



**UvA-DARE (Digital Academic Repository)**

**Riches of the Poor: Using Crummy Entity Linkers for Interactive Search in Digital Humanities**

Olieman, A.M.; Beelen, K.; van Lange, M.; Kamps, J.; Marx, M.J.

[Link to publication](#)

*Citation for published version (APA):*

Olieman, A., Beelen, K., van Lange, M., Kamps, J., & Marx, M. (2017). Riches of the Poor: Using Crummy Entity Linkers for Interactive Search in Digital Humanities. Abstract from 16th Dutch-Belgian Information Retrieval Workshop, Hilversum, Netherlands.

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <http://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

# Riches of the Poor: Using Crummy Entity Linkers for Interactive Search in Digital Humanities \*

Alex Olieman<sup>1,2</sup> Kaspar Beelen<sup>1</sup> Milan van Lange<sup>3</sup> Jaap Kamps<sup>1</sup> Maarten Marx<sup>1</sup>

<sup>1</sup> University of Amsterdam, The Netherlands, {olieman,k.beelen,kamps,maartenmarx}@uva.nl

<sup>2</sup> Stamkracht BV, Amsterdam, The Netherlands, alex@stamkracht.com

<sup>3</sup> NIOD Institute for War, Holocaust, and Genocide Studies, Amsterdam, The Netherlands, m.van.lange@niod.knaw.nl

## 1 INTRODUCTION

The traditional approach for evaluating Entity Linking (EL) systems employs metrics borrowed from Information Retrieval, *precision* and *recall*, to summarize system performance and to aid in their comparison. Consequently, the notion of improving EL technology is closely tied to demonstrating an increased performance in terms of precision and recall. Recent evaluation efforts, utilizing multiple benchmark datasets, have however shown that a good performance on one dataset often does not generalize to others [2]. This limitation of current EL benchmarks, we argue, limits the usefulness of performance statistics for application developers who would like to include EL in their software.

## 2 ENTITY LINKING EVALUATION

There is a poor understanding of the state-of-the-art in Entity Linking [2, 4], which can be partly blamed on the desire for generality in the task definition and its associated objects. How might we then go about demonstrating how effective or useful EL technology can be, in spite of the issues with our metrics and benchmarks? It would be beneficial, in terms of impact and continued funding, to show that EL technology can already be valuable for end-users in its current form.

## 3 SEARCHING FOR ENTITIES BY THE DOZEN

In this paper [3] we introduced WideNet, an entity-based interactive search tool, specifically tailored to assist historians with the exploration of large document collections. Users query the tool with a category of entities and the UI subsequently shows, per subcategory, which entities are mentioned in the corpus, and how frequently, as well as which entities did not occur (see Figure 1). It also displays a list of preview results, showing limited context, to offer quick clues about the relevance of the category.

## 4 CONCLUSION

The main motivation behind the creation of WideNet was to show that (a) even imperfect technologies such as the state-of-the-art Entity Linking systems can support useful niche applications, and (b) tool design should take into account the methodological practices of the scholar. We demonstrated in [3] how the interface supports historians by providing them with a holistic overview of references to complex phenomena such as historical events. Moreover, the grouped presentation of search results with different levels of context dampens frequency in favor of variety: long tail entities become

The screenshot shows the WideNet search interface. On the left, a 'Found' list displays various entities with red circular icons indicating their frequency. The entities listed are: Tachtigjarige Oorlog (15), Remonstranten (10), Zilvervloot (10), Republiek-der-Zeven-Verenigde-Ned. (10), Vrede van Westfalen (10), Tiende Penning (7), Generaliteitslanden (4), Duinkerker kapers (3), Smeekschift der Edelen (3), Eeuwig Edict (1577) (1), Watergeuzen (1), and Niet jaren (Tachtigjarige Oorlog) (1). Below this is a 'Not found' section with 'Abdij Ter Hage' and 'Abdij van Rijnsburg'. On the right, 'Preview results' shows document snippets with entity links. The first result is 'Defensie (begroting)' with a link to 'nl.proc.ob.d.h-1k-20102011-31-42.1.8.1'. The second is '...vooral ook succesvolle strijd tegen de Spaanse overheerser. Hij hervormde de organisatie. ...leken. Daarmee keerde hij het tij in de Tachtigjarige Oorlog. Zo wonnen we de slag bij N...'. The third is 'Grondwetsbepalingen inzake verdediging' with a link to 'nl.proc.ob.d.h-1k-19971998-3243-3279.1.5.3'. The fourth is '...historisch feit dat na beëindiging van de Tachtigjarige Oorlog in 1650 de Staten van Hollan...'. The fifth is '...htigjarige Oorlog in 1650 de Staten van Holland een zeer verregaande alfdanking van mili...'. The sixth is '...r de 350ste verjaardag van de Vrede van Westfalen of van Münster herdenken, mag ook ...'. The seventh is 'Stemmen met stempas' with a link to 'nl.proc.ob.d.h-1k-20082009-7125-7140.1.7.1'. The eighth is '...ntie dat wij met dit stuk wetgeving een eeuwig edict hebben gecreeerd. Het is echter wel...'. The ninth is '15. Europese top' with a link to 'nl.proc.ob.d.h-1k-20122013-84-15.1.11.8'. The tenth is '...en terugnemen? Ik denk dat de Hollandse watergeuzen zich diep schamen voor uw opm...'. The bottom of the interface shows '2 / 31' and 'Tachtigjarige Oorlog' with a green checkmark and a red cross icon.

Figure 1: Assessing the relevance of categories and entities.

more visible while highly frequent but irrelevant results can easily be discarded.

Our general conclusion is that the use case and application crucially determine whether the quality of an EL system is good enough or not: the same EL system may be useless for one task but very useful for another. Hence, paraphrasing a famous paper on early machine translation [1]: there are good applications of crummy entity linkers.

*Acknowledgments.* This research was supported in part by Netherlands Organization for Scientific Research (ExPoSe, NWO CI # 314.99.108; DiLiPaD, NWO Digging into Data # 600.006.014).

## REFERENCES

- [1] Kenneth W. Church and Eduard H. Hovy. 1993. Good Applications for Crummy Machine Translation. *Machine Translation* 8 (1993), 239–258.
- [2] Xiao Ling, Sameer Singh, and Daniel S. Weld. 2015. Design Challenges for Entity Linking. *Transactions of ACL* 3, 2011 (2015), 315–328.
- [3] Alex Olieman, Kaspar Beelen, Milan van Lange, Jaap Kamps, and Maarten Marx. 2017. Good Applications for Crummy Entity Linkers? The Case of Corpus Selection in Digital Humanities. In *Proceedings of SEMANTiCS 2017*. <https://doi.org/10.1145/3132218.3132237> arXiv:1708.01162
- [4] Marieke van Erp, Pablo N. Mendes, Heiko Paulheim, Filip Ilievski, Julien Plu, and Giuseppe Rizzo. 2016. Evaluating Entity Linking: An Analysis of Current Benchmark Datasets and a Roadmap for Doing a Better Job. In *Proceedings of LREC '16*. 4373–4379.

\*This is an extended abstract of [3]