# The Danish Gigaword Project

**Leon Strømberg-Derczynski,**[1] **Rebekah Baglini,**[2] **Morten H. Christiansen,**[2,3]
**Manuel R. Ciosici,**[1] **Jacob Aarup Dalsgaard,**[2] **Riccardo Fusaroli,**[2]
**Peter Juel Henrichsen,**[4] **Rasmus Hvingelby,**[5] **Andreas Kirkedal,**[1,6]
**Alex Speed Kjeldsen,**[7] **Claus Ladefoged,**[8] **Finn Årup Nielsen,**[9]
**Malte Lau Petersen,**[2] **Jonathan Hvithamar Rystrøm,**[2] **Daniel Varab.**[1,10]

[1]:IT University of Copenhagen; [2]:Aarhus University; [3]:Cornell University;
[4]:Danish Language Council; [5]:Alexandra Institute; [6]:Interactions LLC;
[7]:University of Copenhagen; [8]:TV2 Regionerne; [9]:Technical University of Denmark;
[10]:Novo Nordisk.

**Abstract**

Danish is a North Germanic/Scandinavian language spoken primarily in Denmark, a country with a tradition of technological and scientific innovation. However, from a technological perspective, the Danish language has received relatively little attention and, as a result, Danish language technology is hard to develop, in part due to a lack of large or broad-coverage Danish corpora. This paper describes the Danish Gigaword project, which aims to construct a freely-available one billion word corpus of Danish text that represents the breadth of the written language.

**Keywords:** Danish, gigaword, corpus

## 1. Introduction

It is hard to develop good general-purpose language processing tools without a corpus that is broadly representative of the target language. Further, to develop high-performance deep learning models requires hundreds of millions of tokens. Currently, no single such dataset exists for the Danish language. To address this gap, then, we propose an open corpus that can be freely downloaded and used by anyone wishing to work on Danish NLP. This means language researchers' access is not limited to the byproducts of larger projects that are provided bilingually; and also that anyone wanting to work on Danish will be able to, without having to pay large licensing fees, which exclude organisations from developing Danish NLP, thus restricting Danish-speakers from receiving the many benefits of this powerful range of technologies.

This paper details the Danish Gigaword project, DAGW, which aims to build a billion-word corpus representing the language across various dimensions, including modality, time, setting, and place.

Collecting such a corpus automatically is tricky: automatic language identification tools confound closely related texts, especially Danish and Bokmål, and are likely to miss important data. Nevertheless, is it important, and existing representations underperform for Danish: the multilingual FastText embeddings (Joulin et al., 2018) miss core Danish words such as "træls" (Hovy et al., 2015); Multilingual BERT has little to no support for the normal Danish vowel "å".[1]

To remedy this situation, we propose the construction of a Danish Gigaword corpus. The overriding goals are to create a dataset that is (1) representative; (2) accessible; (3) a general-purpose corpus for NLP on Danish.

## 2. Background

Several Danish text corpora of varying size and composition have been compiled during recent decades.

CLARIN-DK offers a variety of individual corpora, of varying genres, annotations, and writing times. However, uncertainty around redistribution makes it non-trivial to bundle this content into a dataset simple to work with at scale in both research and industry; to do so would be an inappropriate repurposing of this high-quality repository. Rather, the CLARIN data – or its accompanying annotations as distributed – is not always licensed for commercial reuse, and so partly at odds with this project's open goals.

Similarly, other huge monolithic datasets such as the Common Crawl Danish data suffer from the inclusion of significant amounts of non-Danish content, in part due to the pervasive confusion between Danish and Norwegian Bokmål by highly-multilingual language ID classifiers (Haas, 2019). Further, these datasets often have a bias toward content from recent years, leaving models built over them sub-optimally prepared to process older Danish.

Some major corpora are related to dictionary production, as is the case for the 56 million word Korpus-DK available for search at the dictionary site ordnet.dk.[2] Leipzig Corpora Collection[3] assembles Danish corpora from the Web, news sites and the Danish Wikipedia (Goldhahn et al., 2012). Distributed under the Creative Commons Attribution license, the largest corpora have one million sentences, each corpus approximately corresponding to 20 million tokens. The Dasem Python package[4] seeks to assemble open corpora and has been used for training word2vec models with the Leipzig Corpora Collection, Gutenberg Project, Europarl and Wikipedia corpora (Nielsen and Hansen, 2017).

---

[1] BotXO maintains a Danish BERT instance at https://github.com/botxo/nordic_bert

[2] http://ordnet.dk
[3] http://wortschatz.uni-leipzig.de/en/download
[4] http://github.com/fnielsen/dasem

The combined size of these corpora is still far from one billion words, and, as a consequence, they do not directly propel research and development in Danish NLP forward with the speed that is required to keep up with the international NLP community. Danish NLP tools lag behind NLP tools for better-resourced languages, such as English, and the gap is continuously increasing.

Researchers and developers have been painfully aware of this deficiency for years, and the problem has been addressed in several language policy reports over the years, such as the Danish METANET White Paper, "Danish in the Digital Age" (Pedersen et al., 2012), and more recently, in the government report "Dansk Sprogteknologi i verdensklasse" ('world-class Danish language technology') (Kirchmeier et al., 2019). More recently, it has been decided to support the production/collection of such basic resources as part of a governmental AI strategy[5], but concrete development plans remain pending.

The first gigaword Corpus was the English Gigaword (Graff et al., 2003). It consisted of roughly one billion ($10^9$) words of English-language newswire text from four major sources: Agence France Press, Associated Press Worldwide, New York Times, and Xinhua English. These, in turn, had largely been previously published as smaller corpora in their own right. The content was single-genre, national and global newswire, published between 1994 and 2002.

Other gigaword corpora emerged later, for French, Arabic, Chinese, and Spanish. And other projects are reaching the required goal for yet more languages; even Icelandic, a language with just over 360,000 speakers, has a healthy gigaword project (Steingrímsson et al., 2018).

## 3. Language diversity

If a corpus dataset is to be useful for a wide range of applications, it must include a wide range of language. This means mixing domains, mixing speakers, and mixing styles. Failing to do this can lead to serious deficiencies in the data. For example, when NLP work started on social media text, the Wall Street Journal-trained part of speech taggers missed key words such as "Internet" – due to the articles being from the late eighties and early nineties – and "bake", due to their domain (Derczynski et al., 2013). This does not form a strong basis for general-purpose NLP, and so it will be crucial to capture and distribute as broad a range of Danish as possible in the Danish Gigaword.

## 4. Dataset construction

The Danish Gigaword Corpus consists of sections, with each section corresponding to a single source of text. Following prior efforts to construct broad-coverage datasets (Derczynski et al., 2016), sections are selected based on how well they help the corpus' coverage of Danish over a variety of dimensions, including: time of authorship; speech situation; modality; domain; register; age of utterer; dialect of utterer; socioeconomic status of utterer. This is an intentional strong departure from editions of English Gigaword that focused on newswire; criterion (1) of the corpus, representativeness – following Biber (1993), requires

the inclusion of sources beyond newswire text. A set of the currently-in-process and complete sections is detailed in Table 1.

### 4.1. Data and metadata unification

Each section is contained in one directory, named after the "prefix" for the section. This prefix is pre-pended to all files in that section, with each file representing a single UTF encoded document (for a definition of document appropriate for that section). Within each section there is: (a) a LICENSE file describing precisely how that section is licensed; (b) a JSONL[6] file describing metadata about each document, including the document's ID (which is also its filename), preferably date information and an origin URI for re-retrieving the document or a relevant API or metadata result, and other optional fields with pre-defined names. For multi-speaker corpus sections, an optional "talere.json" can be included in the section, containing one JSON dictionary that is keyed by speaker ID. This assumes speaker IDs are used consistently through all documents in that section.

Sections are managed individually, as part of a larger repository of the whole Danish Gigaword corpus. A validation script helps make sure that the sections are uniformly represented.

### 4.2. Data protection

The corpus cannot contain "sensitive" data as per the GDPR definition; that means no information identifying sexual orientation, political beliefs, religion, health, etc. Thus, data discussing potentially personally sensitive topics, for example, social media around political discussions, must be stripped of personally identifying information. Further, social media content is supplied as IDs and code for rehydration, avoiding redistribution of this content.

### 4.3. Test/Train partitions

Following the result that fixing test/train splits leads to unreliable results (Gorman and Bedrick, 2019), explicitly no test/train partitions should be set in Danish Gigaword; users are encouraged to randomly select splits.

### 4.4. Licensing

To reach criterion (2), all parts of the corpus must be licensed openly, for free distribution. An example license is something like Creative Commons general license (CC0) or CC-BY. Some parts may be included under tighter licenses, such as CC-NC, which forbids commercial re-use, but the general goal of the dataset is to further all research on NLP for Danish, and so this kind of license is not preferred.

Some older corpora (e.g. Kromann et al. (2003)) used the right under Danish copyright law to cite small excerpts – up to 250 words – from published articles. This is a creative solution to sharing digital language data. For Danish Gigaword, we prefer whole articles, as they are easier to work with, providing the full context.

## 5. Corpus sections

Here we detail some of the sections included in the corpus, specifying what they bring to the dataset to make it a rich

---

| Text source | Date | Modality | Domain | Dialect | Socioeconomic | Tokens |
|---|---|---|---|---|---|---|
| Folketinget | 2009-2019 | Spoken | Parliament speech | Rigsdansk | high | 52M |
| Reddit | 2008+ | Written | Social media | mixed | mixed | 73M |
| DDT | 1983-92 | Written | Newswire | Rigsdansk | medium | 0.1M |
| Retsinformation | | Written | Legal | legal | high | 188M |
| OpenSubtitles | 1980+ | Spoken | Video subtitles | mixed | mixed | 131M |
| Spontaneous speech | 2019 | Spoken | Informal | mixed | mixed | 0.7M |
| Religious texts | | Written | Religious | Rigsdansk | unknown | 0.6M |
| DanAvis 20 | 1999-2003 | Written | Newswire | Rigsdansk | medium | 20M |
| TV2 Regionerne | 2010-2019 | Written | Newswire | Rigsdansk | medium | 10M |
| Wikipedia | 2003-2019 | Written | Encyclopaedic | Standard | mixed | 52M |
| Europarl | 1996-2011 | Spoken | Parliament speech | standard | mixed | 48M |
| Paracrawl | | Written | Web data | mixed | mixed | 103M |
| Twitter | 2019 | Written | Social media | mixed | mixed | 0.26M |
| Common Crawl | Now | Written | Web data | mixed | mixed | 81.5M |
| Bornholmsk | 1900s | Written | Transcribed books | Bornholmsk | mixed | 0.4M |
| **Total** | | | | | | 565M |

Table 1: Text dimensions by text source in the Danish Gigaword corpus.

resource that is able to cover a wide range of lexical, syntactic, and sociolinguistic phenomena expressed by Danish users.

The project considers a genre as a language style recognised by (or used to define) a community, such as news articles, personal letters, or online chat; a domain as a particular topical focus (or set of foci) that are discussed, such as biomedicine, politics, or gaming; and a medium as the means by which communication is conducted, such as writing, online chat, conversation, and so on. There is a natural overlap here between medium and speech situation though delineating this is beyond the scope of this work.

While the goal of DAGW is to cover a range of genres, domains, and media, it is very difficult to measure the prevalence of each of these across all Danish users, let alone then gather and redistribute this data. Therefore, the goal is simply to cover something of everything that can be feasibly included, without letting any particularly monolithic combination dominate (in contrast to e.g. the 100% written newswire content of English Gigaword v1). Not every intersection between genres, domains, and media can be covered, nor represented proportionally, in the first version of this project.

### 5.1. TV2 Regionerne

This corpus section comprises a contemporary sample of Danish newswire. Approximately 50 000 full newswire articles published between 2010 and 2019 are included. This source comprises articles of regional interest, written following editorial standards. The value that this section brings to the corpus is in both its temporal variation, covering a decade of events, and also in its spatial variation, covering many local events across most of the country (TV2 Bornholm is excluded). This means that many local named entities will be represented which might otherwise be missed in a dataset of national news.

### 5.2. Folketinget

The Danish parliament (Folketinget) keeps a record of all meetings in the parliament hall.[7] All records have a

transcript which was produced using a Danish version of SpeechMagic ASR from Nuance that was adapted by Dictus. The ASR system was GMM-based until recently when a neural network-based ASR system named Dictus Sun replaced SpeechMagic.[8]

All transcripts have been post-edited by linguists employed by Folketinget for intelligibility, i.e. dysfluencies, restarts, repairs, and mistakes have been edited out. The transcript is therefore not a representation of the spoken language but rather information content. The transcripts are made available in temporary revisions and continuously updated without notice, but from manual inspection the transcripts are of good quality.

In the parliament hall, the speaker is addressing members of the parliament. Monologues may include rebuttals or other comments to statements in previous monologues. While speakers can read aloud from a prepared statement or speak extemporaneously, we expect no difference to be apparent in the data because of the post-editing.

### 5.3. Retsinformation

The site `https://www.retsinformation.dk` provides access to Danish laws and regulations as well as documents from the Danish parliament (Folketinget). The text is provided by Folketinget, ministries, the ombudsman of the Folketinget and Rigsrevisionen. The legislative texts in this section include a variety of features: Uppercase text, redaction where names and addresses are left out, itemized text with chapter and section numbering, headlines, words with intra-letter spacing.

### 5.4. Spontaneous speech

The conversational corpus included originates from interdisciplinary research conducted within the Interacting Minds Center,[9] and the Puzzle of Danish project[10] at Aarhus University. Transcribed Danish speech is generally a rare kind of data, and spontaneous speech especially so;

---

[7] There are no records of committee meetings or *samråd*.

these manually transcribed conversations thus form a valuable resource to be able to distribute, especially given the careful construction of this data. Spontaneous and pseudo-spontaneous conversations are elicited in a variety of contexts: getting to know each other, solving a puzzle together, making joint decisions, etc. The participants have agreed on releasing anonymized transcripts of their conversations. All conversations involve two speakers sometimes conversing face-to-face, sometimes via a chat tool. Speech is transcribed post-hoc by native speakers. Studies published relying on this data include Fusaroli et al. (2012), Dideriksen et al. (2019), and Tylén et al. (2016).

### 5.5. Danish Wikipedia

This section comprises a dump of Danish Wikipedia[11], stripped of Wikipedia-specific markup. The content is collaboratively written by a broad range of authors, and covers many specific articles that often do not exist in other languages. This makes for a broad range of styles, most of which have been at least roughly checked for syntactic and orthographic canonicity by editors of the Danish Wikipedia, and is a rich source of region-specific named entities, often situated in full, fluent sentences. The content is reproduced verbatim in accordance with the GNU Free Documentation License.

### 5.6. Europarl

The Europarl Parallel Corpus (Koehn, 2005) contains proceedings of the European Parliament in 21 European languages that were automatically extracted and aligned. We include the Danish part of the Europarl corpus and perform no preprocessing other than file format conversions.

### 5.7. OpenSubtitles

OpenSubtitles[12] is a website where a community of people write and share subtitles for mostly big-budget movies. We extract the Danish subtitles from the OpenSubtitles section of OPUS (Lison and Tiedemann, 2016). The corpus is cleaned to fix incorrect use of characters such as capital letter I instead of lower case letter L. Files not containing any characters specific to Danish (i.e., any of the letters $å$, $æ$, or $ø$) are removed.

### 5.8. Religious text

A Danish translation of the Bible from the Massively Parallel Bible corpus (Christodouloupoulos and Steedman, 2015) is included. No pre-processing was performed other than file format conversion.

### 5.9. Danish Twitter

Social media content is rich in unedited text, allowing for a very broad range of expressions. We know that social media users typically vary their language use to afford some representation for what would typically be communicated non-verbally, and while there are corpora for this for e.g. English, there are very few published corpora containing Danish social media text (e.g. (Hovy et al., 2015; Lillie et al., 2019)). This section contains approximately 29 000 tweets in Danish from the #dkpol hashtag collected during the national parliamentary elections of 2019 (Derczynski et al., 2019) as dehydrated content, and a script for rebuilding this part of the corpus, thus permitting GDPR-compliant redistribution.

### 5.10. DanAvis20

Corpus DanAvis20 (20M words) consists of articles from various national Danish (daily) newspapers including Aktuelt, Berlingske Tidende, Dagen, and Weekendavisen. The articles were published during 1999-2003. All texts included have been cleared for distribution under the CC0 license (cf. Section 4.4.). As part of the clearing agreement the papers were slightly edited by limiting all text quotes to 200 words (at most), picking sentences from longer papers at random. Sentences were mildly scrambled (DanAvis20 has no instances left of 4 adjacent sentences). Proper names were pseudonymized (except 'Denmark', 'København', 'USA' and a few others). Infrequent content words (10ppm or less) were replaced in situ by 'statistical cognates', i.e. words of similar frequency and equivalent morpho-syntactic form (e.g. replacing "Der er sardiner i køleskabet." with "Der er skilsmissesager i forsikringsselskabet." while keeping "Ministeren rejser hjem igen"). As overall statistical and lexical properties of DanAvis20 are thus kept invariant, the corpus still provides good material for most NLP training purposes.

### 5.11. The *Bornholmsk Ordbog* Dictionary Project

Fictional texts of various kinds written in Bornholmsk, the dialect spoken on the Danish island of Bornholm,[13] have been digitized (OCR'ed and proofread) by volunteers working within the recently resumed *Bornholmsk Ordbog* dictionary project (Kjeldsen, 2019). Most of the material included is written by Otto J. Lund in the period 1930-48 (novels, short stories and poems), but the Bornholmsk sub-corpus, which in its present state amounts to circa 400K words, also includes folk stories published by J. P. Kuhre in 1938, and by K. M. Kofoed in 1935, fictional letters by various authors published in the 1930s, as well as poems by Alfred Jensen published in 1948 and various other texts from the same period. The non-standardised orthography varies considerably from source to source. The Bornholmsk part of the Danish Gigaword is a significantly extended dataset, well beyond that studied in earlier NLP work on the dialect (Derczynski and Kjeldsen, 2019).

## 6. Project status

### 6.1. Corpus distribution

As mentioned earlier in this paper and other places (Kirchmeier et al., 2019; Kirkedal et al., 2019), one problem that Danish NLP has suffered from is a lack of large accessible corpora in Danish. To address this, as well as maintaining strict licensing standards that permit open and free re-distribution, Danish Gigaword is hosted on

---

[11]https://dumps.wikimedia.org/dawiki/
[12]https://www.opensubtitles.org

[13]The language code for Bornholmsk under IETF BCP-47 is da-bornholm.

GitHub and will be uploaded to major dataset distribution services (e.g. Figshare) at each significant release.

## 6.2. Goals, inclusion, and data diversity

The project aims to provide a living dataset for training machine learning models to process Danish. This means including a broad range of quality data that fits the licensing criteria. So, although not every intersection of each dimension that can be used to describe corpora can be included, any dataset contributions that fit are more than welcome, and the project group is actively seeking these out. For example, the social media Twitter section is currently limited to politics, a domain already well-represented in DAGW; thus, more short social media on other topics is a natural extension.

## 6.3. Project sustainability

DAGW is an intrinsically community-owned, open project. The current group of participants covers academia, industry, and public sector, in a bid to improve and uphold its relevance at a broad level. However, the project is also volunteer-led and volunteer-driven, which brings intrinsic risk. Aside from cross-sector involvement, this project attempts to mitigate that risk through policies in licensing, distribution, membership, community, and data integrity.

The data is licensed CC-BY. This gives it broad reach and applicability, and makes it easier for stakeholders to join than a copyleft or noncommerical license, such as GPL or CC-NC, would. It also improves distribution prospects: because of this licensing choice, DAGW can be hosted at a third-party research data repository like Zenodo or Figshare, shifting the responsibility for effective hosting, housing and provision of the data to third parties specailised in doing this. The project also maintains an open policy, with any qualified stakeholder welcome to join, especially if there is a compatible donation of data. Denmark's size helps keep a manageable community. The Danish Gigaword also fosters community involvement by publishing results – for example, this paper. Finally, a small toolkit is included in the project's Github repository for validating any committed data, ensuring integrity, quality and uniformity of DAGW data and metadata automatically.

## 6.4. Future work

Danish Gigaword is a very active project and we hope to complete the first release in 2020. More sources are constantly being surveyed to add balance to the corpus, including fiction, older works from the 1800s, web fora, and contemporary newswire. After reaching the first billion, the project will continue, providing regular iterations and updates, with data to be released under Creative Commons licensing and freely distributed.

## 7. Conclusion

This paper reports the goal and progress of the Danish Gigaword project, a unified voluntary effort across many institutions and by many Danish speakers to construct a billion-word corpus that is representative of the language and is intended to be useful to a maximally broad and diverse group of users.

In Denmark, natural language processing is nascent and growing faster and faster. Content restrictions and conservative licensing abound. We hope that this concrete and significant contribution serves not only to benefit anyone doing natural language processing or other linguistic activities with Danish, but also as an encouraging example to others to make data available; for research, and also to bring the technological benefits enjoyed by the anglophone world to Danish speakers.

## 8. Bibliographical References

Biber, D. (1993). Representativeness in corpus design. *Literary and linguistic computing*, 8(4):243–257.

Christodouloupoulos, C. and Steedman, M. (2015). A massively parallel corpus: the bible in 100 languages. *Language Resources and Evaluation*, 49(2):375–395, Jun.

Derczynski, L. S. and Kjeldsen, A. S. (2019). Bornholmsk Natural Language Processing: Resources and Tools. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics (NoDaLiDa), September 30-October 2, Turku, Finland*, pages 338–344.

Derczynski, L., Maynard, D., Aswani, N., and Bontcheva, K. (2013). Microblog-genre noise and impact on semantic annotation accuracy. In *Proceedings of the 24th ACM Conference on Hypertext and Social Media*, pages 21–30. ACM.

Derczynski, L., Bontcheva, K., and Roberts, I. (2016). Broad Twitter Corpus: A diverse named entity recognition resource. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1169–1179.

Derczynski, L., Albert-Lindqvist, T. O., Bendsen, M. V., Inie, N., Pedersen, J. E., and Pedersen, V. D. (2019). Misinformation on twitter during the danish national election: A case study. In *Proceedings of the conference for Truth and Trust Online*.

Dideriksen, C., Fusaroli, R., Tylén, K., Dingemanse, M., and Christiansen, M. H. (2019). Contextualizing conversational strategies: Backchannel, repair and linguistic alignment in spontaneous and task-oriented conversations. In *Proceedings of the 41st Annual Conference of the Cognitive Science Society*, pages 261–267. Cognitive Science Society.

Fusaroli, R., Bahrami, B., Olsen, K., Roepstorff, A., Rees, G., Frith, C., and Tylén, K. (2012). Coming to terms: quantifying the benefits of linguistic coordination. *Psychological science*, 23(8):931–939.

Goldhahn, D., Eckart, T., and Quasthoff, U. (2012). Building Large Monolingual Dictionaries at the Leipzig Corpora Collection: From 100 to 200 Languages. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation*, pages 759–765, May.

Gorman, K. and Bedrick, S. (2019). We need to talk about standard splits. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2786–2791, Florence, Italy, July. Association for Computational Linguistics.

Graff, D., Kong, J., Chen, K., and Maeda, K. (2003). En-

glish Gigaword. *Linguistic Data Consortium, Philadelphia*, 4(1):34.

Haas, R. (2019). Discriminating between similar nordic languages using machine learning.

Hovy, D., Johannsen, A., and Søgaard, A. (2015). User review sites as a resource for large-scale sociolinguistic studies. In *Proceedings of the 24th international conference on World Wide Web*, pages 452–461. International World Wide Web Conferences Steering Committee.

Joulin, A., Bojanowski, P., Mikolov, T., Jégou, H., and Grave, E. (2018). Loss in Translation: Learning Bilingual Word Mapping with a Retrieval Criterion. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.

Kirchmeier, S., Henrichsen, P. J., Diderichsen, P., and Hansen, N. B. (2019). *Dansk sprogteknologi i verdensklasse*. The Danish Language Council.

Kirkedal, A., Plank, B., Derczynski, L., and Schluter, N. (2019). The Lacunae of Danish Natural Language Processing. In *Proceedings of the 22nd Nordic Conference on Computional Linguistics (NoDaLiDa)*, pages 356–362.

Kjeldsen, A. S. (2019). Bornholmsk Ordbog, version 2.0. *Mål og Mæle*, 40. årgang:22–31.

Koehn, P. (2005). Europarl: A parallel corpus for statistical machine translation. In *MT Summit*, volume 5, pages 79–86.

Kromann, M. T., Mikkelsen, L., and Lynge, S. K. (2003). Danish Dependency Treebank. In *Proc. TLT*, pages 217–220.

Lillie, A. E., Middelboe, E. R., and Derczynski, L. (2019). Joint rumour stance and veracity prediction. In *Proceedings of the 22nd Nordic Conference on Computional Linguistics (NoDaLiDa), September 30-October 2, Turku, Finland*, pages 208–221.

Lison, P. and Tiedemann, J. (2016). Opensubtitles2016: Extracting Large Parallel Corpora from Movie and TV Subtitles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, pages 923–929.

Nielsen, F. Å. and Hansen, L. K. (2017). Open semantic analysis: The case of word level semantics in Danish. In *Human Language Technologies as a Challenge for Computer Science and Linguistics*, pages 415–419, October.

Pedersen, B., Wedekind, J., Kirchmeier-Andersen, S., Nimb, S., Rasmussen, J.-E., Larsen, L., Bøhm-Andersen, S., Henriksen, P., Kjærum, J., Revsbech, P., Thomsen, H., Hoffensetz-Andresen, S., and Maegaard, B. (2012). *Det danske sprog i den digitale tidsalder*. Springer.

Steingrímsson, S., Helgadóttir, S., Rögnvaldsson, E., Barkarson, S., and Guðnason, J. (2018). Risamálheild: A very large Icelandic text corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.

Tylén, K., Fusaroli, R., Smith, P., and Arnoldi, J. (2016). The social route to abstraction. *Cognitive Science*.