

数据库与信息处理

n-Gram/2L索引结构的存储与时间优化算法

刘凤晨¹, 刘庆文², 胡 玥², 黄 河¹

1.北京航空航天大学 软件学院, 北京 100083

2.北京科技大学 计算机科学系, 北京 100083

收稿日期 2007-6-5 修回日期 2007-8-13 网络版发布日期 2008-1-31 接受日期

摘要 对分词检索算法n-Gram/2L的索引结构作了改进, 在第二级倒排表中加入对文章标识的索引, 提出一种基于Zigzag的分词检索算法*n-Gram/2LZ* (*n-Gram/2L on Zigzag join*)。在对数据量较大的文章进行检索和索引时, 该算法在保留原有算法特性的基础上进一步减少了索引冗余, 降低了索引的存储量, 同时对查询算法的优化降低了查询时的系统开销, 并且减少索引中记录访问次数, 提高了查询效率。

关键词 [算法](#) [索引](#) [n-gram](#) [倒排表](#)

分类号

Space and time optimized algorithm of *n-Gram/2L* index structure

LIU Feng-chen¹, LIU Qing-wen², HU Yue², HUANG He¹

1. College of Software, Beijing University of Aeronautics and Astronautics, Beijing 100083, China

2. Department of Computer Science, Beijing University of Science and Technology, Beijing 100083, China

Abstract

This paper presents an improved algorithm of n-Gram/2L index for text retrieval by adding document identifier index into the secondary level inverted index, and proposes a retrieval algorithm: *n-Gram/2LZ* (*n-Gram/2L on Zigzag join*) based on Zigzag join. This algorithm retains the advantage of former *n-Gram/2L* algorithm and reduces redundancy and storage of the document index, while retrieving and indexing large data. And the optimization of the query algorithm decreases the system overhead when processing query as well as enhances query efficiency by reducing reading the same record repeatedly.

Key words [algorithms](#) [indexing](#) [n-gram](#) [inverted index](#)

DOI:

通讯作者 刘凤晨 aric101@hotmail.com

扩展功能

本文信息

- [Supporting info](#)
- [PDF\(808KB\)](#)
- [\[HTML全文\]\(0KB\)](#)

参考文献

服务与反馈

- [把本文推荐给朋友](#)
- [加入我的书架](#)
- [加入引用管理器](#)
- [复制索引](#)
- [Email Alert](#)
- [文章反馈](#)

浏览反馈信息

相关信息

- [本刊中包含“算法”的相关文章](#)
- [本文作者相关文章](#)
- [刘凤晨](#)
- [刘庆文](#)
- [胡 玥](#)
- [黄 河](#)