優先，避免將成字部件拆分為更小單位的字根，因此需將可成字的字根合併表達。例：「罍」字在檔案中的構字式原為「◻口口尹◻」，而「口」和「口」可合為「𠔁」字，因此構字式修改為「𠔁△尹」。

(2) 構字部件採用異體字根

檔案中有些構字部件採用異體字根，需改為正確字根。例：檔案中的構字式「執△火」，比對原書後，「火」應改為異體字根的「灬」，構字式為「執△灬」。

(3) 手誤之構字符號

有些構字符號因手誤打錯，需修改為正確的構字符號。例：「𠔁△遭」比對原書字形後，構字式中的直連符號應改為橫連符號「𠔁△遭」。

二、異體字問題

《廈英辭典》有不少缺字是漢字構形資料庫內字形的異體字，因此在處理缺字時，會有該使用原書之字形，或以漢字構形資料庫已收之字取代的考量出現，因此處理異體字時最主要的方式是，如只是細微的差異，皆以漢字構形資料庫之字形取代，不再另外造字；而字典

中特別列出的異體字則保留原字形。例：原書缺字字形為「鹿△章」，而漢字構形資料庫已收有「麇」字，兩者只有「鹿」字是直連或包含的細微差異，因此以漢字構形資料庫之「麇」字取代，不再另外造字。

三、錯字問題

校對《廈英辭典》缺字時，發現檔案之字與書上字形明顯不同，並查詢教育部異體字字典確認非異體字，皆歸類為錯字。例：檔案中使用「四△會」表達缺字，而原書之字為「罾」，「四△會」為錯字，以正確字形「罾」取代。

四、「●」缺字問題

《廈英辭典》的「●」符號代表無法辨認或無法用構字式組合之字，經比對原書後，如為異體字，則改成漢字構形資料庫已收之字，而未收的字是為待造字，將進行造字。例：廈英 S 檔案中的「●」符號，比對原書後為「𠄎△絲」字，是資料庫未收之待造字，將進行造字。

參、《廈英辭典》缺字處理流程

《廈英辭典》缺字處理的詳細流程如下：

一、缺字分析與構字式修改

(一) 使用程式分析《廈英辭典》檔案，找出文中所有的構字式，建立《廈英辭典》缺字資料庫。

(二) 利用漢字構形資料庫比對《廈英辭典》缺字資料庫，找出漢字構形資料庫內沒有的構字式，並作以下處理：

1. 修改 Unicode 部件。例：「疒△兆」的「疒」修正為「疒△兆」。
2. 正規化構字式，需修改的部份如下：

(1) 正規化構字組合。例：「□口口尹□」應修改為「四△尹」。

(2) 正規化異體字根。例：「執~~火~~火」實則為「執~~火~~...」，「火」應改為異體字根的「...」。

(三) 以漢字構形資料庫再次比對修改後的構字式，確認漢字構形資料庫未收錄的缺字將與原書進行校對。

二、缺字校對

(一) 《廈英辭典》資料庫中，未收在漢字構形資料庫的缺字構字式與原書一一進行校對，需修改的部份如下：

1. 修改資料庫內錯誤的缺字字形。例：錯誤的缺字字形「四~~會~~會」改為正確的「晉」字。
2. 修改資料庫缺字的異體字字形，使用漢字構形資料庫已收字形。例：「鹿~~章~~章」以漢字構形資料庫已收之「麀」字的「鹿~~章~~章」取代。

(二) 將「●」缺字與原書一一進行校對，如為異體字，則改成漢字構形資料庫已收之字，未收的字將進行造字。

三、修改 word 檔

修改 word 文本上的缺字構字式：

(一) 以《廈英辭典》資料庫修正之結果，修改 word 檔案中構字式的 Unicode 部件為 Big5，並把改過之部件製作成「修正 Unicode 部件」清單，列出修改過的 Unicode 部件，完整清單請參考附件二，欄位說明如下：

1. 原構字式：缺字在檔案中的構字式或字形。
2. Unicode 部件：缺字在檔案中使用的 Unicode 部件。
3. 修正構字式：修改後的構字式。
4. 頻次：缺字出現的次數。
5. 檔名：缺字所在的檔案名稱。

(二) 以《廈英辭典》資料庫修正之結果，正規化 word 檔案中未正

規化之構字式、修改錯誤的缺字構字式和異體字形。將改過的構字式列成「修正構字式」清單，並在備註欄寫上修改原因，以供語言所參考。修改過之字頻共計 11 字，完整清單請參考附件三，欄位說明如下：

1. 頁碼：缺字在原書的頁碼。
2. 原構字式：缺字在檔案中原來的構字式或字形。
3. 修正構字式：修改後的構字式。
4. 頻次：缺字出現的次數。
5. 檔名：缺字所在的檔案名稱。
6. 備註凡例：

(1) **錯字**：與原書字形或字義不符，並查詢教育部異體字字典確認非異體字後，皆歸類為錯字，手動修改為正確字形。修改過之字頻共計 2 字。例：檔案使用「四~~公~~會」表達缺字，而原書使用之字為「罍」，且兩者並非異體字，手動修改為正確字形「罍」。

(2) **異體字，使用漢字構形資料庫之字形**：缺字與漢字構形資料庫已收之字形是異體字，只有細微的差異，皆以漢字構形資料庫之字形取代，不再另外造字。修改過之字頻共計 2 字。例：原書缺字字形為「鹿~~公~~章」，而漢字構形資料庫已收有「麇」字，兩者只有「鹿」字是直連或包含的細微差異，因此以漢字構形資料庫之「麇」字取代，不再另外造字。

(3) **正規化構字式**：檔案使用未正規化的構字式表達缺字時，需修改為漢字構形資料庫正規化的構字式。修改過之字頻共計 5 字。例：「執~~公~~火」正規化成「執~~公~~灬」，「火」改為異體字根的「灬」。

(三) 《廈英辭典》檔案中出現「●」字形，共計 17 字，比對原書

後以字形或構字式手動取代，漢字構形資料庫未收的缺字將進行造字，待造字共計 5 字，並製作「●」字修正表，請參考附件四。欄位說明如下：

1. 頁碼：「●」字在原書的頁碼。
2. 檔案原文摘錄：摘錄檔案「●」符號。
3. 原書字形：「●」在原書中之字形，以此字取代檔案原文的「●」。
4. 檔名：「●」字所在的檔案名稱。
5. 備註：記錄修改相同詞句的次數或是否為待造字。

五、《廈英辭典》缺字進行造字

確認《廈英辭典》缺字資料庫無誤後，與漢字構形資料庫比對，未收的缺字將進行造字，完成後收至漢字構形資料庫。《廈英辭典》缺字共計 208 字，其中漢字構形資料庫已收錄 176 字，待造字共 32 字，將進行造字。為了未來修改的方便性，檔案內缺字仍維持構字式形態，並製作一份包含構字式的造字記錄表供語言所參考。請參考[《廈英辭典》缺字造字記錄](#)。

六、確認《廈英辭典》檔案的缺字處理完成

1. 執行word「工具」中的「構字式轉字形」²，確認檔案修正後的構字式，可否轉為字形。
2. 無法轉換字形之字，需確認構字式是否正確，再比對漢字構形資料庫，未收字將進行造字。

² 如何在 word 文本上看見缺字字形，請參考附錄一

附件一、《廈英辭典》修正Unicode部件表

原構字式	Unicode 部件	修正構字式	頻次	檔名
疒△兆	疒	疒△兆	2	廈英TH
疒△烝	疒	疒△烝	1	廈英Ch
艹△冠	艹冠	𠂇 ^艹 完女 ^口	1	廈英H
艹△耿	艹	艹△耿	1	廈英H
艹△附	艹	艹△附	1	廈英H
忄△勿	忄	忄△勿	1	廈英H
卩△啗	卩	卩△啗	1	廈英H
艹△隔	艹	艹△隔	1	廈英K
艹△翹	艹	艹△翹	1	廈英K
疒△狂	疒	疒△狂	1	廈英K
疒△交	疒	疒△交	1	廈英K
虫△量	量	虫△量	1	廈英K
艹△貝	艹	艹△貝	1	廈英P
艹△實	艹	艹△實	1	廈英S

附件二、《廈英辭典》修正構字式清單

頁碼	原構字式	修正構字式	頻次	檔名	備註
p. 18	口口尹	𠂇 ^口 尹	2	廈英B	正規化構字式，

頁碼	原構字式	修正構字式	頻次	檔名	備註
					「口△口」合為「囗」。
p. 119	興目△	興△目	1	廈英 H	正規化構字式，構字符號位置錯誤。
p. 46	足△遭	足△遭	1	廈英 Ch	正規化構字式，構字符號「△」改為「△」。
p. 181	執△火	執△灬	2	廈英 J	正規化構字式，「火」改為「灬」。
p. 225	鹿△章	鹿△章	1	廈英 K	異體字，使用漢字構形資料庫之字形。
p. 575	𠄎月先先日 𠄎	𠄎月𠄎日	1	廈英 TS TSH	正規化構字式，「先△先」合為「𠄎」。
p. 576	四△會	𠄎	1	廈英 TS TSH	錯字，四改為𠄎，會改為曾。
p. 199	想△共	相△共	1	廈英附 錄 KH-U	錯字，想改為相。
p. 240	言△党	讜	1	廈英附 錄 KH-U	異體字，使用漢字構形資料庫之字

頁碼	原構字式	修正構字式	頻次	檔名	備註
					形。

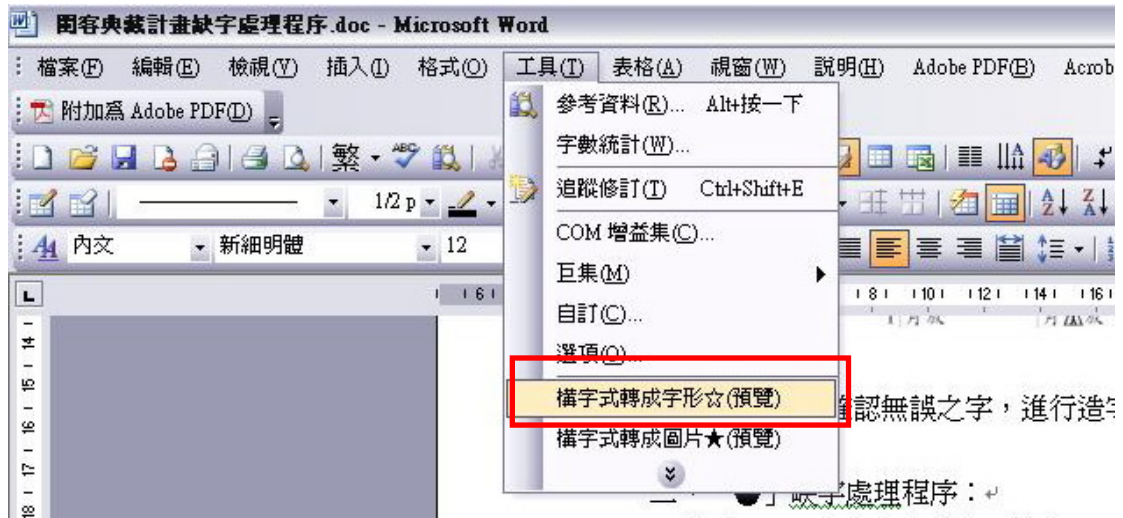
附件三、《廈英辭典》「●」字修正表

頁碼	檔案原文摘錄	原書字形	檔名	備註
p. 69	●	𠄎 𠄎 膝	廈英 Chh	
p. 91	●	𠄎 木 戶 𠄎 一 夕 𠄎 𠄎	廈英 Chh	待造字
p. 186	●	𠄎 女 𠄎 口 𠄎	廈英 J	待造字
p. 228	●	藁	廈英 K	
p. 252	●	茗	廈英 K	
p. 258	●	籬	廈英 Kh	
p. 343	●	鼯	廈英 N	
p. 357	●	瞶	廈英 P	
p. 381	●	金 𠄎 友	廈英 P	改為鉞
p. 402	●	米 𠄎 𠄎	廈英 P	
p. 419	●	𠄎 𠄎 絲	廈英 S	待造字
p. 449	●	櫻	廈英 S	
p. 489	●	目 𠄎 𠄎	廈英 T	
p. 166	●	瞶	廈英附錄 KH-U	
p. 176	●	黼	廈英附錄 KH-U	

頁碼	檔案原文摘錄	原書字形	檔名	備註
p. 214	●	𠄎立	廈英附錄 KH-U	待造字
p. 257	●	魚𠄎	廈英附錄 KH-U	

附錄一：如何在 word 檔案上看見缺字字形

一、安裝漢字構形資料庫後，點選word「工具」表單內的「構字式轉成字形」(如下圖)，即可呈現出漢字構形資料庫內已有之完整字形。其他使用word的相關問題，請參考「[文獻處理實驗室](#)」的「[缺字處理計畫](#)」網頁。



二、更新版本或下載新增字檔

如使用 word「工具」表單內的「構字式轉成字形」，卻無法看到新造的字形，有以下兩種解決方式：

1. 使用的漢字構形資料庫為舊版，無法看到新造字形，可以到「[缺字處理計畫](#)」網頁更新至最新版的漢字構形資料庫，目前最新的為 2.52 版。
2. 如已是最新的版本，卻仍無法看見新造之字，可至「[缺字處理計畫](#)」網頁的「[行事曆](#)」下載更新的造字檔。例：下圖為 2.51 版的更新字檔，需下載「檔名」欄內的檔案，更新漢字構形資料庫。

下載漢字構形資料庫2.51版更新檔案

請參閱下載說明

名稱	檔名	檔案大小	更新原因	更新日期
楷書構形資料庫	cdphanzi.mdb	20.1M	增加95個新字	2008年2月12日
標楷體外字集一	hzcdp01k.ttf	5.0M	修改1個字形	2008年2月12日
細明體外字集一	hzcdp01m.ttf	3.9M	修改1個字形	2008年2月12日
標楷體外字集二	hzcdp02k.ttf	5.8M	修改1個字形	2008年2月1日
細明體外字集二	hzcdp02m.ttf	4.5M	修改1個字形	2008年2月1日
標楷體外字集六	hzcdp06k.ttf	6.8M	修改1個字形	2008年2月1日
細明體外字集六	hzcdp06m.ttf	5.0M	修改1個字形	2008年2月1日
標楷體外字集七	hzcdp07k.ttf	6.9M	修改1個字形	2008年2月1日
細明體外字集七	hzcdp07m.ttf	5.4M	修改1個字形	2008年2月1日
標楷體外字集八	hzcdp08k.ttf	5.7M	修改2個字形	2008年1月17日
細明體外字集八	hzcdp08m.ttf	4.3M	修改2個字形	2008年1月17日
標楷體外字集九	hzcdp09k.ttf	646K	增加95個新字	2008年2月12日
細明體外字集九	hzcdp09m.ttf	520K	增加95個新字	2008年2月12日
漢字構形資料庫2.51	cdphanzi.exe	904K	修正程式錯誤	2008年2月12日