# QUITA

## Quantitative Index Text Analyzer

Miroslav Kubát

Vladimír Matlach

Radek Čech

# Studies in Quantitative Linguistics

## Editors

Fengxiang Fan   (fanfengxiang@yahoo.com)
Emmerich Kelih  (emmerich.kelih@uni-graz.at)
Reinhard Köhler  (koehler@uni-trier.de)
Ján Mačutek     (jmacutek@yahoo.com)
Eric S. Wheeler  (wheeler@ericwheeler.ca)

1. U. Strauss, F. Fan, G. Altmann, *Problems in quantitative linguistics 1*. 2008, VIII + 134 pp.

2. V. Altmann, G. Altmann, *Anleitung zu quantitativen Textanalysen. Methoden und Anwendungen.* 2008, IV+193 pp.

3. I.-I. Popescu, J. Mačutek, G. Altmann, *Aspects of word frequencies*. 2009, IV +198 pp.

4. R. Köhler, G. Altmann, *Problems in quantitative linguistics 2*. 2009, VII + 142 pp.

5. R. Köhler (ed.), *Issues in Quantitative Linguistics.* 2009, VI + 205  pp.

6. A. Tuzzi, I.-I. Popescu, G.Altmann, *Quantitative aspects of Italian texts*. 2010, IV+161 pp.

7. F. Fan, Y. Deng, *Quantitative linguistic computing with Perl.*  2010, VIII + 205 pp.

8.  I.-I. Popescu et al., *Vectors and codes of text*. 2010, III + 162 pp.

9.  F. Fan, *Data processing and management for quantitative linguistics with Foxpro.* 2010, V + 233 pp.

10. I.-I. Popescu, R. Čech, G. Altmann, *The lambda-structure of texts*. 2011,  II + 181 pp

11. E. Kelih et al. (eds.), *Issues in Quantitative Linguistics Vol. 2*. 2011, IV + 188 pp.

12. R. Čech, G. Altmann, *Problems in quantitative linguistics Vol. 3*. 2011, VI + 168 pp.

13. R. Köhler, G. Altmann (eds.), *Issues in Quantitative Linguistics Vol. 3*. 2013, IV + 403 pp.

14. R. Köhler, G. Altmann, *Problems in Quantitative Linguistics Vol. 4.* 2014. VIII+148 pp.

15. Best, K.-H., Kelih, E. (eds.), *Entlehnungen und Fremdwörter: Quantitative Aspekte.* 2014. VI + 163 pp.

16. I.-I. Popescu, K.-H. Best, G. Altmann, *Unified modeling of length in Language.* 2014, VIII+123 pp.

17. G. Altmann, R. Čech, J. Mačutek, L. Uhlířová (eds.), *Empirical Approaches tot ext and Language Analysis dedicated to Luděk Hřebíček on the occasion of his 80[th] birthday.* 2014, VI + 231 pp.

18. M. Kubát, V. Matlach, R. Čech, *QUITA. Quantitative Index Text Analyzer.* 2014, VII+106.

# Contents

# 1    Introduction

Quantitative experimental methods have been increasingly used in the humanities in recent years. We can hardly imagine the disciplines of social science, such as psychology, sociology or economics, without a quantitative approach. On the other hand, the majority of linguists, historians and especially literary critics are still refusing to use quantitative methods. One of the reasons is the fact that those researchers consider quantitative methods, and especially statistical methods, too difficult to apply to their field. QUITA (Quantitative Indicator Text Analyzer) is a tool which aims to help all people who try to analyse texts by quantitative methods.

QUITA is a program to enable researchers from various disciplines (linguistics, criticism, history, sociology, psychology, politics, biology, etc.) to analyse texts using quantitative methods. There are many indicators which measure various characteristics of text. Although the authors of QUITA focused mainly on indicators connected to the frequency structure of a text, there are also functions for several other characteristics. Since QUITA is designed especially for researchers outside quantitative linguistics, it includes functions for the most basic and common indicators.

Given that the main purpose of QUITA is to provide a user-friendly tool of quantitative text analysis for researchers without a deeper knowledge of quantitative linguistics, statistics or programming, QUITA also provides simple statistical comparisons and the ability to create charts. There is no need to use any additional software such as spreadsheet applications or special statistical programs. QUITA is therefore the program that combines all the essential parts of any quantitative research effort: obtaining results, statistical testing and graphical visualization.

The QUITA manual is written with step-by-step instructions. All tools are concisely described and accompanied by screen-shots. Every indicator is briefly presented (complete with references), and mathematically defined. There are also examples of computation and statistical comparison.

Although the manual provides users with all the essential information about QUITA, it was not possible to cover most topics in deeper detail. For this purpose, we highly recommend the book *Word frequency studies* (Popescu et al. 2009) which is a comprehensive overview about quantitative analysis using indicators based on the frequency structure of a text. The book *Aspects of Word Frequencies* (Popescu et al. 2009) is also well worth reading. Detailed examples of computing most indicators used in QUITA can be found in *Metody kvantitativní analýzy (nejen) básnických textů* (Čech et al. 2014).

Since we aim to help as many researchers as possible, QUITA is distributeed as freeware. Thus anyone can use QUITA without any restrictions. The latest version of the software is available on the website https://code.google.com/p/oltk/.  In published work, acknowledgement of QUITA would be appropriate and appreciated.

# 2    System requirements

Supported Systems: Windows XP, Windows Vista, Windows 7, Windows 8.

System requirements: NET Framework 3.5

Optional system requirements: Python, Perl

NOTE:

Users are automatically notified that it is necessary to install the system requirements with the download links.