



## NRC Publications Archive Archives des publications du CNRC

### **Editorial: Computational linguistics and literature** Szpakowicz, Stan; Feldman, Anna; Kazantseva, Anna

This publication could be one of several versions: author's original, accepted manuscript or the publisher's version. / La version de cette publication peut être l'une des suivantes : la version prépublication de l'auteur, la version acceptée du manuscrit ou la version de l'éditeur.

For the publisher's version, please access the DOI link below. / Pour consulter la version de l'éditeur, utilisez le lien DOI ci-dessous.

#### **Publisher's version / Version de l'éditeur:**

<https://doi.org/10.3389/fdigh.2018.00024>

*Frontiers in Digital Humanities*, 5, 2018-09-27

#### **NRC Publications Record / Notice d'Archives des publications de CNRC:**

<https://nrc-publications.canada.ca/eng/view/object/?id=39bf37f2-e732-4e3a-b4de-c614e522e7d9>

<https://publications-cnrc.canada.ca/fra/voir/objet/?id=39bf37f2-e732-4e3a-b4de-c614e522e7d9>

Access and use of this website and the material on it are subject to the Terms and Conditions set forth at

<https://nrc-publications.canada.ca/eng/copyright>

READ THESE TERMS AND CONDITIONS CAREFULLY BEFORE USING THIS WEBSITE.

L'accès à ce site Web et l'utilisation de son contenu sont assujettis aux conditions présentées dans le site

<https://publications-cnrc.canada.ca/fra/droits>

LISEZ CES CONDITIONS ATTENTIVEMENT AVANT D'UTILISER CE SITE WEB.

**Questions?** Contact the NRC Publications Archive team at

PublicationsArchive-ArchivesPublications@nrc-cnrc.gc.ca. If you wish to email the authors directly, please see the first page of the publication for their contact information.

**Vous avez des questions?** Nous pouvons vous aider. Pour communiquer directement avec un auteur, consultez la première page de la revue dans laquelle son article a été publié afin de trouver ses coordonnées. Si vous n'arrivez pas à les repérer, communiquez avec nous à PublicationsArchive-ArchivesPublications@nrc-cnrc.gc.ca.





# Editorial: Computational Linguistics and Literature

Stan Szpakowicz<sup>1\*</sup>, Anna Feldman<sup>2</sup> and Anna Kazantseva<sup>3</sup>

<sup>1</sup> School of Electrical Engineering and Computer Science, University of Ottawa, Ottawa, ON, Canada, <sup>2</sup> Department of Linguistics and Department of Computer Science, Montclair State University, Montclair, NJ, United States, <sup>3</sup> National Research Council Canada (NRC-CNRC), Ottawa, ON, Canada

**Keywords:** computational linguistics, literary data, machine learning, machine translation, post-editing, rhetorical figures, themes and motifs in poetry, Middle High German epic poetry

## Editorial on the Research Topic

### Computational Linguistics and Literature

Computational Linguistics—or, more technically, Natural Language Processing—has made great strides in the past several years. Machine learning (ML) is the technology of choice; deep learning, in particular, has pushed the envelope. These methods work best in narrow domains, given vast amounts of data and asked for information rather than for interpretation. Literary data are not limited by topic, and they are hardly ever plentiful enough. That is why work on such data has been, as yet, on the periphery of research and development in Computational Linguistics. This Research Topic aims to bring the processing of literary data to the attention of a broader audience.

The early versions of the articles we present here have been published at ACL workshops on Computational Linguistics for Literature held in 2015 and 2016<sup>1</sup>. The papers offer a representative enough sample of the research brought to those workshops.

Toral et al. tackle a weakness of an ostensibly successful major application, Machine Translation (MT). The raw results—especially long machine-translated texts, novels above all—are so error-ridden that substantial post-editing is required. Here is the big question: is MT even worthwhile in such cases? In an instructive experiment with a chapter of a novel, the authors show that the overall cost of translation from scratch significantly outstrips the cost of post-editing. They also show that neural MT beats the more traditional phrase-based statistical MT in terms of productivity gains.

Dubremetz and Nivre look at three rhetorical figures in which words are repeated to a strong stylistic effect: chiasmus, epanaphora, and epiphora<sup>2</sup>. While it is trivial to find repeated lexical material, and indeed such repetition is not infrequent, the actual rhetorical devices are very rare and impossible to recognize automatically without sufficient training data for a ML system. The authors run a multi-faceted experiment with literary texts, scientific texts and quotations. They show that various devices predominate in various types of texts. For example, and perhaps a little surprisingly, the seemingly very artistic chiasmus is more likely to appear in scientific texts<sup>3</sup>.

Navarro-Colorado processes a large corpus of Golden Age Spanish sonnets. He uses latent Dirichlet allocation (LDA), a ML algorithm widely accepted as good at determining topics in a text. In addition to a few fairly technical findings, the paper shows that LDA recognizes either a theme in a sonnet, or a poetic motif. The immediate conclusions apply to texts of the same type—moderately short well-structured poems—but the paper can spur similar investigation on similar literary data.

<sup>1</sup> A previous special issue builds upon selected papers from the workshops held in 2012, 2013, and 2014.

<sup>2</sup> See the paper for definitions.

<sup>3</sup> The paper actually studies a type of chiasmus called antimetabole, e.g., “eat to live, not live to eat”.

## OPEN ACCESS

### Edited and reviewed by:

Jean-Gabriel Ganascia,  
Université Pierre et Marie Curie,  
France

### \*Correspondence:

Stan Szpakowicz  
szpak44@gmail.com

### Specialty section:

This article was submitted to  
Digital Literary Studies,  
a section of the journal  
Frontiers in Digital Humanities

**Received:** 23 August 2018

**Accepted:** 06 September 2018

**Published:** 27 September 2018

### Citation:

Szpakowicz S, Feldman A and  
Kazantseva A (2018) Editorial:  
Computational Linguistics and  
Literature. *Front. Digit. Humanit.* 5:24.  
doi: 10.3389/fdigh.2018.00024

Hench and Estes work with highly specialized data: Middle High German epic poetry. We will not steal the authors' thunder; they do a great job of presenting this esoteric matter in a clear and entertaining fashion. In a nutshell, they examine stress patterns in verses, and study how medieval poets used such patterns for semantic and stylistic purposes. Once again, ML is at play. The authors rely on a conditional random field, a method useful in analyzing sequential data.

## AUTHOR CONTRIBUTIONS

All authors listed have made a substantial, direct and intellectual contribution to the work, and approved it for publication.

## ACKNOWLEDGMENTS

This Research Topic owes everything to the authors and the reviewers—our thanks to all. We also thank Cecilia Ovesdotter Alm for an initial review of one of the papers.

**Conflict of Interest Statement:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

*Copyright © 2018 Szpakowicz, Feldman and Kazantseva. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.*