

14th Workshop on Building and Using Comparable Corpora (BUCC 2021)

In Conjunction with the 13th International Conference on Recent
Advances in Natural Language Processing (RANLP 2021)

Online
6 September 2021

Editors:

**Richard Rapp
Serge Sharoff
Pierre Zweigenbaum**

ISBN: 978-1-7138-4098-5

Printed from e-media with permission by:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571



Some format issues inherent in the e-media version may also appear in this print version.

Copyright© (2021) by the Association for Computational Linguistics
All rights reserved.

Printed with permission by Curran Associates, Inc. (2022)

For permission requests, please contact the Association for Computational Linguistics
at the address below.

Association for Computational Linguistics
209 N. Eighth Street
Stroudsburg, Pennsylvania 18360

Phone: 1-570-476-8006
Fax: 1-570-476-0860

acl@aclweb.org

Additional copies of this publication are available from:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571 USA
Phone: 845-758-0400
Fax: 845-758-2633
Email: curran@proceedings.com
Web: www.proceedings.com

Table of Contents

<i>Machine Translation in Low Resource Setting</i> Pushpak Bhattacharyya	1
<i>Mining Bilingual Word Pairs from Comparable Corpus using Apache Spark Framework</i> Sanjanasri JP, Vijay Krishna Menon, Soman KP and Krzysztof Wolk	2
<i>Effective Bitext Extraction From Comparable Corpora Using a Combination of Three Different Approaches</i> Steintor Steingrímsson, Pintu Lohar, Hrafn Loftsson and Andy Way	8
<i>Syntax-aware Transformers for Neural Machine Translation: The Case of Text to Sign Gloss Translation</i> Santiago Egea Gomez, Euan McGill and Horacio Saggion	18
<i>Employing Wikipedia as a resource for Named Entity Recognition in Morphologically complex under-resourced languages</i> Aravind Krishnan, Stefan Ziehe, Franziska Pannach and Caroline Sporleder	28
<i>Semi-Automated Labeling of Requirement Datasets for Relation Extraction</i> Jeremias Bohn, Jannik Fischbach, Martin Schmitt, Hinrich Schuetze and Andreas Vogelsang ...	40
<i>Majority Voting with Bidirectional Pre-translation For Bitext Retrieval</i> Alexander Jones and Derry Tanti Wijaya	46
<i>EM Corpus: a comparable corpus for a less-resourced language pair Manipuri-English</i> Rudali Huidrom, Yves Lepage and Khogendra Khomdram	60
<i>On Pronunciations in Wiktionary: Extraction and Experiments on Multilingual Syllabification and Stress Prediction</i> Winston Wu and David Yarowsky	68
<i>A Dutch Dataset for Cross-lingual Multilabel Toxicity Detection</i> Ben Burtenshaw and Mike Kestemont	75