

<https://helda.helsinki.fi>

Extracting Knowledge from Parliamentary Debates for Studying Political Culture and Language

Tamper, Minna

CEUR-WS.org

2022-08-11

Tamper , M , Leal , R , Sinikallio , L , Leskinen , P , Tuominen , J & Hyvönen , E 2022 ,
Extracting Knowledge from Parliamentary Debates for Studying Political Culture and
Language . in S Tiwari , N Mihindukulasooriya & F Osborne, et al. (eds) , Proceedings of the
1st International Workshop on Knowledge Graph Generation From Text and the 1st
International Workshop on Modular Knowledge (TEXT2KG 2022 and MK2022) . CEUR
Workshop Proceedings , vol. 3184 , CEUR-WS.org , Aachen , pp. 70-79 , International
Workshop on Knowledge Graph Generation From Text and the International Workshop on
Modular Knowledge , Hersonissos , Greece , 30/05/2022 .

<http://hdl.handle.net/10138/348159>

cc_by

publishedVersion

Downloaded from Helda, University of Helsinki institutional repository.

This is an electronic reprint of the original article.

This reprint may differ from the original in pagination and typographic detail.

Please cite the original version.

Extracting Knowledge from Parliamentary Debates for Studying Political Culture and Language

Minna Tamper^{1,2}, Rafael Leal², Laura Sinikallio^{1,2}, Petri Leskinen^{2,1},
Jouni Tuominen^{1,2} and Eero Hyvönen^{1,2}

¹University of Helsinki (HELDIG and HSSH), Finland. <https://heldig.fi>

²Aalto University (Semantic Computing Research Group - SeCo), Finland. <https://seco.cs.aalto.fi>

Abstract

This paper presents knowledge extraction and natural language processing methods used to enrich the knowledge graph of the plenary debates (textual transcripts of speeches) of the Parliament of Finland. This knowledge graph includes some 960 000 speeches (1907–2021) interlinked with a prosopographical knowledge graph about the politicians. A recent subset of the speeches was used to extract named entities and topical keywords for semantic searching and browsing the data and for data analysis. The process is based on linguistic analysis, named entity linking, and automatic subject indexing. The results were included into the PARLIAMENTSAMPO knowledge graph in a SPARQL endpoint. This data can be used for studying parliamentary language and culture in Digital Humanities research and for developing applications, such as the PARLIAMENTSAMPO portal.

Keywords

parliamentary studies, natural language processing, linked data, digital humanities

1. Introduction

Parliaments enact new laws, oversee the work of the government, and decide on the state budget. Parliamentary data are used in many areas of research [1], as they provide a wealth of information on the state and functioning of democratic systems, political life and, more generally, language and culture. For these reasons, a lot of parliamentary materials have been digitized in recent decades [2]. Digitized parliamentary materials offer a wide range of perspectives on different research topics and have been used in a variety of fields, such as linguistics, political science, economics, and history. A most important research material for parliament studies are the debates in the parliaments, i.e., sequences of transliterated speeches (minutes) of Members of Parliament (MP) and other politicians, through which one can study the language and its changes itself as well as the underlying societal phenomena at large [3].

This paper argues and shows that by enriching textual parliamentary speeches with linked


Text2KG 2022: International Workshop on Knowledge Graph Generation from Text, Co-located with the ESWC 2022, May 05-30-2022, Crete, Hersonissos, Greece

✉ minna.tamper@aalto.fi (M. Tamper); rafael.leal@aalto.fi (R. Leal); laura.sinikallio@helsinki.fi (L. Sinikallio); petri.leskinen@aalto.fi (P. Leskinen); jouni.tuominen@helsinki.fi (J. Tuominen); eero.hyvonen@aalto.fi (E. Hyvönen)

🆔 0000-0002-3301-1705 (M. Tamper); 0000-0001-7266-2036 (R. Leal); 0000-0001-7398-6585 (L. Sinikallio); 0000-0003-2327-6942 (P. Leskinen); 0000-0003-4789-5676 (J. Tuominen); 0000-0003-1695-5840 (E. Hyvönen)



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

data using knowledge extraction methods [4], it is possible to support Digital Humanities (DH) research and enhance the usability of the data in applications, such as semantic search, browsing, and data analysis. As a case study, a part the ca. 960 000 speeches of the system *ParliamentSampo – Finnish Parliament on the Semantic Web* [5, 6] are used. In an earlier work, the speeches covering the whole history 1907–2021 of the Parliament of Finland (PoF) were extracted from original heterogeneous data sources and transformed into a speech knowledge graph (S-KG) [7] (and also into the Parla-CLARIN format¹). At the same time, the S-KG was interlinked with a prosopographical KG (P-KG) representing detailed biographical data and networks of the ca. 2800 MPs and politicians involved in the PoF activities [8]. Both graphs were published as a LOD service on the Linked Data Finland platform LDF.fi [9], including a SPARQL endpoint². In this paper, the textual speeches of this PARLIAMENTSAMPO dataset are enriched further using knowledge extraction techniques in order to support DH [10] analysis and for further development of the semantic portal PARLIAMENTSAMPO on top of the endpoint.

In this paper, we first shortly overview the related work (Section 2) followed by the description of the speech data and then focus on the new data enrichments using Natural Language Processing (NLP) methods (Section 3). Section 4 discusses how the new data can be utilized in the PARLIAMENTSAMPO portal. Lastly, the contributions of this work are summarized and discussed (Section 5).

2. Related Work

Regarding the digitization of parliamentary data, plenary debates have been in central role, e.g., [11] and the CLARIN list of parliamentary corpora³ in different countries. Parliamentary materials have also been transformed into linked data, too. A prominent example of this is the LinkedEP [12] system on the European Parliament’s data. Linked data has also been used in the Italian Parliament⁴, and the LinkedSaeima for the Latvian parliament [13] in addition to the Finnish ParliamentSampo system [5, 6] whose data [7, 8] was re-used in the paper.

Knowledge extraction has been applied to enrich datasets to enable distant reading approaches to studying parliamentary debates. For example, the Latvian LinkedSaeima dataset has utilized named entity linking (NEL) to enrich their metadata. Similarly, the Dutch parliamentary debates dataset [14] has been enriched with named entities (NE). The Slovenian siParl corpus [15] includes with linguistic information about the parliamentary debates in CONLL-U format.

The NLP methods used in this work have been developed mainly for handling Finnish texts. With respect to NER, some of the most relevant tools are StanfordNER [16], FiNER [17] and the FinBERT based NER tool, of which the last one is currently estimated to be the most accurate [18, 19, 20]. Similarly, there are morphological analyzers for Finnish besides the Turku Neural parser, such as the two used in this paper, Voikko and uralicNLP, the latter of which employs Omorfi [21]. Regarding entity linking, there are few tools available for Finnish, such as ARPA [22]. Various tools have been created also for other languages to link NEs to different

¹<https://github.com/clarin-eric/parla-clarin>

²The data will be published openly using the CC BY 4.0 license by the end of 2022.

³<https://www.clarin.eu/resource-families/parliamentary-corpora>

⁴<http://data.camera.it>

datasets, such as [23, 24, 25].

In Finland, parliamentary materials have been digitized and utilized to some extent in DH and social science research. For example, [26] examines the differences in political speech between parties throughout the parliamentary period 1907–2018. In [27], the content of the plenary speeches given in Parliament in 1999–2014 were studied by using topic modeling. Also, in [28] the debates were examined. However, data have so far been used only in a few studies that deploy methods from corpus linguistics, language technology, or computer science [2].

Previous search applications for the Finnish parliamentary speech data are based mostly on traditional text search. However, search applications have been developed for other digitized and enriched Cultural Heritage datasets [29, 30]. The data analysis tools to examine the results are few, such as the concordance analysis of the Language Bank of Finland⁵, where the words are visualized in their textual contexts and show some statistics of occurrences in the search results. The Language Bank’s tool has many corpora and one small corpus covering a small part of the entire time series of the Finnish parliamentary speeches.

3. Datasets and Knowledge Extraction

3.1. Core Datasets

The PARLIAMENTSAMPO system includes data about the MPs, parliamentary speeches, and political organizations within the PoF. The data covers also the comments of the Speaker (President) of the PoF and all other small comments recorded in the minutes, e.g., in connection with voting proceedings. The PARLIAMENTSAMPO data contains two major parts:

1. The Prosopographic Knowledge Graph The Prosopographic Knowledge Graph [8] covers all MPs of Finland since the year 1907. At its core lies a RDF conversion of data about MPs from the originally XML-formatted Open Data service⁶ of PoF. In addition to basic information, such as times and places of birth and death, the data includes detailed information about politicians’ life events, such as studies, working life, political career, and their written publications. In addition to people, the graph contains information about organizations, professions, and positions, as well as places. Organizations include, e.g., parties, ministries, parliamentary groups, committees, and constituencies, as well as schools, organizations, and companies outside the political community.

2. The Parliamentary Speeches Knowledge Graph The knowledge graph of parliamentary speeches contains speeches collected from all the minutes of the plenary sessions of the PoF since 1907 [7]. This knowledge graph was compiled from the documents available on the Open Data services⁷ and web sites⁸ of the PoF. Depending on the time period they covered, the documents were available in different formats: PDF, HTML, or XML. PDF documents were transformed into text with OCR.

In addition to the actual speeches, the speech graph contains all the relevant metadata attached to the minutes, such as interjections, information about the session where the speech

⁵<https://www.kielipankki.fi/support/access/>

⁶<https://avoindata.eduskunta.fi/#/fi/dbsearch>

⁷<https://avoindata.eduskunta.fi/#/fi/home>

⁸<https://www.eduskunta.fi/fi/Sivut/default.aspx>

was given (time, date, serial number, etc.), speaker information (name, role, party) and possible topic of discussion, and supporting documents (e.g. committee report). Based on the metadata, the speeches were linked to the MPs P-KG. For example, speakers and the parties they represent are resources with URI identifiers described in the P-KG.

3.2. Knowledge Extraction

In our work, the speech knowledge graph was enriched using various NLP methods. The toolset that was used in enriching the BiographySampo dataset [31, 32] was re-used together with new methods for NER, lemmatization, and automatic subject indexing. The parliamentary debates were enriched with named entity recognition (NER) and linking, subject indexing, and by creating a linguistic knowledge graph containing linguistic details for the speeches. Here, NLP methods were used on a subset of the speeches dataset, consisting of speeches from parliamentary session 2015 to the end of parliamentary session 2021, totaling in a little over 114 000 speeches, covering about 12% of the speeches dataset.

Lemmatization and Subject Indexing The *Secompling*⁹ was used for the tasks of lemmatization and subject indexing. It is an under-development library, which aims at integrating different Finnish NLP tools.

Lemmatization can be seen of as a kind of text normalization, especially for a language as morphologically rich as Finnish, which has 15 inflectional cases and a rich system for derivative words. Lemmatization enables exact term-based search instead of wildcard-based stemming. Lemmatization allows word count-based algorithms, such as TF-IDF, to work with more precision. *Secompling* employs the Turku Neural parser pipeline [33, 34] for lemmatization, and *Voikko*¹⁰ and *uralicNLP* [35] to check and possibly fix errors regarding these base forms. The *Secompling* lemmatization module has not been formally evaluated yet.

Subject indexing allows texts to be described succinctly by focusing on keywords that best characterize their contents. In our work, the subject indexing tool *Annif* [36], developed by the National Library of Finland, is used for this task. As *Annif* is capable of using machine-learning-based correlational backends such as *Parabel*, it may sometimes suggest NEs not mentioned in the texts. Since we focus on entities that are actually mentioned in the speeches, NEs from *Annif* were ignored. The other subject keywords are filtered out according to their weight. The keywords provided by *Annif* are entities from the General Finnish Ontology YSO¹¹ – which is part on the national Finnish LOD infrastructure [37] – with ready-to-use URIs for data linking.

A total of 10467 keywords were identified for this dataset, with an average of 23.23 keywords per text, a maximum of 78, a minimum of 1 and a standard deviation of 8.46. The most common subjects were *poliitikot* ‘politicians’ (around 49% of the texts), *ministerit* ‘ministers’ (ca. 45%), *kunnat* ‘municipalities’ (ca. 42%) and *lainsäädäntö* ‘legislation’ (ca. 39%).

Named Entity Recognition and Linking NEL was performed on the speeches to improve data browsing and searching in the PARLIAMENTSAMPO portal. Similarly to the analytics done for the textual biographies in the BiographySampo [32] system, NEL enables more detailed data analytics in the PARLIAMENTSAMPO dataset, too. NEs were extracted using the *Nelli* [38, 39]

⁹<https://version.aalto.fi/gitlab/seco/secompling>

¹⁰<https://voikko.puimula.org/>

¹¹<https://finto.fi/yso/fi/new?clang=en>

tool and its results linked using the ARPA. Unlike in the BiographySampo dataset, here Nelli was configured to use FinBERT’s combined NER model [18], Reksi [38], and the Turku Neural parser pipeline. FinBERT’s NER tool that is coupled with Reksi to pick up links to legislation, references to various dates, and identifiers (e.g., URLs). The Turku Neural parser was selected for morphological analysis based on its performance [33]. These tools extracted entities that were later linked using ARPA to the PARLIAMENTSAMPO dataset and to other external datasets, such as the Kanto ontology of Finnish actors¹², the Place Name Register PNR ontology of contemporary Finnish places¹³, and the YSO places¹⁴ ontology that contains also historical places mentioned in the speeches.

The tools used for NER managed to extract NEs from 89% of the speeches of which 30% contained people, 19% mentions of time, 12% organizations, and 7% of places. However, the linking of entities requires still some work. For example, the full name references to MPs were linked while the surname references were not. The place mentions were linked mostly correctly, however, the target ontologies lacked some mentioned place names like *Wuhan*.

Morpholinguistic Knowledge Graph Lastly, the speeches were transformed into a separate morpholinguistic knowledge graph (MLKG) containing detailed linguistic and morphological information about the speeches using a pipeline previously used for BiographySampo [40]. This graph can be used for linguistic analysis of the parliamentary speeches similarly to the work done in BiographySampo [32]. For example, in the BiographySampo dataset, it was noticed that biographies of women contained more family-related terminology while biographies about men used more words related to war and religion. In order to apply same methods to parliamentary speeches, a similar pipeline was used, updated to use the Turku Neural parser pipeline, and adjusted to handle larger datasets in smaller chunks of text. In this case, it was configured to process data by year. The results are also linked to the PARLIAMENTSAMPO speeches dataset to enable analysis of speeches using the speech metadata.

4. Using the Enriched Data in PARLIAMENTSAMPO

The enriched PARLIAMENTSAMPO data is used in the development of the PARLIAMENTSAMPO portal [6] which is based on the Sampo model [41] and the Sampo-UI framework [42]. The portal demonstrates how the data service can be used for developing applications for DH research. In this application the data can be browsed using ontology-based faceted search, and the results can then be analyzed with the integrated visualization and data analysis tools.

The enriched data is initially used to boost the browsing and searching capabilities of the portal. For instance, NEs and keywords can be used via facets to find speeches that mention a specific topic or a NE, such as a place or an organization. Coupled with the speeches, their metadata, and the prosopographical data, this enables studying, e.g., how MPs talk about matters related to their constituency. At the moment mentioned organizations and places have already been added into portal as facets to test the data.

Currently, the MLKG about the speeches is limited to a few years. It is not yet included in the

¹²<https://finto.fi/finaf/fi/>

¹³<https://www.ldf.fi/dataset/pnr/>

¹⁴<https://finto.fi/yso-paikat/en/>

PARLIAMENTSAMPO portal, but we plan to add it and develop similar linguistic analysis views as in BiographySampo. This enables, e.g., to study the vocabulary used by the MPs and parties in their speeches. It is also possible to compare differences in vocabulary of men and women.

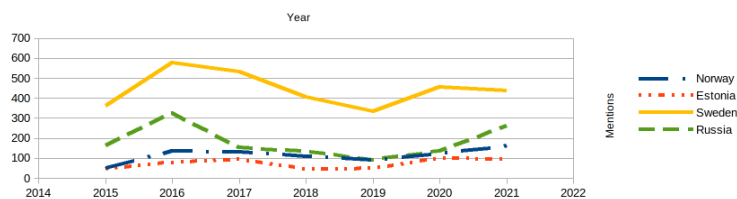


Figure 1: Frequency of place mentions in parliamentary debates from the 2015 parliamentary session until end of 2021.

In addition, the SPARQL endpoint underlying the PARLIAMENTSAMPO portal can be used for querying, analyzing, and visualizing the enriched data. In Fig. 1, e.g., the speeches mentioning Finland’s neighbouring countries Norway, Sweden, Estonia, and Russia are counted on a yearly basis and plotted from 2015 to 2021. The plot shows that Sweden appears more frequently than the other neighbours. Russia is also mentioned increasingly in 2020 and 2021. Based on initial analysis of the speeches mentioning Russia and Sweden, the discussions and their frequencies are related to topics such as the annexing of Crimea, managing good relations, defence and security of Finland and its nearby areas, such as the Baltic sea. These mentions reflect the working order of the parliament, the domestic and world events described in the media at the time. It remains as future work to study the context of these mentions in more detail. Similarly, by linking to place ontologies it is possible to leverage the benefit of organized information to create visualizations that cluster all, e.g., Russia-related place names as mentions about Russia. It also enables the use of map-based visualizations.

5. Discussion

In this paper, we presented work for enriching the Finnish parliamentary debate corpus to support data browsing and using it for DH research. This is ongoing work that still requires adjustments and extensive evaluation, similarly, the PARLIAMENTSAMPO portal is still under development. The tools used in the enrichment have been previously evaluated with different corpora, but not for the parliamentary data. The FinBERT NER tool has achieved an accuracy of 93.11% using the combined model in cross-corpus evaluation [18]. Similarly, the Turku Neural Parser pipeline is evaluated based on CoNLL 2018 UD Shared Task [43] with accuracy of LAS¹⁵ 86.60%, UPOS¹⁶ 96.66%, and XPOS¹⁷ 97.63% [33]. Subject indexing is difficult to evaluate, however based on evaluation of the Annif tool, its accuracy is 30–50% depending on the test corpus [44]. These results have been produced on formal Finnish language texts similar to the Finnish parliamentary debates corpus.

¹⁵Labeled attachment score (LAS) is the proportion of words that have connected correctly the head word with the right dependency relation.

¹⁶Universal part- of-speech tagging

¹⁷Language-specific part-of-speech tagging

The data has been partially added to the PARLIAMENTSAMPO knowledge graph and utilized already in the facets of the semantic portal. The enriched data enables DH research through topics and NEs. The enrichments help to find interesting phenomenon in the PARLIAMENTSAMPO dataset. Similarly to the Finnish dataset, NEs have been added, e.g., into the Dutch and Latvian parliamentary debate corpora. The linked NEs and keywords enable data analytics and search optimization in the faceted search application. The MLKG contains millions of triples of morphological and linguistic information as linked data. Unlike, e.g., Slovenian debate corpora, the Finnish dataset can be queried directly using SPARQL to analyze speeches using also the metadata. However, due to size of dataset, there is still much work to be done to speed up the queries. It remains future work to create applications for the DH community to enable to study the debates in more detail.

Acknowledgements Our work is part of the Semantic Parliament project¹⁸, funded by the Academy of Finland and is also related to the EU project InTaVia¹⁹ and the EU COST action Nexus Linguarum²⁰. The project uses the computing resources of the CSC – IT Center for Science.

References

- [1] C. Benoît, O. Rozenberg (Eds.), *Handbook of Parliamentary Studies: Interdisciplinary Approaches to Legislatures*, Edward Elgar Publishing, 2020. doi:10.4337/9781789906516.
- [2] M. Andrushchenko, K. Sandberg, R. Turunen, J. Marjanen, M. Hatavara, J. Kurunmäki, T. Nummenmaa, M. Hyvärinen, K. Teräs, J. Peltonen, J. Nummenmaa, Using parsed and annotated corpora to analyze parliamentarians' talk in Finland, *Journal of the Association for Information Science and Technology* 185 (2021) 1–15. doi:10.1002/asi.24500.
- [3] K. Elo, J. Karimäki, Luonnonsuojelusta ilmastopoliittikkaan: Ympäristöpoliittisen käsitteistön muutos parlamenttipuheessa 1960–2020, *Politiikka* 63 (2021). doi:10.37452/politiikka.109690.
- [4] J. L. Martinez-Rodriguez, A. Hogan, I. Lopez-Arevalo, Information Extraction Meets the Semantic Web: A Survey, *Semantic Web – Interoperability, Usability, Applicability* 11 (2020) 255–335. doi:10.3233/SW-180333.
- [5] E. Hyvönen, L. Sinikallio, P. Leskinen, S. Drobac, J. Tuominen, K. Elo, M. L. Mela, M. Koho, E. Ikkala, M. Tamper, R. Leal, J. Kesäniemi, Parlamenttisampo: eduskunnan aineistojen linkitetyn avoimen datan palvelu ja sen käyttömahdollisuudet, *Informaatiotutkimus* 40 (2021). doi:10.23978/inf.107899.
- [6] E. Hyvönen, L. Sinikallio, P. Leskinen, M. L. Mela, J. Tuominen, K. Elo, S. Drobac, M. Koho, E. Ikkala, M. Tamper, R. Leal, J. Kesäniemi, Finnish parliament on the semantic web: Using parlamentsampo data service and semantic portal for studying political culture and language, in: *Digital Parliamentary data in Action (DiPaDa 2022)*, Workshop at the 6th Digital Humanities in Nordic and Baltic Countries Conference, long paper, CEUR Workshop Proceedings, Vol. 3133, 2022. URL: <http://ceur-ws.org/Vol-3133/paper05.pdf>.

¹⁸<https://seco.cs.aalto.fi/projects/sem parl/en/>

¹⁹<https://intavia.eu>

²⁰<https://nexuslinguarum.eu>

- [7] L. Sinikallio, S. Drobac, M. Tamper, R. Leal, M. Koho, J. Tuominen, M. La Mela, E. Hyvönen, Plenary Debates of the Parliament of Finland as Linked Open Data and in Parla-CLARIN Markup, in: 3rd Conference on Language, Data and Knowledge (LDK 2021), volume 93, 2021, pp. 8:1–8:17. doi:10.4230/OASICS.LDK.2021.8.
- [8] P. Leskinen, E. Hyvönen, J. Tuominen, Members of Parliament in Finland Knowledge Graph and Its Linked Open Data Service, in: Further with Knowledge Graphs. Proceedings of the 17th International Conference on Semantic Systems, 6-9 September 2021, Amsterdam, The Netherlands, 2021, pp. 255–269. doi:10.3233/SSW210049.
- [9] E. Hyvönen, J. Tuominen, M. Alonen, E. Mäkelä, Linked Data Finland: A 7-star Model and Platform for Publishing and Re-using Linked Datasets, in: The Semantic Web: ESWC 2014 Satellite Events, Revised Selected Papers, Springer-Verlag, 2014, pp. 226–230. URL: https://doi.org/10.1007/978-3-319-11955-7_24.
- [10] E. Gardiner, R. G. Musto, The Digital Humanities: A Primer for Students and Scholars, Cambridge University Press, New York, NY, USA, 2015. <https://doi.org/10.1017/CBO9781139003865>.
- [11] E. Lapponi, M. G. Søyland, E. Velldal, S. Oepen, The Talk of Norway: a richly annotated corpus of the Norwegian parliament, 1998–2016, Lang Resources & Evaluation 52 (2018) 873–893. doi:10.1007/s10579-018-9411-5.
- [12] A. Van Aggelen, L. Hollink, M. Kemman, M. Kleppe, H. Beunders, The debates of the European Parliament as Linked Open Data, Semantic Web – Interoperability, Usability, Applicability 8 (2017) 271–281. doi:10.1007/s42001-019-00060-w.
- [13] U. Bojārs, R. Dargis, U. Lavrinovičs, P. Paikens, LinkedSaeima: A Linked Open Dataset of Latvia’s Parliamentary Debates, in: Semantic Systems. The Power of AI and Knowledge Graphs. SEMANTiCS 2019, Springer, 2019, pp. 50–56. doi:10.1007/978-3-030-33220-4_4.
- [14] D. Juric, L. Hollink, G.-J. Houben, Bringing Parliamentary Debates to the Semantic Web., in: Proceedings of the Workshop on Detection, Representation, and Exploitation of Events in the Semantic Web (DeRiVE 2012), 2012, pp. 51–60.
- [15] A. Pancur, T. Erjavec, The siParl corpus of Slovene parliamentary proceedings, in: Proceedings of the Second ParlaCLARIN Workshop, 2020, pp. 28–34.
- [16] J. R. Finkel, T. Grenager, C. D. Manning, Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling, in: Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05), June, 25-30, 2005, University of Michigan, Ann Arbor, Michigan, USA, Association for Computational Linguistics, 2005, pp. 363–370. doi:10.3115/1219840.1219885.
- [17] T. Ruokolainen, P. Kauppinen, M. Silfverberg, K. Lindén, A Finnish news corpus for named entity recognition, Language Resources and Evaluation 54 (2020) 247–272. doi:10.1007/s10579-019-09471-7.
- [18] J. Luoma, M. Oinonen, M. Pyykönen, V. Laippala, S. Pyysalo, A Broad-coverage Corpus for Finnish Named Entity Recognition, in: Proceedings of The 12th Language Resources and Evaluation Conference, 2020, pp. 4615–4624.
- [19] A. Virtanen, J. Kanerva, R. Ilo, J. Luoma, J. Luotolahti, T. Salakoski, F. Ginter, S. Pyysalo, Multilingual is not enough: BERT for Finnish, 2019. arXiv:1912.07076.
- [20] T. Ruokolainen, K. Kettunen, À la recherche du nom perdu–Searching for Named Entities

- with Stanford NER in a Finnish Historical Newspaper and Journal Collection, in: 13th IAPR International Workshop on Document Analysis Systems, 2018.
- [21] T. A. Pirinen, Development and Use of Computational Morphology of Finnish in the Open Source and Open Science Era: Notes on Experiences with Omorfi Development., *SKY Journal of Linguistics* 28 (2015) 381–393. doi:10.23978/inf.107890.
- [22] E. Mäkelä, Combining a REST Lexical Analysis Web Service with SPARQL for Mashup Semantic Annotation from Text, in: *The Semantic Web: ESWC 2014 Satellite Events - ESWC 2014 Satellite Events*, Anissaras, Crete, Greece, May 25-29, 2014, Revised Selected Papers, Springer International Publishing, 2014, pp. 424–428. doi:10.1007/978-3-319-11955-7_60.
- [23] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, Z. Ives, DBpedia: A Nucleus for a Web of Open Data, in: *The Semantic Web: 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007 + ASWC 2007*, Busan, Korea, November 11-15, 2007, Springer, Berlin, Heidelberg, 2007, pp. 722–735. doi:10.1007/978-3-540-76298-0_52.
- [24] D. Damjanovic, K. Bontcheva, Named Entity Disambiguation using Linked Data, in: *The Semantic Web: Research and Applications - 9th Extended Semantic Web Conference, ESWC 2012*, Heraklion, Crete, Greece, May 27-31, 2012. Proceedings, Springer-Verlag Berlin Heidelberg, 2012, pp. 231–240.
- [25] L. Derczynski, D. Maynard, G. Rizzo, M. Van Erp, G. Gorrell, R. Troncy, J. Petrak, K. Bontcheva, Analysis of named entity recognition and linking for tweets, *Information Processing and Management* 51 (2015) 32–49. doi:10.1016/j.ipm.2014.10.006.
- [26] S. Simola, A century of partisanship in Finnish political speech, 2020. Part of PhD thesis: *Essays in Labor and Political Economics*, Aalto University.
- [27] K. Makkonen, P. Loukasmäki, Eduskunnan täysistunnon puheenaiheet 1999–2014: Miten käsitellä LDA-aihemalleja?, *Politiikka* 61 (2019) 127–159.
- [28] E. Lillqvist, I. K. Kavonius, M. Pantzar, “Velkakello tikittää”: Julkisyhteisöjen velka suomalaisessa mielikuvastossa ja tilastoissa 2000–2020, *Kansantaloudellinen Aikakauskirja* 116 (2020) 581–607.
- [29] E. Hyvönen, E. Ikkala, M. Koho, R. Leal, H. Rantala, M. Tamper, How to search and contextualize scenes inside videos for enriched watching experience: Case stories of the second world war veterans, 2022. Under peer review.
- [30] A. Brandsen, S. Verberne, K. Lambers, M. Wansleeben, Can BERT Dig It?–Named Entity Recognition for Information Retrieval in the Archaeology Domain, arXiv preprint arXiv:2106.07742 (2021).
- [31] E. Hyvönen, P. Leskinen, M. Tamper, H. Rantala, E. Ikkala, J. Tuominen, K. Keravuori, BiographySampo - publishing and enriching biographies on the semantic web for digital humanities research, in: P. Hitzler, M. Fernández, K. Janowicz, A. Zaveri, A. J. Gray, V. Lopez, A. Haller, K. Hammar (Eds.), *The Semantic Web. ESWC 2019*, Springer-Verlag, 2019, pp. 574–589. doi:10.1007/978-3-030-21348-0_37.
- [32] M. Tamper, P. Leskinen, E. Hyvönen, R. Valjus, K. Keravuori, Analyzing Biography Collection Historiographically as Linked Data: Case National Biography of Finland, *Semantic Web – Interoperability, Usability, Applicability* (2021). Accepted.
- [33] J. Kanerva, F. Ginter, N. Miekka, A. Leino, T. Salakoski, Turku Neural Parser Pipeline:

- An End-to-End System for the CoNLL 2018 Shared Task, in: Proceedings of the CoNLL 2018 Shared Task: Multilingual parsing from raw text to universal dependencies, 2018, pp. 133–142. doi:10.18653/v1/K18-2013.
- [34] J. Kanerva, F. Ginter, T. Salakoski, Universal Lemmatizer: A sequence-to-sequence model for lemmatizing Universal Dependencies treebanks, *Natural Language Engineering* (2020) 1–30. doi:10.1017/S1351324920000224.
- [35] M. Hämäläinen, UralicNLP: An NLP library for Uralic languages, *Journal of Open Source Software* 4 (2019) 1345. doi:10.21105/joss.01345.
- [36] O. Suominen, Annif: DIY automated subject indexing using multiple algorithms, *LIBER Quarterly* 29 (2019) 1–25. doi:10.18352/lq.10285.
- [37] E. Hyvönen, How to create a national cross-domain ontology and linked data infrastructure and use it on the semantic web (2021). URL: <https://seco.cs.aalto.fi/publications/2021/hyvonen-dcmi-2021.pdf>, keynote presentation for the DCMI 2021 conference.
- [38] M. Tamper, A. Oksanen, J. Tuominen, A. Hietanen, E. Hyvönen, Automatic Annotation Service APPI: Named Entity Linking in Legal Domain, in: *The Semantic Web: ESWC 2020 Satellite Events*, Springer-Verlag, 2020, pp. 110–114. doi:10.1007/978-3-030-62327-2_36.
- [39] M. Tamper, E. Hyvönen, P. Leskinen, Visualizing and analyzing networks of named entities in biographical dictionaries for digital humanities research, in: *Proceedings of the 20th International Conference on Computational Linguistics and Intelligent Text Processing (CICling 2019)*, Springer, 2019. Accepted.
- [40] M. Tamper, P. Leskinen, K. Apajalahti, E. Hyvönen, Using Biographical Texts as Linked Data for Prosopographical Research and Applications, in: *Digital Heritage. Progress in Cultural Heritage: Documentation, Preservation, and Protection. 7th International Conference, EuroMed 2018, Nicosia, Cyprus*, Springer-Verlag, 2018, pp. 125–137. doi:10.1007/978-3-030-01762-0_11.
- [41] E. Hyvönen, Digital humanities on the Semantic Web: Sampo model and portal series, 2021. Submitted.
- [42] E. Ikkala, E. Hyvönen, H. Rantala, M. Koho, Sampo-UI: A full stack JavaScript framework for developing semantic portal user interfaces, *Semantic Web – Interoperability, Usability, Applicability* 13 (2022) 69–84. doi:10.3233/SW-210428.
- [43] D. Zeman, J. Hajič, M. Popel, M. Potthast, M. Straka, F. Ginter, J. Nivre, S. Petrov, CoNLL 2018 shared task: Multilingual parsing from raw text to Universal Dependencies, in: *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, Association for Computational Linguistics, Brussels, Belgium, 2018, pp. 1–21. doi:10.18653/v1/K18-2001.
- [44] O. Suominen, M. Lehtinen, J. Inkinen, Annif and Finto AI: Developing and Implementing Automated Subject Indexing, *Jlis.it* 13 (2022) 265–282. doi:10.4403/jlis.it-12740.