**Slovenská akadémia vied**
Jazykovedný ústav Ľudovíta Štúra

# Natural Language Processing, Corpus Linguistics, E-learning

Seventh International Conference
Bratislava, Slovakia, 13–15 November 2013
Proceedings

Editors
Katarína Gajdošová
Adriána Žáková

**The articles have been reviewed by members of the Program Committee.**

The articles can be used under the
Creative Commons Attribution-ShareAlike 3.0 Unported License

# Table of Contents

# Foreword

***Slovko 2013 – Natural Language Processing, Corpus Linguistics, E-learning*** will be again held in Bratislava. The organizers – Slovak National Corpus Department of Ľ. Štúr Institute of Linguistics, Slovak Academy of Sciences are honoured to host participants from eight countries: Bulgaria, Czech Republic, Germany, Greece, Hungary, Poland, Slovakia and Slovenia.

Over three days participants will be able to benefit from 29 presentations, including 3 plenary talks. Unfortunately, one third of submitted papers on given topics has not been recommended by the Programme Committee members. We thank to all reviewers for their constructive suggestions and their help to make the conference even more successful.

The 7[th] edition of the biennial conference increased the presence of the linguistically-oriented (corpus-based and corpus-driven) studies. The more technically oriented papers provide information on effectiveness of the approaches applied, experimenting and innovative methods. Latest trends and tendencies in enhancing the corpus data can be found also in the papers written by Slovak authors.

We wish all participants of the conference Slovko 2013 profitable time and positive inspiration for further cooperation in the field of natural language processing, corpus linguistics and similar research.

*Mária Šimková*
*Translated by Adriána Žáková*