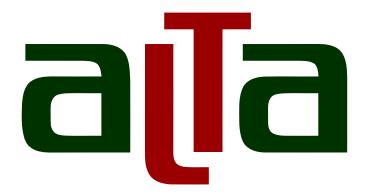
Australasian Language Technology Association Workshop 2017

Proceedings of the Workshop



Editors: Jojo Sze-Meng Wong Gholamreza Haffari

6–8 December 2017 Queensland University of Technology Brisbane, Australia

Australasian Language Technology Association Workshop 2017 (ALTA 2017)

http://alta2017.alta.asn.au

Online Proceedings:

http://alta2017.alta.asn.au/proceedings

Gold Sponsors:



Silver Sponsors:



Bronze Sponsors:



Volume 15, 2017 ISSN: 1834-7037

ALTA 2017 Workshop Committees

Workshop Chairs

- Stephen Wan (Data61)
- Jojo Sze-Meng Wong (Monash University)

Workshop Programme Chairs

- Jojo Sze-Meng Wong (Monash University)
- Gholamreza Haffari (Monash University)

Programme Committee

- Oliver Adams, University of Melbourne
- Timothy Baldwin, University of Melbourne
- Benjamin Borschinger, Google
- Julian Brooke, University of Melbourne
- Lawrence Cavedon, RMIT University
- Trevor Cohn, The University of Melbourne
- Nathalie Colineau, DST Australia
- Mark Dras, Macquarie University
- Dominique Estival, Western Sydney University
- Gabriela Ferraro, CSIRO Data61
- Hamed Hassanzadeh, The Australian e-Health Research Centre
- Nitin Indurkhya, University of New South Wales
- Sarvnaz Karimi, CSIRO Data61
- Mac Kim, CSIRO Data61
- Yitong Li, The University of Melbourne
- Teresa Lynn, Dublin City University
- Andrew MacKinlay, IBM Research
- Diego Mollá-Alliod, Macquarie University
- Anthony Nguyen, The Australian e-Health Research Centre
- Bahadorreza Ofoghi, The University of Melbourne
- Sylvester Orimaye, East Tennessee State University
- Cécile Paris, CSIRO Data61
- Matthias Petri, The University of Melbourne
- Lizhen Qu, CSIRO Data61
- Will Radford, Red Marker
- Abeed Sarker, University of Pennsylvania
- Rolf Schwitter, Macquarie University
- Ehsan Shareghi, Monash University
- Laurianne Sitbon, Queensland University of Technology
- Hanna Suominen, The Australian National University
- Karin Verspoor, The University of Melbourne
- Wei Wang, University of New South Wales
- Ingrid Zukerman, Monash University

Preface

This volume contains the papers accepted for presentation at the Australasian Language Technology Association Workshop (ALTA) 2017, held at Queensland University of Technology in Brisbane, Australia on 6–8 December 2017.

The goals of the workshop are to:

- bring together the Language Technology (LT) community in the Australasian region and encourage interactions and collaboration;
- foster interaction between academic and industrial researchers, to encourage dissemination of research results:
- provide a forum for students and young researchers to present their research;
- facilitate the discussion of new and ongoing research and projects;
- increase visibility of LT research in Australasia and overseas and encourage interactions with the wider international LT community.

This year's ALTA Workshop presents 13 peer-reviewed papers, including 10 long papers and 3 short papers. We received a total of 23 submissions for long and short papers. Each paper was reviewed by three members of the program committee, using a double-blind protocol. Great care was taken to avoid all conflicts of interest.

ALTA 2017 includes a presentations track, following the workshops since 2015 when it was first introduced. This aims to encourage broader participation and facilitate local socialisation of international results, including work in progress and work submitted or published elsewhere. Presentations were lightly reviewed by the ALTA chairs to gauge overall quality of work and whether it would be of interest to the ALTA community. Offering both archival and presentation tracks allows us to grow the standard of work at ALTA, to better showcase the excellent research being done locally.

ALTA 2017 continues the tradition of including a shared task, this year on correcting OCR errors. Participation is summarised in an overview paper by organisers Diego Mollá-Alliod and Steve Cassidy. Participants were invited to submit a system description paper, which are included in this volume without review.

We would like to thank, in no particular order: all of the authors who submitted papers; the programme committee for the time and effort the put into maintaining the high standards of our reviewing process; the co-chair Stephen Wan for coordinating the logistics that go into running the workshop, from arranging the space, catering, budgets, sponsorship and more; the shared task organisers Diego Mollá and Steve Cassidy; our keynote speakers Lewis Mitchell and Robert Dale for agreeing to share their perspectives on the state of the field; and the tutorial presenter Ben Hachey for his efforts towards the three parts of the tutorial. We would like to acknowledge the constant support and advice of the ALTA Executive Committee.

Finally, we gratefully recognise our sponsors: Capital Markets CRC, Sintelix, Google, CSIRO/Data61 and Queensland University of Technology. Importantly, their generous support enabled us to offer travel subsidies to all students presenting at ALTA, and helped to subsidise conference catering costs and student paper awards.

Jojo Sze-Meng Wong Gholamreza Haffari

ALTA Programme Chairs

ALTA 2017 Programme

Wednesday, 6 December 2017

*Tutorial Session 1 (Monash Caulfield, B214)		
13:00–17:00	Tutorial: Ben Hachey Active Learning and Beyond!	
13:00-14:15	Part 1: From Zero to Hero	
14:15-14:30	Break	
14:30–15:45	Part 2: Live Shared Task	
15:45–16:00	Break	
16:00-17:00	Part 3: Wild Blue Yonder	

Thursday, 7 December 2017

Opening & Key	note (Room P421)
9:00-9:15	Opening
9:15–10:15	Keynote 1 (from ADCS): Dan Russell What do you really need to know? Learning and knowing in the age of the Internet
10:15-10:45	Morning tea
Session 1: Macl	nine Learning and Applications (Room P521)
10:45–11:05	Paper: Leonardo Dos Santos Pinheiro and Mark Dras Stock Market Prediction with Deep Learning: A Character-based Neural Language Model for Event-based Trading
11:05–11:25	Paper: Fei Liu, Trevor Cohn and Timothy Baldwin Improving End-to-End Memory Networks with Unified Weight Tying
11:25–11:45	Paper Shivashankar Subramanian, Trevor Cohn, Timothy Baldwin and Julian Brooke Joint Sentence-Document Model for Manifesto Text Analysis
11:45–12:05	Paper: Ming Liu, Gholamreza Haffari, Wray Buntine and Michelle Ananda-Rajah Leveraging Linguistic Resources for Improving Neural Text Classification
12:05–12:15	Paper: Hamideh Hajiabadi, Diego Molla-Aliod and Reza Monsefi On Extending Neural Networks with Loss Ensembles for Text Classification
12:15–13:15	Lunch
Session 2 & Key	ynote (Room P421)
13:15–14:15	Keynote 2: Lewis Mitchell What do you really need to know? Learning and knowing in the age of the Internet
14:15–14:30	Paper: Shiwei Zhang, Xiuzhen Zhang and Jeffrey Chan (ADCS short paper) A Word-Character Convolutional Neural Network for Language-Agnostic Twitter Sentiment Analysis
14:30–14:45	Paper: Lance De Vine, Shlomo Geva and Peter Bruza (ADCS short paper) Efficient Analogy Completion with Word Embedding Clusters
14:45–15:05	Paper: Aili Shen, Jianzhong Qi and Timothy Baldwin A Hybrid Model for Quality Assessment of Wikipedia Articles
15:05–15:15	Paper: Diego Molla-Aliod Towards the Use of Deep Reinforcement Learning with Global Policy For Query-based Extractive Summarisation
Session 3: Trans	slation and Low Resource Languages (Room P521)
15:45–16:05	Presentation: Inigo Jauregi Unanue, Lierni Garmendia Arratibel, Ehsan Zare Borzeshi and Massimo Piccardi English-Basque Statistical and Neural Machine Translation
16:05–16:25	Presentation: Euna Kim Study on the Role of Machine Translation in Social Network Services in terms of User Centered Orientation: A Case Study of Instagram
16:25–16:45	Presentation: Yunsil Jo Study on Documentary Translation for Dubbing
16:45–17:05	Paper: Oliver Adams, Trevor Cohn, Graham Neubig and Alexis Michaud Phonemic Transcription of Low-Resource Tonal Languages
17:05–17:25	Presentation: Hanieh Poostchi, Ehsan Zare Borzeshi and Massimo Piccardi BiLSTM-CRF for Persian Named-Entity Recognition
17:25	End of Day 1
19:00	Dinner

Friday, 8 December 2017

9:15-10:15	Keynote 3 (fromADCS): Victor Kovalev, Redbubble (Room P421)
	Solving hard problems at massive scale – applied data science research approach at Redbubble
Session 4: Com	putational Linguistics and Information Extraction (Room P521)
10:45–10:55	Paper: Dat Quoc Nguyen, Thanh Vu, Dai Quoc Nguyen, Mark Dras and Mark Johnson From Word Segmentation to POS Tagging for Vietnamese
10:55–11:15	Paper: Shunichi Ishihara A Comparative Study of Two Statistical Modelling Approaches for Estimating Multivariate Like lihood Ratios in Forensic Voice Comparison
11:15–11:35	Paper: Katharine Cheng, Timothy Baldwin and Karin Verspoor Automatic Negation and Speculation Detection in Veterinary Clinical Text
11:35–11:55	Paper: Xiang Dai, Sarvnaz Karimi and Cecile Paris Medication and Adverse Event Extraction from Noisy Text
11:55–12:15	Paper: Maria Myunghee Kim Incremental Knowledge Acquisition Approach for Information Extraction on both Semi-structured and Unstructured Text from the Open Domain Web
12:15	Lunch
13:15–14:15	Keynote 4: Robert Dale, Language Technology Group Pty Ltd (Room P421)
	Commercialised NLP: The state of the art
14:15–15:00	Poster Session (ALTA & ADCS)
15:00-15:30	Afternoon Tea
Shared Task Ses	ssion (Room P521)
15:30–15:40	Diego Molla-Aliod and Steve Cassidy Overview of the 2017 ALTA Shared Task: Correcting OCR Errors
15:40-15:50	Gitansh Khirbat OCR Post-Processing Text Correction using Simulated Annealing (OPTeCA)
15:50–16:00	Yufei Wang SuperOCR for ALTA 2017 Shared Task
Final Session (R	Room P521)
16:00-16:15	Best Paper and Poster Presentation Awards
16:15–16:45	Business Meeting
16:45-17:00	Closing
17:00	End of Day 2

Contents

Invited talks	1
Tutorials	3
Long papers	5
Stock Market Prediction with Deep Learning: A Character-based Neural Language Model for Event-based Trading Leonardo Dos Santos Pinheiro and Mark Dras	6
Improving End-to-End Memory Networks with Unified Weight Tying Fei Liu, Trevor Cohn and Timothy Baldwin	16
Joint Sentence-Document Model for Manifesto Text Analysis Shivashankar Subramanian, Trevor Cohn, Timothy Baldwin and Julian Brooke	25
Leveraging linguistic resources for improving neural text classification Ming Liu, Gholamreza Haffari, Wray Buntine and Michelle Ananda-Rajah	34
A Hybrid Model for Quality Assessment of Wikipedia Articles Aili Shen, Jianzhong Qi and Timothy Baldwin	43
Phonemic Transcription of Low-Resource Tonal Languages Oliver Adams, Trevor Cohn, Graham Neubig and Alexis Michaud	53
A Comparative Study of Two Statistical Modelling Approaches for Estimating Multivariate Likelihood Ratios in Forensic Voice Comparison Shunichi Ishihara	61
Automatic Negation and Speculation Detection in Veterinary Clinical Text Katharine Cheng, Timothy Baldwin and Karin Verspoor	70
Medication and Adverse Event Extraction from Noisy Text Xiang Dai, Sarvnaz Karimi and Cecile Paris	79
Incremental Knowledge Acquisition Approach for Information Extraction on both Semi-structured and Unstructured Text from the Open Domain Web Maria Myunghee Kim	88
Short papers	97
On Extending Neural Networks with Loss Ensembles for Text Classification Hamideh Hajiabadi, Diego Mollá-Alliod and Reza Monsefi	98

Towards the Use of Deep Reinforcement Learning with Global Policy For Query-based Extr Summarisation	
Diego Mollá-Alliod	103
From Word Segmentation to POS Tagging for Vietnamese Dat Quoc Nguyen, Thanh Vu, Dai Quoc Nguyen, Mark Dras and Mark Johnson	108
ALTA Shared Task papers	114
Overview of the 2017 ALTA Shared Task: Correcting OCR Errors Diego Mollá-Alliod and Steve Cassidy	115
OCR Post-Processing Text Correction using Simulated Annealing (OPTeCA) Gitansh Khirbat	119
SuperOCR for ALTA 2017 Shared Task Yufei Wang	124

Invited talks

Characterising Information and Happiness in Online Social Activity

Lewis Mitchell (Lecturer in Applied Mathematics, University of Adelaide)

Abstract. Understanding the nature of influence and information propagation in social networks is of clear societal importance, as they form the basis for phenomena like "echo chambers" and "emotional contagion". However, these concepts remain surprisingly ill-defined. In studies of large online social networks, proxies for influence and information are routinely employed, leading to confusion as to whether the phenomena they underlie actually exist. In this talk I will demonstrate how online social media streams can be used as proxies for population-level health characteristics such as obesity and happiness, and introduce information-theoretic tools for constructing social networks from underlying information flows between individuals. I will present results relating individual predictability to popularity and contact volume, and introduce a paradigmatic mathematical model of information flow over social networks.

Bio. Lewis's research focusses on large-scale methods for extracting useful information from online social networks, and on mathematical techniques for inference and prediction using these data. He works on building tools for real-time estimation of social phenomena such as happiness from written text, and prediction of population-level events like disease outbreaks, elections, and civil unrest.

Commercialised NLP: The State of the Art

Robert Dale (Principal Consultant, Language Technology Group Pty Ltd)

Abstract. The last few years have seen a tremendous surge in commercial interest in Artificial Intelligence, and with it, a widespread recognition that technologies based on Natural Language Processing can support valuable commercial applications. In this talk, I'll aim to give a comprehensive picture of the commercial NLP landscape, focusing on what I see as the key categories of activity: [1] virtual assistants, including chatbots; [2] text analytics and text mining technologies; [3] machine translation; [4] natural language generation; and [5] text correction technologies. In each case my goal is to sketch the history of work in the area, to identify the major players, and to give a realistic appraisal of the state of the art.

Bio. Robert Dale runs the Language Technology Group, an independent consultancy providing unbiased advice to corporations and businesses on the selection and deployment of NLP technologies. Until recently, he was Chief Technology Officer of Arria NLG, where he led the development of a cloud-based natural language generation tool; prior to joining Arria in 2012, he held a chair in the Department of Computing at Macquarie University in Sydney, where he was Director of that university's Centre for Language Technology. After receiving his PhD from the University of Edinburgh in 1989, he taught there for several years before moving to Sydney in 1994. He played a foundational role in building up the NLP community in Australia, and was editor in chief of the Computational Linguistics journal from 2003 to 2012. He writes a semi-regular column titled 'Industry Watch' for the Journal of Natural Language Engineering.

Tutorials

Active Learning ... and Beyond!

Ben Hachey (The University of Sydney)

This half-day session will take participants through situations they might face applying Natural Language Processing to real-world problems. We'll choose a canonical task (text classification) and focus on the main issue that faces practitioners in green fields projects — where does the data come from? Our aim is to equip participants with the theoretical background and practical skills to quickly build high-quality text classification models.