

网络、通信与安全

## 基于FFT的网页正文提取算法研究与实现

李 蕾<sup>1,2</sup>, 王劲林<sup>1</sup>, 白 鹤<sup>1,2</sup>, 胡晶晶<sup>1,2</sup>

1.中国科学院 声学研究所 DSP中心,北京 100080

2.中国科学院 研究生院,北京 100039

收稿日期 修回日期 网络版发布日期 2007-10-11 接受日期

**摘要** 主要研究“正文式”网页的有效信息提取算法。该种底层网页真正含有Web页面所表达的主题信息,通常包含一大段的正文信息,正文信息的前后是一些格式信息(例如导航信息、交互信息、JavaScript脚本等)。分析了此种网页的页面结构特征,将问题转化为——给定一个底层网页的HTML源文件,求解最佳的正文区间;从而提出了一种基于快速傅立叶变换的网页正文内容提取算法。采用窗口分段的方法,利用统计学原理和FFT,得出每个可能区间的权值,从而求解出最佳正文区间。实验结果表明,此种方法能比较准确的对“正文式”网页的有效信息进行提取。

**关键词** [中文信息处理](#) [Web页面](#) [信息提取](#) [页面结构](#) [FFT](#) [区域分割](#)

分类号

## Research and implementation of FFT-based extraction algorithm of webpage content main body

LI Lei<sup>1,2</sup>, WANG Jin-lin<sup>1</sup>, BAI He<sup>1,2</sup>, HU Jing-jing<sup>1,2</sup>

1.DSP and Network Research Center, Institute of Acoustics, Chinese Academy of Sciences, Beijing 100080, China

2.Graduate University of Chinese Academy of Sciences, Beijing 100039, China

### Abstract

This paper studies the extraction algorithm of the effective information of “Content-Dominated” Web pages. This kind of Web pages contains the major content information of the Web sites. It includes a long paragraph of content main body, and format information in the beginning and the ending (e.g. navigation information, interaction information, JavaScript and so on). This paper analyzes the structural characteristics of this kind of Web page, and transformed the problem as: given an HTML source file of a “Content-Dominated” Webpage, to find the best range of the content main body. Presents an FFT-based extraction algorithm of webpage content main body. By applying window-segmentation, statistics theory and FFT, this method calculates the weight of every possible range; and thereby selects the best one as solution. The experimental result proves that this algorithm can efficiently extract the effective information of “Content-Dominated” Web pages.

**Key words** [Chinese information processing](#) [Web page](#) [information extraction](#) [Web page structure](#)  
[Fast Fourier Transformation \(FFT\)](#) [page segmentation](#)

DOI:

### 扩展功能

#### 本文信息

- ▶ [Supporting info](#)
- ▶ [PDF\(1704KB\)](#)
- ▶ [\[HTML全文\]\(0KB\)](#)

#### 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“中文信息处理”的相关文章](#)

#### 本文作者相关文章

- [李 蕾](#)
- [王劲林](#)
- [白 鹤](#)
- [胡晶晶](#)
-