

文章编号: 1003-0077(2008)04-0109-05

基于 DUCET 的藏文排序方法

黄鹤鸣¹，契嘎·德熙嘉措²（赵晨星）

(1. 青海师范大学 物理系,青海 西宁 810008; 2. 青海藏文信息研究所,青海 西宁 810008)

摘 要: DUCET 为每个藏文字符规定了排序码,但藏文音节的拼写复杂性使得藏文排序不能直接应用这些排序码,提出了基于 DUCET 的藏文音节排序方法,主要思想是:首先,将二维的藏文音节转化成一维的字母串;其次,从 DUCET 中查出每个字母的排序码,得到藏文音节对应的排序码串;最后,通过比较排序码串实现藏文音节间的排序。还讨论了藏文音节与一般藏文字母串以及藏文字符串与外文字符串间的比较规则。

关键词: 计算机应用; 中文信息处理; 藏文字符串; 藏文音节; DUCET; 排序

中图分类号: TP391.1

文献标识码: A

A DUCET-based Tibetan Sorting Algorithm

HUANG He-ming¹, ZHAO Chen-xing²

(1. Physics Department , Qinghai Normal University , Xining , Qinghai 810008 , China ;

2. Qinghai Institute of Tibetan Information and Technology, Xining, Qinghai 810008, China)

Abstract : DUCET (Default Unicode Collation Element Table) is an international standard of character collation. This paper proposes a method of DUCET-based Tibetan sorting algorithm. It first expands two-dimensional Tibetan scripts into a one-dimension string of Tibetan letters. Then it locates the collation code of each Tibetan letter from DUCET. Finally, by comparing any two distinctive collation code strings, including Tibetan scripts and non-Tibetan scripts, a correct DUCET-based Tibetan order will be achieved.

Key words : computer application ; Chinese information processing ; Tibetan scripts ; DUCET ; Tibetan strings ; collation

从音节层面上看，藏文音节串（例如：འགྲོ་བཤིན་པུས་ལ།）的排序和英文、汉文字符串的排序相同：一个音节串的顺序由串中各个音节依次决定。但和汉字或者英文字母不同的是，每个藏文音节的顺序还需由构成它的各个字母依次决定。实际上，藏文字符排序的关键就是将音节中各个字母对音节顺序的影响用简洁通用的算法表示出来，如果这个问题解决了，藏文字符排序问题就解决了，因为音节层面上的排序可以利用汉文或英文字符串的排序算法。

和汉文或者英文字符排序比起来,藏文排序较烦琐,因为:虽然对于结构完整的现代藏文音节

而言,字典顺序是由基本辅音、前加辅音、上加辅音、下加辅音、元音、后加辅音和又后加辅音依次决定的,但对于一般现代藏文音节而言,不能直接应用这个规则来排序; 和现代藏文音节比起来,梵音藏文音节在字母系统、拼写规则等方面差异较大,排序也相应较困难; 藏文语法(例如:格助词)会影响相关音节的拼写。

到目前为止,讨论藏文音节排序的文献主要有文献[1~3],其中文献[1]由于没有考虑像 ཅེ 这种有上加辅音而没有前加辅音的音节在字典顺序中的特殊性,排序结果不完全符合字典顺序;文献[2]、[3]通过对藏文字典顺序的分析,建立了藏文音节的

收稿日期：2007-10-18 定稿日期：2008-01-14

基金项目：信息产业部电子信息产业发展基金资助项目（信部运[2002]393号）

作者简介: 黄鹤鸣(1969—),男,硕士,副教授,研究方向为藏文信息技术,模式识别;契嘎·德熙嘉措(1946—),男,教授,硕士生导师,研究方向为藏文信息技术。

