

Workshop on Arabic Natural Language Processing

(ANLP 2014))

**Held at the Conference on Empirical Methods in Natural
Language Processing (EMNLP 2014)**

**Doha, Qatar
25 October 2014**

ISBN: 978-1-63439-479-6

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© (2014) by the Association for Computational Linguistics
All rights reserved.

Printed by Curran Associates, Inc. (2014)

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: 845-758-2634
Email: curran@proceedings.com
Web: www.proceedings.com

Introduction

Welcome to Arabic Natural Language Processing Workshop at EMNLP 2014 in Doha, Qatar.

There has been a lot of progress in the last 15 years in the area of Arabic Natural Language Processing (NLP). In particular, the TIDES, GALE, and BOLT programs provided a significant boost to Arabic NLP, both in generating new language and speech resources on a large scale, and in advancing the state-of-the-art in morphological processing, parsing, named entity recognition, information retrieval, speech recognition, and machine translation. The substantial investment done through these projects reflect the fact that the Middle East continues to grow in its political and economic importance. We also observe that countries in the Middle East invest substantially into higher education and into building an ecosystem, which fosters new research initiatives. This creates the hope that our own research field, NLP, and especially Arabic NLP will continue to grow, will continue to attract students both in the region and in top international universities, and that new job opportunities will open up not only in the well established language service providers, but also through start-ups offering innovative solutions.

A number of Arabic NLP (or Arabic NLP-related) workshops and conferences have taken place, both in the Arab World and in association with international conferences. The Arabic NLP workshop at EMNLP 2014 follows in the footsteps of these previous efforts to provide a forum for researchers to share and discuss their ongoing work. The Arabic NLP workshop also includes a shared task on Automatic Arabic Error Correction, which was designed in the tradition of high profile NLP shared tasks such as CONLL's grammar/error detection and correction shared tasks in 2011-2013, and numerous machine translation campaigns by NIST/WMT/MEDAR, among others. The challenge chosen for the shared task is highly relevant, not only to spelling correction while composing a text, but also to developing techniques for automatically correcting errors in the far-from-perfect outputs of NLP technologies such as speech recognition or machine translation.

We are happy to have received 40 submissions. Unfortunately, not all the papers could be included in the workshop due to time limitations. The acceptance rate is 50%. The papers cover a wide range of topics: building language resources for standard and dialectal Arabic, language identification, sentiment analysis, named entity disambiguation, and machine translation for dialectal Arabic, etc. Twelve papers were selected for oral presentation and were organized under the general topics Corpora (four papers), Text Mining (four papers), Translation & Transliteration (three papers) and one paper describing the shared task. The remaining eight papers were selected to be presented in a poster session. There is no difference in quality between the oral and poster presentations.

The shared task was a success. We received 18 systems submissions from nine teams in six countries, representing a diverse set of approaches. Nine shared task system description (short) papers are included in the proceedings to document the shared task systems, but were not reviewed with the rest of the papers of the main workshop. These papers will be presented as posters.

The quantity and quality of the contributions to the main workshop, as well as the shared task, are strong indicators that there is a continued need for this kind of dedicated Arabic NLP workshop. We would like to acknowledge all the hard work of the submitting authors and thank the reviewers for their diligent work and for the valuable feedback they provided. We are also thankful to the work of the shared task committee, website committee and the publication co-chairs.

It has been an honor to server as program co-chairs. We hope that the reader of these proceedings will find them stimulating and beneficial.

Nizar Habash and Stephan Vogel, Arabic NLP Workshop, EMNLP 2014.

Organizers:

Program Co-Chairs

Nizar Habash, New York University Abu Dhabi
Stephan Vogel, Qatar Computing Research Institute

Publication Co-chairs

Nadi Tomeh, Université Paris 13, Sorbonne Paris Cité
Houda Bouamor, Carnegie Mellon University Qatar

Website Committee

Noura Farra, Columbia University
Kareem Darwish, Qatar Computing Research Institute

Shared Task Committee

Behrang Mohit (co-chair), Carnegie Mellon University Qatar
Alla Rozovskaya (co-chair), Columbia University
Wajdi Zaghouani, Carnegie Mellon University Qatar
Ossama Obeid, Carnegie Mellon University Qatar
Nizar Habash (advisor), New York University Abu Dhabi

Program Committee:

Abdelmajid Ben-Hamadou, University of Sfax, Tunisia
Abdelhadi Soudi, Ecole Nationale de l'Industrie Minérale, Morocco
Abdelsalam Nwesri, University of Tripoli, Libya
Achraf Chalabi, Microsoft Research, Egypt
Ahmed Ali, Qatar Computing Research Institute, Qatar
Ahmed Rafea, The American University in Cairo, Egypt
Alexis Nasr, University of Marseille, France
Ali Farghaly, Monterey Peninsula College, USA
Almoataz B. Al-Said, Cairo University, Egypt
Alon Lavie, Carnegie Mellon University, USA
Aly Fahmy, Cairo University, Egypt
Azadeh Shakery, University of Tehran, Iran
Azzeddine Mazroui, University Mohamed I, Morocco
Bassam Haddad, University of Petra, Jordan
Bayan Abu Shawar, Arab Open University, Jordan
Behrang Mohit, Carnegie Mellon University Qatar, Qatar
Eric Atwell, University of Leeds, UK
Farhad Oroumchian, University of Wollongong, Australia
Ghassan Mourad, Université Libanaise, Lebanon
Hassan Sawaf, eBay Inc., USA
Hazem Hajj, American University of Beirut, Lebanon
Hend Alkhalifa, King Saud University, Saudi Arabia
Houda Bouamor, Carnegie Mellon University Qatar, Qatar
Imed Zitouni, Microsoft Research, USA

Joseph Dichy, Université Lyon 2, France
Karim Bouzoubaa , Mohammad V University, Morocco
Karine Megerdoomian, The MITRE Corporation, USA
Katrin Kirchhoff, University of Washington, USA
Kemal Oflazer, Carnegie Mellon University Qatar, Qatar
Khaled Shaalan, The British University in Dubai, UAE
Khaled Shaban, Qatar University, Qatar
Khalil Sima'an, Universiteit van Amsterdam, Netherlands
Lamia Hadrich Belguith, University of Sfax, Tunisia
Michael Rosner, University of Malta, Malta
Mohamed Elmahdy, Qatar University, Qatar
Mohsen Rashwan, Cairo University, Egypt
Mona Diab, George Washington University, USA
Mustafa Jarrar, Bir Zeit University, Palestine
Nada Ghneim, Higher Institute for Applied Sciences and Technology, Syria
Nadi Tomeh, Université Paris 13, Sorbonne Paris Cité, France
Ossama Emam, IBM, USA
Otakar Smrž, Charles University , Czech Republic
Owen Rambow, Columbia University, USA
Preslav Nakov, Qatar Computing Research Institute, Qatar
Ramzi Abbes, TECHLIMED, France
Salwa Hamada, Cairo University, Egypt
Shahram Khadivi, Tehran Polytechnic, Iran
Sherri Condon , The MITRE Corporation, USA
Taha Zerrouki, University of Bouira, Algeria
Violetta Cavalli-Sforza, Al Akhawayn University, Morocco

Table of Contents

<i>Using Twitter to Collect a Multi-Dialectal Corpus of Arabic</i> Hamdy Mubarak and Kareem Darwish	1
<i>The International Corpus of Arabic: Compilation, Analysis and Evaluation</i> Sameh Alansary and Magdy Nagi	8
<i>Building a Corpus for Palestinian Arabic: a Preliminary Study</i> Mustafa Jarrar, Nizar Habash, Diyam Akra and Nasser Zalmout	18
<i>Annotating corpus data for a quantitative, constructional analysis of motion verbs in Modern Standard Arabic</i> Dana Abdulrahim	28
<i>The First QALB Shared Task on Automatic Text Correction for Arabic</i> Behrang Mohit, Alla Rozovskaya, Nizar Habash, Wajdi Zaghoulani and Ossama Obeid	39
<i>A Framework for the Classification and Annotation of Multiword Expressions in Dialectal Arabic</i> Abdelati Hawwari, Mohammed Attia and Mona Diab	48
<i>Al-Bayan: An Arabic Question Answering System for the Holy Quran</i> Heba Abdelnasser, Maha Ragab, Reham Mohamed, Alaa Mohamed, Bassant Farouk, Nagwa El-Makky and Marwan Torki	57
<i>Automatic Arabic diacritics restoration based on deep nets</i> Ahmad Al Sallab, Mohsen Rashwan, Hazem M. Raafat and Ahmed Rafea	65
<i>Combining strategies for tagging and parsing Arabic</i> Maytham Alabbas and Allan Ramsay	73
<i>Named Entity Recognition System for Dialectal Arabic</i> Ayah Zirikly and Mona Diab	78
<i>Semantic Query Expansion for Arabic Information Retrieval</i> Ashraf Mahgoub, Mohsen Rashwan, Hazem Raafat, Mohamed Zahran and Magda Fayek	87
<i>Transliteration of Arabizi into Arabic Orthography: Developing a Parallel Annotated Arabizi-Arabic Script SMS/Chat Corpus</i> Ann Bies, Zhiyi Song, Mohamed Maamouri, Stephen Grimes, Haejoong Lee, Jonathan Wright, Stephanie Strassel, Nizar Habash, Ramy Eskander and Owen Rambow	93
<i>Tunisian dialect Wordnet creation and enrichment using web resources and other Wordnets</i> Rihab Bouchlaghem, Aymen Elkhelifi and Rim Faiz	104
<i>A Pipeline Approach to Supervised Error Correction for the QALB-2014 Shared Task</i> Nadi Tomeh, Nizar Habash, Ramy Eskander and Joseph Le Roux	114
<i>Arabic Spelling Correction using Supervised Learning</i> Youssef Hassan, Mohamed Aly and Amir Atiya	121
<i>Autocorrection of arabic common errors for large text corpus</i> Taha Zerrouki, Khaled Alhawiti and Amar Balla	127

<i>Automatic Correction of Arabic Text: a Cascaded Approach</i>	
Hamdy Mubarak and Kareem Darwish	132
<i>CMUQ@QALB-2014: An SMT-based System for Automatic Arabic Error Correction</i>	
Serena Jeblee, Houda Bouamor, Wajdi Zaghouni and Kemal Oflazer	137
<i>Fast and Robust Arabic Error Correction System</i>	
Michael Nawar and Moheb Ragheb	143
<i>GWU-HASP: Hybrid Arabic Spelling and Punctuation Corrector</i>	
Mohammed Attia, Mohamed Al-Badrashiny and Mona Diab	148
<i>TECHLIMED system description for the Shared Task on Automatic Arabic Error Correction</i>	
Djamel MOSTEFA, Omar ASBAYOU and Ramzi ABBES	155
<i>The Columbia System in the QALB-2014 Shared Task on Arabic Error Correction</i>	
Alla Rozovskaya, Nizar Habash, Ramy Eskander, Noura Farra and Wael Salloum	160
<i>A Large Scale Arabic Sentiment Lexicon for Arabic Opinion Mining</i>	
Gilbert Badaro, Ramy Baly, Hazem Hajj, Nizar Habash and Wassim El-Hajj	165
<i>Evaluating Distant Supervision for Subjectivity and Sentiment Analysis on Arabic Twitter Feeds</i>	
Eshrag Refaee and Verena Rieser	174
<i>Arabic Native Language Identification</i>	
Shervin Malmasi and Mark Dras	180
<i>AIDArabic A Named-Entity Disambiguation Framework for Arabic Text</i>	
Mohamed Amir Yosef, Marc Spaniol and Gerhard Weikum	187
<i>Domain and Dialect Adaptation for Machine Translation into Egyptian Arabic</i>	
Serena Jeblee, Weston Feely, Houda Bouamor, Alon Lavie, Nizar Habash and Kemal Oflazer .	196
<i>Unsupervised Word Segmentation Improves Dialectal Arabic to English Machine Translation</i>	
Kamla Al-Mannai, Hassan Sajjad, Alaa Khader, Fahad Al Obaidli, Preslav Nakov and Stephan Vogel	207
<i>Arabizi Detection and Conversion to Arabic</i>	
Kareem Darwish	217

Author Index