

基于XML语言甲骨文语料库元数据抽取的研究

吴琴霞^{1,2}, 高峰^{1,2}, 刘永革^{1,2}

(1. 安阳师范学院 计算机与信息工程学院, 河南 安阳 455002;

2. 河南省甲骨文信息处理重点实验室, 河南 安阳 455002)

摘要:计算机辅助甲骨文考释首先要解决的是以统一的数据形式来管理甲骨文信息,文中提出了一种基于XML语言的元数据抽取方法,提出了将抽取出的元数据存放到元数据库中,使用视图对元数据库进行动态地增加或删除,来弥补元数据抽取不全或错误的现象。同时通过XML Schema文档检测保存的元数据的语法格式,为以后映射到甲骨文语料库中做了语法的检测。该方法为甲骨文语料库的建设提出了一种新方案,也为后期语料信息语义化和专家系统推理打下了基础。

关键词:语料库;元数据;信息抽取;XML;XML Schema

中图分类号:TP393.01

文献标识码:A

文章编号:1673-629X(2012)05-0216-03

Study of Oracle Bone Inscriptions Corpus Metadata Extraction Based on XML

WU Qin-xia^{1,2}, GAO Feng^{1,2}, LIU Yong-ge^{1,2}

(1. School of Computer and Information Engineering, Anyang Normal University, Anyang 455002, China;

2. Henan Province Key Laboratory of Oracle-bone Inscriptions Information Processing, Anyang 455002, China)

Abstract: In order to assist expert makes explanations for oracle bone inscriptions, first of all, to be solved based on a unified data format to manage the oracle information. It presents a meta-data extraction method based on XML to extract the information is converted to a unified XML format, extracted metadata stored in metadata database, used a view on metadata database for dynamically inserting or deleting, to compensate for metadata extraction not all or the wrong phenomenon. And at the same time through XML Schema document detect grammar format of preservation metadata, is later mapped to oracle bone inscriptions corpus for the grammar of corpus detection. The method for oracle corpus built a new program, as well as laid the groundwork for the later corpus semantic information and expert system reasoning.

Key words: corpus; metadata; information extraction; XML; XML Schema

0 引言

文字是民族的文化之根,甲骨文是中国迄今发现的最早的成系统的文字,是汉字的源头。但是,在四千多甲骨文字中,已考释出的甲骨文字不到一千字。因此甲骨文考释在甲骨学研究领域一直是一个热门课题^[1]。甲骨文经过近一百年的考释,容易考释的字所剩无几,所以现在的甲骨文考释已经成了一项极富挑战性同时也非常困难的工作。用计算机辅助甲骨文考释与甲骨文专家相比具有以下优势:如海量的存储数据、严密的逻辑推理等。因此随着甲骨文信息数字化研究层次的深入,目前使用计算机去辅助甲骨文考释

未释字已经可行。语料库是指大量文本的集合,是使用计算机技术辅助研究语言学而产生的,甲骨文语料库是研究甲骨文的数字化平台,让计算机辅助甲骨文考释未释字首要任务就是甲骨文语料库的建设。语料库的建设首要任务是元素据的抽取,元数据对语料库的研究具有重要的意义。XML是万维网(W3C)推出的一种具有结构性的标记语言,可以组织数据、标记数据、存放数据和定义数据类型。它允许用户自定义标签,具有良好的可伸缩性和灵活性,XML不仅能提供对资源内容的表示,同时也能提供资源所具有结构信息。因此文中提出了一种基于XML语言的甲骨文语料库元数据的抽取技术。

收稿日期:2011-09-22;修回日期:2011-12-26

基金项目:国家自然科学基金资助项目(60875081)

作者简介:吴琴霞(1980-),女,河南扶沟人,硕士研究生,讲师,CCF会员,研究方向为语义Web、中文信息处理;刘永革,教授,研究方向为数据库与知识库系统、数据挖掘、中文信息处理。

1 XML与语料库

XML指可扩展的标记语言(Extensible Markup Language, XML)。XML是一套用来描述数据语义的标

记规则,其焦点是数据的内容。这些标签将文档分成许多部件,并对这些部件加以语义标识。XML 是描述万维网上数据的统一语言,已成为未来的 Web 语言和 Web 上数据交换的标准,XML 将会成为数据处理和数据传输的最常用工具。XML 具有良好的扩展性和灵活性,不仅能对资源内容采用统一的数据格式进行表示,同时也能对资源所具有的结构信息进行表示,XML 适合表示各种信息数据因而被广泛接受。

语料库(Corpus)是语料的仓库或者语料的集合,甲骨文语料库是针对甲骨文领域的多个应用而专门收集有一定结构可被计算机程序检索到的、具有一定规模的语料信息的集合。甲骨文语料库不仅是原始语料的集合,而是具有数据结构的、标注了语法语义和语用的语言信息的语料集合。语料库在各个语言学研究领域广泛地应用,它是语言研究数字化的重要基础,使用甲骨文语料库辅助甲骨文考释可以克服传统的甲骨文考释存在的问题,如:1)语料的真实性和客观性;2)语料所涉及的广泛性;3)知识共享等等^[2]。

采用XML语言存储语料信息具有很多优点。

第一,XML语言具有通用的国际标准,方便数据的集成、共享和交换;

第二,XML是各种应用程序之间进行数据交换的最常用的工具,可以减少语料库的程序和语料库数据的依赖性,提高语料库数据的独立性,实现机器可读;

第三,XML标记语言采用树形结构,语法形式简单,数语言工作者便于掌握,可以方便的组织和整理语料库;

第四,XML标记语言具有可扩展性的特点,允许用户自定义标签,方便元数据的添加。

2 语料仓库元数据的抽取

甲骨文拓片上保存的数据为非结构化信息,计算机无法实现自动化分析和抽取。目前为止发现的甲骨片大约有十五万片左右,信息量大并且信息结构复杂。然而甲骨片上记录的信息都是有规律可寻的,这就对实现元数据的自动抽取带来了可能性。因此可以采用XML标签语言对甲骨片上的卜辞信息进行描述,对抽取出的元数据来定义对应的XML Schema文档。

2.1 元数据

元数据(Metadata)是用于描述数据及其环境的数据^[3]。元数据一经建立,便可以共享。通过元数得出语料库中语料的分类信息、文本信息等各种各样的信息。这将为辅助甲骨文考释提供更多的可能性。元数据是对数据的一种描述方式,使用元数据可以提高系统的查全性和查准性。采用元数据以结构化的方式表达甲骨片上存放的卜辞信息,方便了信息的共享。这

为后期考释甲骨字提供丰富的语料信息。目前关于元数据的抽取和标注方法中存在一系列问题。文中提出了将抽取出的元数据存放到元数据库中,使用视图对元数据库进行增加或删除,这样方便用户自由的添加或删除元数据。元数据抽取工作分为以下两步:

1)对元数据存储。语料库设计初期,由于设计者自身的知识水平有限,元数据的抽取很难达到其合理性^[4]。因此对抽取出的元数据采用数据库中字段的信息进行存储。

2)保存模板。依据抽取出的元数据生成XML Schema文档,提供采用XML格式保存元数据的模板。

抽取的元数据可以表示为如图1的树形结构。图中的虚线框代表可以省略、*代表可以有零个或多个。依据此树形结构,可以定义出元数据对应的XML Schema文档。

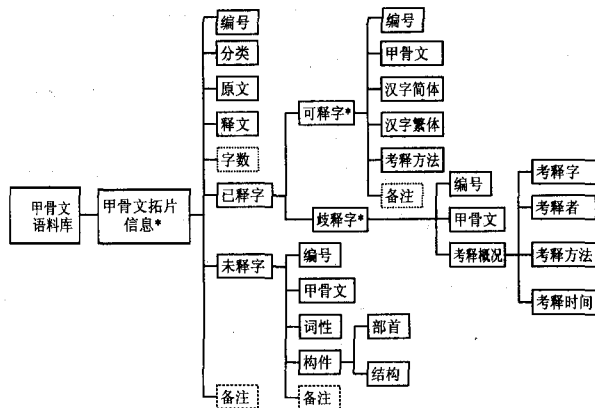


图1 元数据树型结构图

2.2 甲骨文语料库元数据抽取模型

针对甲骨拓片上的卜辞信息进行元数据抽取,依据抽取出的结果定义出XML保存元数据的语法结构即XML Schema,采用XML Schema对采用XML存储甲骨卜辞信息进行验证。最后保存为语料库(即XML存储到关系数据的存储)。元数据抽取模型如图2所示。各模块所对应的功能如下:

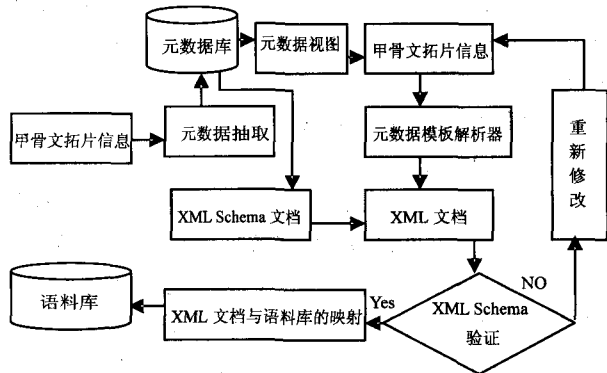


图2 基于XML甲骨文语料库元数据抽取模型

甲骨文元数据的抽取模块:输入甲骨文拓片上的信息根据语料库的需求抽取元数据,根据抽取的

所有案件。

(6)办结通知:当前审批申请者最近申办的办件经过相关某个部门的审批处理完成后,以通知形式告知申办者。

5 结束语

文中根据达县的实际情况,在分析和参考了众多现行的网上行政审批系统的基础之上,对达县网上行政审批系统进行了详细的需求分析,给出了核心业务的用例图和流程图。接着对系统进行了总体设计和数据库设计,给出了系统的体系结构图及数据库的 E-R 图,最后基于 B/S 架构,利用 ASP.NET 等技术实现了达县网上行政审批公共服务平台。该系统的实现对达县的电子政务建设和政府信息化建设有重要的现实意义,为审批系统的开发提供了参考和借鉴。

参考文献:

- [1] 任雪松. 电子政务现状与发展趋势[J]. 职业技术教育研究, 2006(4): 3-4.
- [2] 张雷. 大庆市政府电子政务行政审批系统研究[D]. 哈尔滨: 哈尔滨工程大学, 2006.
- [3] 吴静, 曹元大. 针对网上审批系统的电子政务解决方案[J]. 微机发展(现更名: 计算机技术与发展), 2004, 14(7): 24-27.
- [4] 郑刚. 基于浏览器服务器模式的网上教学系统[J]. 安徽机电学院学报, 2001, 16(1): 58-60.
- [5] 陆冬云, 温浩, 许志宏. 以客户为中心的浏览器/服务器(B/S)网络计算模型[J]. 计算机与应用化学, 2001, 18(4): 336-337.
- [6] 梁娜, 禹农, 杨国青. 基于 B/S 计算模型的 Web 技术在电子商务中的应用[J]. 山东科技大学学报, 2003, 22(1): 65-66.
- [7] 杜根远. 基于 B/S 模式的 Web 三层应用开发[J]. 河南城建高等专科学校学报, 2002, 11(1): 37-39.
- [8] 李书杰, 李志刚. B/S 三层体系结构模式[J]. 河北理工大学学报, 2002, 24(S1): 25-28.
- [9] 肖丁, 吴建林, 周春燕, 等. 软件工程模型与方法[M]. 北京: 北京邮电大学出版社, 2008.
- [10] 鲁博, 柴跃廷. 关于统一建模语言-UML[J]. 计算机工程与科学, 2000(4): 13-19.
- [11] 曲美霞. 企业自动审批流程管理系统的实现[J]. 计算机技术与发展, 2006, 16(2): 67-68.
- [12] Lloyd D B. Integrating reporting services into ASP. net[J]. Dr-Dows Journal, 2005(2): 32-33.
- [13] Anley C. Advanced SQL injection in SQL server applications[J]. IEEE Internet Computing, 2004(8): 8-10.
- [14] Sadiq S W. Handling dynamic schema change in process models[C]//Australasian Database Conference. [s. l.]: [s. n.], 2000.

(上接第 218 页)

以统一的 XML 格式进行组织,方便了甲骨文信息的管理和共享,更为以后的甲骨卜辞本体的生成和语义标注打下了坚实的基础^[11]。为下一步甲骨文考释的知识表示与推理打下了语法基础^[12]。

参考文献:

- [1] 张德劭. 甲骨文考释研究[D]. 上海: 华东师范大学, 2002.
- [2] 何婷婷. 语料库研究[D]. 武汉: 华中师范大学, 2003.
- [3] Chan L M, Zeng M L. Metadata interoperability and standardization—a study of methodology part I: achieving interoperability at the schema level[J]. D-Lib Magazine, 2006, 12(6): 121-123.
- [4] 王宁, 郑振峰. 甲骨文字构形系统研究[M]. 上海: 上海教育出版社, 2006.
- [5] 亓祥波, 南琳, 张福顺. 基于元数据和 XML 的信息抽取与集成技术研究[J]. 信息与控制, 2008, 37(1): 52-57.
- [6] 廖乐健, 曹元大, 李新颖. 基于 Ontology 的信息抽取[J]. 计算机工程与应用, 2002, 38(23): 110-113.
- [7] 吴琴霞, 刘永革. 基于 XML/Schema 甲骨文语料库标注的研究[J]. 科学技术与工程, 2009, 17(9): 5186-5188.
- [8] World-Wide Web Consortium Xquery1. 0: An XML Query-Language[EB/OL]. 2005-11-03. <http://www.w3.org/TR/Xquery>, W3C Working Draft.
- [9] 仲华, 催志明. 基于 XML 的信息抽取和多层向量空间技术研究[J]. 计算机技术与发展, 2007, 17(7): 49-52.
- [10] Shamsfard M, Barforoush A A. Learning ontologies from natural language texts[J]. International Journal of Human-Computer Studies, 2004, 60(1): 17-63.
- [11] 梁晓涛, 谢荣传. 基于 OWL 描述本体的语义信息抽取[J]. 计算机技术与发展, 2006, 16(1): 63-65.
- [12] Li Man, Du Xiaoyong, Wang Shan. Learning ontology from relational database[C]//Proceedings of 2005 International Conference on Machine Learning and Cybernetics. [s. l.]: [s. n.], 2005: 3410-3415.