## Journal Resources

Editorial Info
Abstracting and Indexing
Release Schedule
Advertising Info

## Author Resources

# Orthographic Errors in Web Pages: Toward Cleaner Web Corpora

Christoph Ringlstetter, Klaus U. Schulz and Stoyan Mihov

👁 **Download Options**    〉

**Abstract**    Authors

Since the Web by far represents the largest public repository of natural language texts, recent experiments, methods, and tools in the area of corpus linguistics often use the Web as a corpus. For applications where high accuracy is crucial, the problem has to be faced that a non-negligible number of orthographic and grammatical errors occur in Web documents. In this article we investigate the distribution of orthographic errors of various types in Web pages. As a by-product, methods are developed

for efficiently detecting erroneous pages and for marking orthographic errors in acceptable Web documents, reducing thus the number of errors in corpora and linguistic knowledge bases automatically retrieved from the Web.

## Forthcoming

## Most Read                                        See More

**Lexicon-Based Methods for Sentiment Analysis (14087 times)**
Maite Taboada et al.
Computational Linguistics Volume: 37, Issue: 2, pp. 267-307

**Computational Linguistics and Deep Learning (10542 times)**
Christopher D. Manning
Computational Linguistics Volume: 41, Issue: 4, pp. 701-707

**Near-Synonymy and Lexical Choice (3675 times)**
Philip Edmonds et al.
Computational Linguistics Volume: 28, Issue: 2, pp. 105-144

(Note that the Most Read numbers are based on the number of full text downloads over the last 12 months.)

## Most Cited                                       See More

**Lexicon-Based Methods for Sentiment Analysis (436 times)**
Maite Taboada et al.
Computational Linguistics Volume: 37, Issue: 2, pp. 267-307

**A Systematic Comparison of Various Statistical Alignment Models (174 times)**
Franz Josef Och et al.
Computational Linguistics Volume: 29, Issue: 1, pp. 19-51

**Opinion Word Expansion and Target Extraction through Double Propagation (147 times)**
Guang Qiu et al.
Computational Linguistics Volume: 37, Issue: 1, pp. 9-27

(Note that the Most Cited numbers are based on Crossref's Cited-by service and reflect citation information for the past 24 months. )
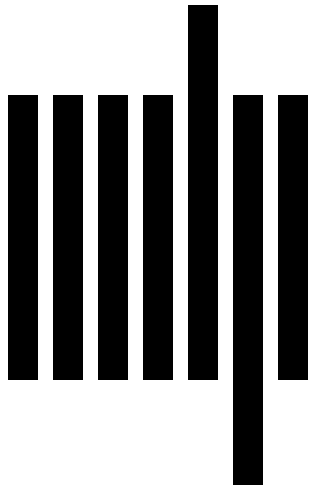
◎ **Download Options** ＞

Favorite ♡
Sign up for Alerts 🔔

Download Citation ↓
RSS TOC 📶

RSS Citation 📶
Submit your article

Support OA at MITP

Journals

Books

US

One Rogers Street Cambridge MA 02142-1209

Terms & Conditions

UK

Suite 2, 1 Duchess Street London, W1W 6AN, UK

Privacy Statement

Connect

Contact Us

Site Help