

藏文音节规则模型及应用

珠杰^{1,2,†} 李天瑞¹ 格桑多吉² 仁青诺布² 乔少杰¹

1 西南交通大学信息科学与技术学院, 成都 610031; 2 西藏大学计算机科学系, 藏文信息工程研究中心, 拉萨 850000;
† E-mail: trocky.jie@gmail.com

摘要 藏文音节具有独特的构造方法, 不同的构造位上有不同的藏文字符。根据藏文字符不同的组合, 构成了千变万化的藏文音节。由于字母的语音特性, 藏文组合形式上有诸多限制。作者以藏文音节为研究对象, 借助藏文语法规则, 建立了现代藏文音节的简化模型和相应的规则库, 介绍了其应用领域, 提出了一种基于音节模型的藏文音节自动拼写算法, 实验验证了规则方法的有效性。

关键词 藏文音节; 藏文规则; 规则库; 音节拼写

中图分类号 TP391

Tibetan Syllable Rule Model and Applications

ZHU Jie^{1,2,†}, LI Tianrui¹, GE Sangduoji², REN Qingnuobu², QIAO Shaojie¹

1. School of Information Science and Technology, Southwest Jiaotong University, Chengdu 610031;
2. Department of Computer Science, Tibetan University, Lhasa 850000; † E-mail: trocky.jie@gmail.com

Abstract Tibetan syllable has a unique configuration approach. Different positions of configuration compose different Tibetan characters. Based on the combination of Tibetan characters, there are a large number of ever-changing Tibetan syllables. Because of the pronunciation features of the letters, there are several limitations of the styles for Tibetan combination. This paper focuses on the study of Tibetan characters. A simplified model for modern Tibetan syllables and its corresponding rule base are established by using Tibetan grammar rules. Its applications are extensively analyzed. Algorithms for automatic spelling of Tibetan syllable are proposed as well. Experiments and case studies validate the rule base of Tibetan.

Key words Tibetan syllable; rule of Tibetan; rule base; syllable spelling

藏文音节作为构词的成分, 有其自身的特征, 特别是字母组合上有很多拼写规则。从书面藏文的信源属性来看, 藏文文本中的音节有 72% 的冗余度, 只有 28% 是可自由选择的^[1], 这说明藏文音节中 3/4 的藏文音节是保证依据语法规则来拼写的。在藏文信息处理的应用上, 藏文规则在藏文排序中有重要的应用价值。江荻等^[2]利用藏文语法提出了构造序、构造级(拼写序)和字符序相结合的排序算法, 建立藏文排序模型, 为藏文排序在计算机中的实现提供了理论基础。Chilton 等^[3]利用藏文规则, 对藏文编码国际标准 ISO/IEC 10646 音节进行了排序, 通过

“collation element”的概念, 建立一个“collation element”表, 该表通过对藏文规则建立权重分级来解决藏文的排序问题, 虽然需要排序的“字符”数量多了许多, 但是算法简单并易于实现, 在 Mysql、MIMER SQL 和 OpenOffice 2.0 等系统得到了成功应用, 但是该方法没有从藏文规则模型的角度来进行讨论。本文依据藏文语法探讨了藏文规则的数学模型, 并建立藏文规则库, 然后将其应用到藏文音节自动拼写和拼写检查等领域中, 说明由规则模型建立的藏文规则方法能够有效解决藏文信息处理研究中的若干基础性问题。

国家自然科学基金(61165013, 60763010, 61100045, 61262058)和西藏自治区自然科学基金(2010 年 41 号)资助

收稿日期: 2012-06-05; 修回日期: 2012-10-01; 网络出版时间: 2012-10-26 17:54

网络出版地址: <http://www.cnki.net/kcms/detail/11.2442.N.20121026.1754.016.html>

1 藏文音节结构

藏文音节结构以基字为核心,既有横向拼写又有纵向拼写,前加字、基字、后加字和再后加字是横向拼写,上加字、基字、下加字和元音符是纵向拼写,因此具有十分复杂的音节结构。字符在音节中的特定位置可以称为“构造位”。根据藏文的语法,各个构造位上出现的字符的性质与数量均有一定的限制,相互之间形成一种约束关系。

藏文音节中不包括梵音转写藏文,藏文音节的基本结构中构造位共有7个,如图1所示。

定义 1^[4] 构造位上的字符称为构件,根据不同位置分别称为前加字、上加字、基字、下加字、元音、后加字和再后加字,如图2所示,我们称之为藏文音节模型-1(简称模型-1)。

每个构造位在藏文音节中的表示为:1是前加位,2是上加位,3是基字位,4是下加位,5是元音位,6是后加位,7是再后加位,分别由前加字、上加字、基字、下加字、元音、后加字和再后加字来表示在字中的位置。

定义 2^[5] 在一个音节中的纵向单位(上下叠加的组合体)叫字丁或叠加字符。例如:“ལྷོ”,“ལྷོ”称为叠加字符。

根据藏文语法,有5个前加字,分别是ག, ད, བ, མ和འ; 3个上加字,分别是ར་མགོ།, ལ་མགོ།和ས་མགོ།; 30个基字是30个藏文字母; 4个下加字,分别是ལ་བདགས།, ར་བདགས།, ལ་བདགས།和ཕ་བདགས།; 5个元音符号中,有3个是上元音符号,分别是ཨ།, ཨེ།, ཨོ།, 1个下元音符号ཨ།, 一个隐含元音ཨ།; 10个后加字,分别是ག, ད, ད, ཅ, བ, མ, འ, ར, ལ, བ; 2个再后加字,分别是ད, ས།。结合模型-1,对每一个构造位

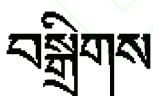


图1 藏文音节的基本结构
Fig. 1 Structure of Tibetan syllable

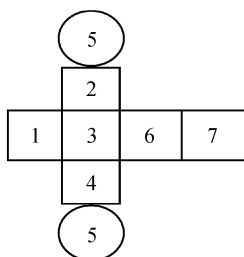


图2 构造位及模型-1
Fig. 2 Model-1 and position of construct

的元素集合描述如下。

设基字集合用 B 来表示:

$$B = \{ག, ར, ཇ, ད, ཅ, ཆ, ཇ, ཉ, ཏ, ཐ, ད, ཅ, བ, མ, འ, ར, ལ, བ, མ, འ, ཅ, ཆ, ཇ, ཉ, ཏ, ཐ, ད, ཅ, བ, མ, འ, ར, ལ, བ, མ, འ, ཅ, ཆ, ཇ, ཉ, ཏ, ཐ\}$$

设前加字集合用 Pr 来表示:

$$Pr = \{ག, ད, བ, མ, འ\}$$

设上加字集合用 U 来表示:

$$U = \{ར་མགོ།, ལ་མགོ།, ས་མགོ།\}$$

设下加字集合用 D 来表示:

$$D = \{ལ་བདགས།, ར་བདགས།, ལ་བདགས།, ཕ་བདགས།\}$$

设后加字集合用 S 来表示:

$$S = \{ག, ད, ད, ཅ, བ, མ, འ, ར, ལ, བ, མ\}$$

设再后加字集合用 SS 来表示:

$$SS = \{ད, ས\}$$

设元音字符集合用 $Tvowel$ 来表示:

$$Tvowel = \{\emptyset, ཨ, ཨེ, ཨོ, ཨ\}, \text{ 其中 } \emptyset \text{ 表示隐含元音字符。}$$

2 藏文音节规则模型

2.1 藏文音节模型的建立与简化

模型 3 是根据藏文的音节结构建立的一个模型,该模型中以基字为核心,在元音和后加字的作用下构成一个音节,在实际写法中除了构造位 3 不能空之外,其余位置均可以为空。当一个基字构成音节时,该音节隐含了元音“ཨ”和后加字“འ”。模型-1的笛卡尔积如下式:

$$Pr \times U \times B \times D \times Tvowel \times S \times SS = \{ \langle p, u, b, d, v, s, ss \rangle | p \in Pr, u \in U, b \in B, v \in Tvowel, s \in S, ss \in SS \}$$

藏语的语音理论体系中将藏语语音分为元音和辅音。根据藏文的语音特性,对于30个辅音字母进行了字性分类,为阳性、中性和阴性3种,其中阴性包括准阴性、极阴性和纯阴性3种,共计5种分类。辅音字母中提取出来的前加字、后加字构件也进行了上述5种的分类。根据每个构件的发音特性,字母组合上有很多拼写限制。从上节中的描述可知,模型-1中的各个构造位的元素是有限的,可以看出经过适当的迭代,藏文全部音节能够被拼写出来,但也产生了很多不符合语法的冗余音节。为了不产生冗余音节,我们通过预组合方式对模型-1进行简化。

假设 1 根据上加字与基字、下加字与基字、上加字+基字+下加字的固定组合关系,模型-1中的

表 1 藏文规则表
Table 1 Table of Tibetan syllable rules

ག	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ	
གསལ་ཕྱིར་དུ་																						
ར་མགོ་ཅན།	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ
ལ་མགོ་ཅན།	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ
ས་མགོ་ཅན།	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ
ཡ་བརྟགས་ཅན།	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ
ར་བརྟགས་ཅན།	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ
ལ་བརྟགས་ཅན།	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ
ས་བརྟགས་ཅན།*	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ	ཐ	ད	འ	ཨ	མ	ཤ	ཧ	ལ	ཚ	ཛ	ཞ	ཟ	འ
ག་མངོན་འབྱུག	གཏ	གཉ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ	གཏ
དཀ་མངོན་འབྱུག	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ	དཀ
བཀ་མངོན་འབྱུག	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ	བཀ
འ་མངོན་འབྱུག	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ	འཀ
མ་མངོན་འབྱུག	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ	མཀ
བརྟགས་པ་ཅན།	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི	ཀྱི

参考文献

- [1] 江荻. 书面藏语的熵值及相关问题 // 中文信息处理国际会议论集. 北京: 清华大学出版社, 1998
- [2] 江荻, 康才峻. 书面藏语排序的数学模型及算法. 计算机学报, 2004, 27(4): 524 - 529
- [3] Chilton R R. Sorting unicode Tibetan using a multi-weight collation algorithm[EB/OL]. [2012-08-10]. <https://collab.itc.virginia.edu/access/wiki/site/26a34146-33a6-48ce-001e-f16ce7908a6a/sorting%20tibetan.html>
- [4] 扎西次仁. 《中华大藏经·丹珠尔》藏文对勘本字频统计分析. 中国藏学, 1997, 2: 122-133
- [5] 王维兰, 陈万军. 藏文字丁、音节频度及其信息熵. 术语标准化与信息技术, 2004, 2: 27-31
- [6] 珠杰. TSRM 的藏文拼写检查算法. 软件学报, 审稿中
- [7] 才旦夏茸. 才旦夏茸全集. 北京: 民族出版社, 2007: 44-45