

Ground truth annotation	Model prediction	Example 1
<p><u>[</u>Watching Whoopi : the politics and ethics of the ethics of witnessing. Harris, G. 05/2009 In: Performance Paradigm. 5, 1, p. n/a. Journal article</p> <p><u>[</u>Susan and Darren : the appearance of authenticity. Harris, G. 12/2008 In: Performance Research. 13, 4, p. 4-15. 12 p. Journal article</p>	<p>Watching Whoopi : the politics and ethics of the ethics of witnessing. <u>[</u>Harris, G. 05/2009 In: Performance Paradigm. 5, 1, p. n/a. Journal article</p> <p>Susan and Darren : the appearance of authenticity. <u>[</u>Harris, G. 12/2008 In:] Journal article</p>	

		Example 2
<p>Recent Presentations Keynote, Teaching History in the 21st Century conference, U.C. Berkeley (May 2017) "Digital Storytelling in Higher Education," Wellesley College (March 2017)</p> <p>Dissertations Supervised Sarah Sweetman, Forging Family: Creating and Perpetuating Collective Memory in Families with Children Adopted</p>		
Ground truth annotation		
Model prediction		
<p>Recent Presentations Keynote, Teaching History in the 21st Century conference, U.C. Berkeley (May 2017) <u>[</u>"Digital Storytelling in Higher Education," Wellesley College (March 2017).]</p> <p>Dissertations Supervised <u>[</u>Sarah Sweetman, Forging Family: Creating and Perpetuating Collective Memory in Families with Children Adopted from China (2013).]</p>		

Ground truth annotation	Model prediction	Example 3
<p><u>[</u>Susan and Darren : the appearance of authenticity. Harris, G. 12/2008 In: Performance Research. 13, 4, p. 4-15. 12 p. Journal article</p> <p><u>[</u>How to shop Harris, G. 2007 In: Bobby Baker. London : Routledge p. 191-195. 5 p. Chapter (peer-reviewed)</p>	<p><u>[</u>Susan and Darren : the appearance of authenticity. Harris, G. 12/2008 In: Performance Research. 13, 4, p. 4-15. 12 p. Journal article</p> <p><u>[</u>How to shop Harris, G. 2007 In: Bobby Baker. London : Routledge p. 191-195. 5 p. Chapter (peer-reviewed)</p>	

		Example 4
<p><u>[</u>Rhetoric Roer, H. 2013 The Oxford Guide to the Historical Reception of Augustine. Pollmann, K. (ed.). Oxford: Oxford University Press, Vol. 3, p. 1650-1657 7 p.]</p>		
Ground truth annotation		
Model prediction		
<p>Rhetoric <u>[</u>Roer, H. 2013 The Oxford Guide to the Historical Reception of Augustine. Pollmann, K. (ed.). Oxford: Oxford University Press, Vol. 3, p. 1650-1657 7 p.]</p>		

		Example 5
<p><u>[</u>Burkitt Lymphoma With Pancreatic Involvement JOURNAL OF PEDIATRIC HEMATOLOGY ONCOLOGY Aftandilian, C. C., Friedmann, A. M. 2010; 32 (8): E338-E340 More</p>		
Ground truth annotation		
Model prediction		
<p>Burkitt Lymphoma With Pancreatic Involvement JOURNAL OF PEDIATRIC HEMATOLOGY ONCOLOGY Aftandilian, C. C., Friedmann, A. M. 2010; 32 (8): E338-E340 More</p>		

		Example 6
<p>Books <u>[</u>Unprotected Texts: The Bible's Surprising Contradictions About Sex and Desire By Jennifer Wright Knust HarborOne February 7, 2012 Buy it now from Amazon.com!</p>		
Ground truth annotation		
Model prediction		
<p>Books Unprotected Texts: The Bible's Surprising Contradictions About Sex and Desire By Jennifer Wright Knust HarborOne February 7, 2012 Buy it now from Amazon.com!</p>		

Figure 2: Predicted examples and corresponding ground truth annotations. Underlined tokens are labeled as I (publication). “[” and “]” are not part of the text. They are used to highlight the boundary of publication strings.

string, since it fails to capture dependency relationships in different lines.

Typical error of the webpage-level model: Example 3 is given by the webpage-level model. We see that the webpage-level model can make a more accurate prediction for multi-line publications. However, it may make false positive predictions for short lines (e.g., “Chapter (peer-reviewed)”), while the line-level model seldom makes such mistakes. This is the motivation for us to integrate both the line-level and the webpage-level models.

Typical error of the PubSE model: PubSE can avoid most of the errors shown in Examples 1 to 3. Nevertheless, PubSE still makes mistakes in some challenging cases. Examples 4 to 6 show such cases.

In Example 4, PubSE does not recognize that “Rhetoric” is a publication title. A possible explanation is that such a short publication title is less common.

In Examples 5 and 6, PubSE did not recognize both multi-line publication strings. For Example 5, the reason may be the venue name in upper-case letters and the page number with a letter “E”. For Example 6, the reason may be the unconventional style to present the book, e.g., an exact date of publication and the word “by”, which does not often appear in publication strings. We plan to investigate these problems in the future.

C Choosing Parameter Values: Batch Size of the Webpage-level Model

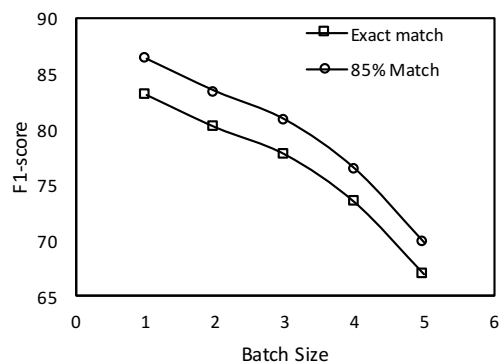


Figure 3: F1-scores of the webpage-level model on the HomePub dataset with different batch sizes

We conduct experiments on a computer with an NVIDIA Tesla K40 GPU which has 12GB GPU memory. Due to the limited GPU memory and

the large size of each webpage (732.1 tokens per webpage on average, with a standard deviation of 1583.3), it is only feasible to set the batch size as one. If the batch size is larger than one, there will be an out-of-memory error. To test the effect of using different batch sizes, we conduct experiments using CPU with 64GB RAM memory, on which we are able to set the batch size to be up to five. The experimental results are shown in Figure 3. The prediction accuracy decreases when the batch size increases. This is mainly because the lengths of different webpages vary greatly, leading to high variance within the same mini-batch. Recent work (Masters and Luschi, 2018) has also shown that small mini-batch size helps to increase the prediction accuracy. As a result, we set the batch size of the webpage-level model to one in the PubSE model.

References

Dominic Masters and Carlo Luschi. 2018. Revisiting small batch training for deep neural networks. *arXiv preprint arXiv:1804.07612*.