

Topic Model Robustness to Automatic Speech Recognition Errors in Podcast Transcripts

RALUCA ALEXANDRA FETIC*, Podimo, Denmark

MIKKEL JORDAHN*, DTU Compute, Technical University of Denmark, Denmark

LUCAS CHAVES LIMA*, Podimo, Denmark

RASMUS ARPE FOGH EGEBAEK, DTU Compute, Technical University of Denmark, Denmark

MARTIN CARSTEN NIELSEN, DTU Compute, Technical University of Denmark, Denmark

BENJAMIN BIERING, Podimo, Denmark

LARS KAI HANSEN, DTU Compute, Technical University of Denmark, Denmark

For a multilingual podcast streaming service, it is critical to be able to deliver relevant content to all users independent of language. Podcast content relevance is conventionally determined using various metadata sources. However, with the increasing quality of speech recognition in many languages, utilizing automatic transcriptions to provide better content recommendations becomes possible. In this work, we explore the robustness of a Latent Dirichlet Allocation topic model when applied to transcripts created by an automatic speech recognition engine. Specifically, we explore how increasing transcription noise influences topics obtained from transcriptions in Danish; a low resource language. First, we observe a baseline of cosine similarity scores between topic vectors from automatic transcriptions and the descriptions of the podcasts written by the podcast creators. We then observe how the cosine similarities between topic vectors decrease as transcription noise increases and conclude that even when automatic speech recognition transcripts are erroneous, it is still possible to obtain high-quality topic vectors.

CCS Concepts: • **Information systems** → **Information retrieval**; • **Computing methodologies** → **Natural language processing**.

Additional Key Words and Phrases: Podcasts, Automatic Speech Recognition, Topic modeling, Recommendation Systems

ACM Reference Format:

Raluca Alexandra Fetic, Mikkel Jordahn, Lucas Chaves Lima, Rasmus Arpe Fogh Egebaek, Martin Carsten Nielsen, Benjamin Biering, and Lars Kai Hansen. 2021. Topic Model Robustness to Automatic Speech Recognition Errors in Podcast Transcripts. In *RecSys '21: ACM Conference Series on Recommender Systems, September 27–October 05, 2021, Amsterdam, NL*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Podcasts have become an increasingly popular audio format in recent years. Podcasts encompass a variety of on-demand audio, such as radio, news, and entertainment in the form of informal discussions, interviews, or even narrated content similar to audiobooks, found in several different categories. Despite the growth in popularity of podcasts, an open challenge is how to find a new podcast to listen to. Research in Podcast recommendation has not yet been able to follow

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

Manuscript submitted to ACM

up and provide users with efficient and high-quality recommendations which are key to ensuring a high quality of streaming services [18]. Podcast recommendation is a challenging task considering the very large amount of podcast episodes that lack metadata on both podcast and episode levels. Previous work has shown that transcription-based topic modeling plays a crucial role in podcast recommendation, as users tend to focus on the topics of the podcasts instead of podcast audio style [22, 35].

Topic modeling techniques such as Latent Dirichlet Allocation (LDA) [6, 7] and Probabilistic Latent Semantic Indexing (PLSI) [17] are widely used to discover topics over high-quality texts (e.g. news, blog posts, etc.). As podcasts are usually not scripted, a transcript generated from an Automatic Speech Recognition (ASR) system is necessary to perform topic modeling. Podcasts represent a particularly challenging audio format for speech recognition systems as a large number of artifacts are commonly present. Some of the challenges are that multiple speakers and speech overlap, background music and jingles, audio effects, and "real-life" recording conditions (e.g. background noise), which makes the ASR-generated transcripts from podcasts prone to errors.

While ASR systems have been widely explored and developed for high-resource languages such as English, there is a significant lack in many low-resource languages¹. As such the systems available in such setups can be expected to produce more errors than their high-resource counterparts. This motivates research for downstream tasks in the low-resource setup to ensure multilingual application. In this paper, we study the robustness of an LDA topic model used on podcast transcriptions generated by an ASR engine in Danish; a low resource language. We utilize podcast episodes with author descriptions and assume that these descriptions contain a good representation of topics. Hence, a high similarity between topic vectors from author descriptions and automated transcripts also indicates a good representation of topics in the automated transcript. To evaluate the robustness of the topic model, we first construct a baseline by computing a cosine similarity between the topic vector representations of the author descriptions and automated transcripts. We then introduce noise to the automated transcripts, controlled by a variable noise injection rate, determining how often words should be replaced, and observe how the cosine similarity changes. We experiment with two types of noise; simulated ASR noise sampled from a conditional distribution derived from a transcription error dataset (see Section 3.2), and words sampled uniformly from the topic model vocabulary for reference. Empirically, over a dataset of 587 episodes from 24 Danish podcasts, we find that the topic model is much more robust to simulated ASR noise than it is to noise from a uniform distribution. We present evidence that the LDA topic model is robust and captures an informative representation of topics, even in the face of imperfect transcriptions.

The remainder of this paper is structured as follows. In Section 2, we present the necessary background on topic modeling and robustness of downstream NLP systems against ASR errors. Section 3 defines the problem and details the methods used by each component to test the podcast topic modeling robustness. Section 4 presents the experimental setup. In Section 5, we present and discuss the results. Lastly, in Section 6, we conclude the paper and propose future research directions.

2 BACKGROUND

2.1 Automatic Speech Recognition

Converting speech in audio to raw text is done using an Automatic Speech Recognition (ASR) system. Speech recognition systems have drastically improved during recent years with advancements such as various data augmentation techniques [28], pre-training procedures on unlabelled speech data [2] and noisy student training [29]. State of the art performance

¹By low resource language, we refer to a smaller language with a limited amount of training data for NLP tasks such as speech recognition.

on the common English benchmark dataset LibriSpeech [27] is as low as 1.4% and 3.3% of Word Error Rate (WER) on the test-clean and the test-other partitions, respectively [36]. Labeled speech recognition training data is highly accessible in English [19, 27] but less so in many other languages. Even with crowd-sourcing data initiatives such as CommonVoice from Mozilla [1], the recognition performance gap between low and high resource languages remains quite high.² Pre-training speech recognition models with cross-lingual data has helped bridge the gap significantly [8]. However, transcripts from ASR systems for low resource languages are likely to be error-prone, especially for complex audio data such as podcasts.

2.2 Topic Modeling and Evaluation

Topic models are used to explore and structure a large set of documents according to latent semantic content. To improve downstream tasks such as search and recommendation of podcasts, a promising method is to utilize a topic model to extract the relevant topics of the podcast and hence enhance the podcast representations [22, 35]. Topic modeling has been extensively studied, and various approaches exist. For instance, Latent Semantic Indexing (LSI) [10] uses Singular Value Decomposition (SVD) or Non-negative Matrix Factorization (NMF) [20] on a term by document matrix to construct a latent space representation, which can be queried for comparison of documents. An extension of LSI, Probabilistic Latent Semantic Indexing (PLSI), models topics as distributions over words, and documents as a probabilistic mixture of those topics [17]. A similar, and very popular, approach is known as Latent Dirichlet Allocation (LDA). LDA differs from PLSI by utilizing prior Dirichlet distributions to model the topic-word distributions making it more robust to unseen data [6, 15]. Numerous extensions of the LDA model have been studied. One such extension, the Correlated Topic Model (CTM) [5], explores the correlations among topics generated by the LDA model. Other extensions include Collective LDA [31], combining multiple corpora during the training of the models, and approaches that explore the influence of the age of the documents in the topics [23, 34].

The quality of a topic model can be evaluated in different ways. A common practice is to evaluate a trained model in terms of perplexity [33] or topic coherence. Topic coherence, as opposed to a perplexity, is more similar to how humans judge the quality of topics [25]. Examples of topic coherence measures include UCI-coherence [24], U_{mass} coherence [21], and coherence based on word embeddings [12].

2.3 Robustness to Noise

Downstream NLP tasks on ASR transcripts need to be robust to noise due to the commonality of transcription errors. Robustness to noise is often evaluated by constructing a baseline result with clean text and then injecting varying degrees of noise to the text and examining how the result changes [4, 32]. To investigate the effects of noise, it is necessary to select different ways of injecting plausible ASR noise into transcripts. A recent study explored the feasibility of improving the robustness of speech-enabled systems with three methods of noise [9]; rule-based substitution which randomly substitutes a candidate word with a phonetically similar one, statistic-based confusion substitution which samples replacement words from a pre-constructed ASR confusion matrix and finally, model-based substitution utilizing a generative GPT model to directly produce ASR-like text. Another study investigated the stability of topics over noisy sources, by testing for topic model agreements [14], after subjecting the training data to insertion of frequent words, deletion and rule-based phonetic substitution errors [32].

²<https://paperswithcode.com/dataset/common-voice>

The robustness of a downstream task varies a lot depending on the specific task and the type of noise. For instance, topic modeling has previously been shown to be robust to deletion of random words, whereas the insertion of new words and phonetic substitution errors has a larger negative impact on topic stability [32]. Another relevant downstream task, neural machine translation with character-based models, has shown to struggle even with small perturbations to the input data [4].

3 METHODS

3.1 Transcript Generation

We produce podcast transcripts by parsing podcast audio through a danish transcription system developed at the Technical University of Denmark (DTU) as part of the Danspeech project.³ The system is based on the wav2vec 2.0 framework [3]. The model was pre-trained using approximately 945 hours of podcast episodes and 400 hours of audiobooks, and fine-tuned using the Connectionist Temporal Classification [13] (CTC) loss function on 200 hours of labeled data from the Nordisk Språkteknologi (NST) danish training dataset⁴ and 267 hours of aligned audiobook data. The Fairseq library [26] was used for both the pre-training procedure as well as fine-tuning. During inference, the ASR engine performs prefix beam search [16] with an open-source danish 3-gram language model⁵ when decoding the probabilities emitted from the wav2vec 2.0 model.

3.2 Noise Injection

We inject noise into the ASR transcripts by means of a noise injection rate parameter, β , that determines the frequency at which we substitute words in the transcript. More specifically, for each word in a given transcript we independently decide if a substitution should take place with probability β . Examples of substitutions at various levels of β are presented in Table 1.

Table 1. Examples of how a transcription changes with automatic speech recognition statistics-based substitutions at varying levels of the noise injection rate parameter β .

β	Transcription	Description
0	historien er rig på spændende fortællinger om drama krig voldelig politiske omvæltninger og fyldt med mystik hemmeligheder og fascinerende menneskeskæbner	Historien er rig på spændende fortællinger om drama, krig, voldelige politiske omvæltninger og er fyldt med mystik, hemmeligheder og fascinerende menneskeskæbner. {...}
0.3	kane er rig på så fortællinger om drama krig — - politiske omvæltninger har fyldt ved mystik er og fascineren menneskeskæbner	Episode 4: Røde agenter Revolutionen i Rusland i 1917 blev startskuddet til en international politisk kamp for at udbrede socialismen til hele verden { ... } I programmet medvirker historikerne Niels Erik Rosenfeldt og Morten Møller.
1.0	historie har rik p så fortælling er — strama til — — er film som — er var fascineren —	

When a word is selected for substitution we apply one of two methods for the substitution. The first method samples a word uniformly at random from the topic model vocabulary, and the second method uses a statistics-based confusion matrix approach (see Section 4.2). When performing statistics-based confusion substitutions, we sample replacement

³<https://danspeech.github.io/danspeech/html/index.html>

⁴<https://www.nb.no/sprakbanken/ressurskatalog/oai-nb-no-sbr-55/>

⁵https://danspeech.github.io/danspeech/html/lms.html#danspeech.language_models.DSLWikiLeipzig3gram

words \bar{w} from a conditional error distribution with the candidate probability given as

$$P(\bar{w}|w) = \frac{W(\bar{w})}{\sum_{\bar{w} \in V(w)} W(\bar{w})}. \quad (1)$$

Here $W(\bar{w})$ denotes the weighting of a candidate word and $V(w)$ denotes the candidate set for word w .

3.3 Topic Modeling and Document Vector Representation Similarity

To produce topic representations we use a general-purpose danish topic model trained on an external, multi-domain text corpus using LDA. Our motivation for doing so is two-fold; most prominently we choose to rely on an external text corpus because the topic modeling process requires large amounts of textual data, which is something that is not readily available to us in a pure podcast domain setup. Secondly, we choose LDA as the modeling framework due to its inherent probabilistic approach of representing topics, which have been shown to be robust to unseen data.

The LDA model can be seen as a probabilistic function $\gamma(d) \mapsto \mathbf{d}_T$, which maps a given document d to a topic vector $\mathbf{d}_T = [p_{t_1}, \dots, p_{t_{|T|}}]$, where p_t represents the probability of each topic t for the document d . Thus to compute the similarity of topics present in a pair of documents, d_1 and d_2 , we parse each document through the topic model and compute the Similarity between the resulting document-level topic vectors as follows,

$$\text{Similarity}(d_1, d_2) = \cos(\gamma(d_1), \gamma(d_2)) \quad (2)$$

where \cos is the cosine similarity. When comparing two ordered sets of document pairs, (S_1, S_2) , we denote the the average similarity between the sets as the CorpusSimilarity (CS) computed as,

$$\text{CS}(S_1, S_2) = \frac{1}{|S_1|} \sum_i \text{Similarity}(S_1(i), S_2(i)) \quad (3)$$

where $S_k(i)$ is document d_i in set S_k .

4 EXPERIMENTS

4.1 Podcast Dataset

The podcast dataset we use to investigate the robustness of topic modeling consists of 587 episodes from 24 podcasts shows in Danish. The 24 podcasts shows belong to 8 categories, assigned by the content creators, such as “Culture & leisure”, “Health & personal development”, “History & religion” and “True crime & mysteries”. The podcast shows are all single speaker with a limited amount of audio artifacts such as background music and jingles. We extend the descriptions with the episode title, podcast title, podcast description, and podcast category. We only include episodes with a high-quality author description, defining high-quality descriptions as any description-transcription pairings that have an initial cosine similarity above 0.5. The number of episodes for each podcast included in the dataset is presented in Figure 2 in the Appendix, in which it can be seen that the distribution of episodes across the shows is imbalanced, with some podcasts containing the majority of the episodes. We produce ASR transcripts for all the episodes by leveraging the ASR engine described in Section 3.1.

4.2 Statistics Based Automatic Speech Recognition Noise

To create the word-level conditional ASR error distributions used for statistics-based confusion substitution (see Section 3.2), we first construct a dataset containing clean text and ASR transcription pairs. We use the same ASR engine as

described in Section 3.1 with the exception that the wav2vec 2.0 model was only fine-tuned on the training dataset from NST. The data consists of approximately 77 hours of data from the NST test dataset and 267 hours of audiobook data resulting in 229,499 pairs of reference-transcript pairs. The probability of a candidate word for a given word is then obtained by counting the frequency at which the word is wrongfully transcribed as the candidate word, normalized by the number of candidates as shown in Equation (1). Examples of how candidate distributions may look for specific words are presented in Table 2 and Table 3, where the *random* candidate is a grouping of many potential errors with very low probability. We can see from the tables that typical ASR errors tend to retain semantic meaning to some degree which gives rise to the hypothesis that topic models are likely to be robust to ASR noise. If an unknown word is encountered then it is simply deleted. Across the transcripts of the podcasts, unknown words occur 20% of the time.

Table 2. Candidate distribution for the word "nogensinde".

Word	Error candidates	Probability
Nogensinde	Nogen sinde	0.917
	Sinde	0.053
	Nogen	0.024
	Nogen sider	0.003
	Står	0.003

Table 3. Candidate distribution for the word "lavet".

Word	Error candidates	Probability
Lavet	Lave	0.409
	Lade	0.136
	Lavede	0.136
	Ladet	0.091
	<i>Random</i>	0.227

4.3 Topic model

We train an LDA topic model on a Danish Wikipedia dataset consisting of 264,505 documents from The Danish Gigaword Corpus [11] using the Gensim framework [30]. All documents are preprocessed by first performing word tokenization and removing punctuation and other special characters, including numbers. Next, we apply Part of speech (POS) filtering, keeping adjectives, nouns and verbs. Finally, we perform lemmatization, lowercase all letters, before finally vectorizing the documents into bag of word (BOW) representations. The BOW vocabulary may contain n-grams and is limited by removing uncommon n-grams that appear in less than ten documents and very common n-grams that occur in more than 90% of the documents. We fix the LDA parameters α and η as $\alpha = \frac{1}{|T|}$ and $\eta = 0.1$, and choose our remaining hyper-parameters by performing grid search optimizing for topic coherence, using the U_{mass} -coherence metric to choose the best model. We tune the following hyper-parameters: Topics [10, 20, 30, 40, 50, 60, 70, 80, 90, 100], Vocabulary of BOW [Unigrams, Unigrams + Bigrams] and Variational Bayes iterations [5, 10, 15, 20]. The underlined values are the values that yielded the best model in terms of U_{mass} coherence score.

4.4 Evaluation of Topic Robustness over Noisy Sources

For the experiments, we construct a baseline by computing the CS between topic vectors of two document sets (S_1, S_2) as described in Section 3.3. We then measure the robustness to noise as the average change in magnitude of cosine similarity scores after injecting noise to documents in S_2 at varying values of β as described in Section 3.2. We alter the noise substitution method between experiments to allow for a comparison between uniform and statistics based noise.

We conduct two complementary experiments:

- (1) Testing for similarity between podcast descriptions (S_1) and noisy ASR transcripts (S_2) at varying levels of noise. This allows us to identify if the topic model is robust to transcription errors, and whether ASR transcripts contain enough information for the topic model to produce meaningful topic vectors.

(2) Testing for similarity between raw ASR transcripts (S_1) and noisy ASR transcripts (S_2) at varying levels of noise.

This further investigates how robust the topic model is to transcription errors, under the assumption that the transcript provides meaningful topic vectors.

For each experiment, we vary β in $[0, 1]$ with steps of 0.1. To reduce variance, we repeat this procedure 50 times and report the average and the standard error.

5 RESULTS AND DISCUSSION

The results of the two series of experiments are presented in Figure 1a and Figure 1b, respectively.

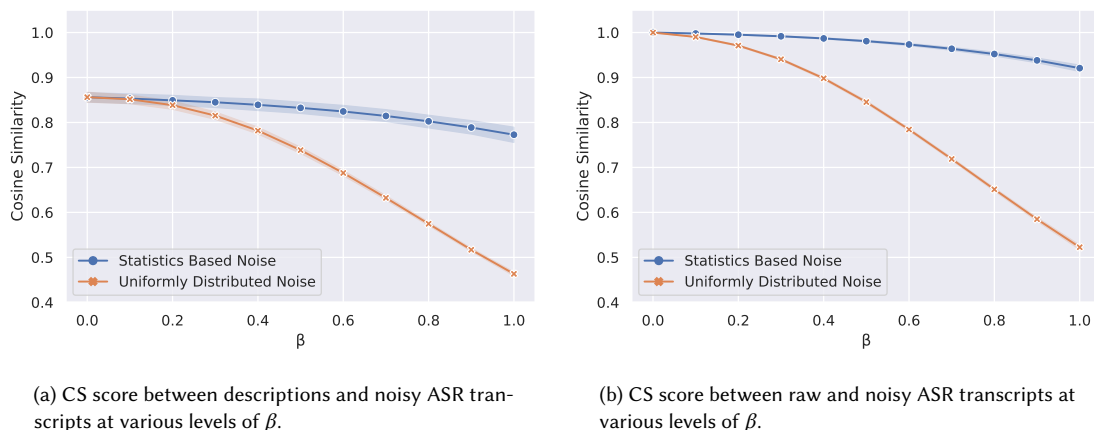


Fig. 1. Analysis of topic model robustness using cosine similarity, shaded area showing standard error of the mean.

Table 4. The number of different podcast shows that are present in the cosine similarity deciles at $\beta = 0$.

Decile	1	2	3	4	5	6	7	8	9	10
No. Unique podcasts	8	12	14	15	14	11	13	14	11	5

As can be seen from Figure 1, injecting noise in ASR transcripts has an effect on the CS for both noise distributions. We observe 45.9% (0.85-0.46) and 49% (1.00-0.51) relative decreases in CS when injecting uniform noise, while limited impact of 8.8% (0.85-0.775) and 9% (1.00-0.91) is found for statistics based ASR noise. This is evidence that topic vectors are significantly more robust to statistics-based ASR noise than uniformly distributed noise. We also observe that the mean CS score between the descriptions and transcripts at $\beta = 0$ is 0.85 which is significantly higher than what we observe as the lower bound, which can be roughly estimated to be 0.45 (found at $\beta = 1$ after injecting uniform noise), which is evidence that ASR transcripts can be used to produce meaningful topic vectors. We suspect that the reason for the small change in topic distributions under the statistics based noise is mainly due to the type of errors produced by the ASR engine which, as seen in Table 2 and Table 3, have a tendency to retain semantic meaning. Note that β can be seen as an estimation of the WER in a podcast transcript (see Figure 3 in Appendix A), which suggests that even when an ASR engine produces a transcript with a high word error rate, the transcript will still be viable for topic modeling because the errors will have little impact on the topic distribution after LDA preprocessing.

When comparing transcripts and author descriptions for the baseline with $\beta = 0$, the standard deviation has a value of 0.128. Furthermore, the difference between the minimum value of 0.500 and the max value of 0.997 is large. We investigate why this occurs by splitting the podcast episodes into deciles based on their cosine similarity and counting the number of unique podcast in each decile as seen in Table 4. The lowest amount of unique podcast shows are in the top and bottom deciles. Furthermore, in the 1st decile, 42 out of the 59 episodes comes from the same podcast "Sagen om Amagermanden" (note, it is also the podcast show that is the most represented) and in the 10th decile, 52 of the 58 episodes are from podcast show "Hvor er mit ansigt?". This result suggests that either the quality of the podcast descriptions or the performance of the ASR engine is very dependent on the specific podcast show. However, since all podcast shows in the dataset are single speaker with limited amount background music, we suspect that quality of the author descriptions is the primary reason why episodes originating from the same podcast show consistently perform poorly in some cases.

6 CONCLUSION

In this work, we conducted experiments to evaluate the quality and robustness of topic representations produced by a general LDA topic model when exposed to noisy ASR transcripts in a low-resource language (Danish). More specifically, we investigated how injecting two different noise profiles to raw ASR transcripts influence the topic distributions across a podcast dataset at varying levels of noise. We created the dataset by leveraging an ASR system to obtain podcast transcripts. We choose to only include podcast episodes that had a high-quality author-written description as part of the meta-data in the dataset, to allow for comparison between the description and transcript topic vectors, relying on the assumption that a well-written description is very similar in terms of topic distribution to that of the episode content. To obtain topic vectors we trained a general-purpose LDA model on Danish Wikipedia data. We experimented with injecting two types of noise into the raw ASR transcripts, namely uniformly random noise and simulated ASR noise. We obtained similarity baselines by computing vector similarities of the raw transcripts with their respective descriptions and the raw transcripts with themselves at various levels of noise injection. We found that injecting random noise to the transcripts significantly influenced the CS score (45.9% and 49% relative to baseline) whereas injecting simulated ASR noise only slightly influenced the CS score (8.8% and 9% relative to baseline). Given our findings, we conclude that even when an ASR engine produces podcast transcripts at higher WERs, we can still obtain meaningful topic representations. We hypothesize that this is because even when an ASR engine produces erroneous text, the majority of the word-level errors will carry similar semantic meaning to the underlying truth. This encourages the use of ASR engines and transcriptions for podcast recommendations. Even in cases where the ASR engine is challenged by complex audio scenery, the transcripts could be sufficient to obtain topic representations that are relevant when representing the semantic content of a podcast episode. Furthermore, we also found implications that using ASR transcripts would be generally more robust than relying on human tinkered descriptions which are largely reliant on the individual content creators when it comes to performing topic modelling on podcast content. We conclude this based on the fact that episodes originating from the same podcast shows either consistently showed very high or very low topic overlap between their transcripts and their descriptions.

The ASR noise investigated in this paper was based on ASR errors from clean recordings without common podcast artifacts such as the overlap of speech and ambient sounds. In the future, it would be worth investigating whether topic models also are robust to ASR noise stemming from these types of errors. Furthermore, there are indications that other downstream tasks might not be as robust to noise as topic modeling, hence analyzing the robustness of other relevant downstream tasks to ASR noise remains an open research question.

7 ACKNOWLEDGMENTS

We gratefully acknowledge support from Innovation Fund Denmark in the forms of their Innobooster and Innoexplorer grants (Grant numbers 0173-00670B and 0160-00023 respectively).

REFERENCES

- [1] Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. 2019. Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670* (2019).
- [2] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *arXiv preprint arXiv:2006.11477* (2020).
- [3] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations. *CoRR abs/2006.11477* (2020). arXiv:2006.11477 <https://arxiv.org/abs/2006.11477>
- [4] Yonatan Belinkov and Yonatan Bisk. 2017. Synthetic and natural noise both break neural machine translation. *arXiv preprint arXiv:1711.02173* (2017).
- [5] D. Blei and J. Lafferty. 2006. Correlated Topic Models. *ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS* 18 (2006), 147. <http://www.cs.cmu.edu/~lafferty/pub/ctm.pdf>
- [6] David Blei, Andrew Ng, and Michael Jordan. 2002. Latent Dirichlet Allocation. In *Advances in Neural Information Processing Systems*, T. Dietterich, S. Becker, and Z. Ghahramani (Eds.), Vol. 14. MIT Press, Vancouver, British Columbia, Canada. <https://proceedings.neurips.cc/paper/2001/file/296472c9542ad4d4788d543508116cbc-Paper.pdf>
- [7] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet Allocation. *J. Mach. Learn. Res.* 3, null (March 2003), 993–1022.
- [8] Alexis Conneau, Alexei Baevski, Roman Collobert, Abdelrahman Mohamed, and Michael Auli. 2020. Unsupervised cross-lingual representation learning for speech recognition. *arXiv preprint arXiv:2006.13979* (2020).
- [9] Tong Cui, Jinghui Xiao, Liangyou Li, Xin Jiang, and Qun Liu. 2021. An Approach to Improve Robustness of NLP Systems against ASR Errors. (March 2021). arXiv:2103.13610 [cs.CL]
- [10] Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. 1990. Indexing by latent semantic analysis. *Journal of the American Society for Information Science* 41, 6 (1990), 391–407. [https://doi.org/10.1002/\(SICI\)1097-4571\(199009\)41:6<391::AID-ASII>3.0.CO;2-9](https://doi.org/10.1002/(SICI)1097-4571(199009)41:6<391::AID-ASII>3.0.CO;2-9) arXiv:<https://asistdl.onlinelibrary.wiley.com/doi/pdf/10.1002/%28SICI%291097-4571%28199009%2941%3A6%3C391%3A%3AAID-ASI%3E3.0.CO%3B2-9>
- [11] Leon Derczynski, Manuel R. Ciosici, Rebekah Baglini, Morten H. Christiansen, Jacob Aarup Dalsgaard, Riccardo Fusaroli, Peter Juel Henriksen, Rasmus Hvingelby, Andreas Kirkedal, Alex Speed Kjeldsen, Claus Ladefoged, Finn Årup Nielsen, Jens Madsen, Malte Lau Petersen, Jonathan Hvithamar Rystrom, and Daniel Varab. 2021. The Danish Gigaword Corpus. In *Proceedings of the 23rd Nordic Conference on Computational Linguistics*. NEALT.
- [12] Anjie Fang, Craig Macdonald, Iadh Ounis, and Philip Habel. 2016. Using Word Embedding to Evaluate the Coherence of Topics from Twitter Data. In *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval (Pisa, Italy) (SIGIR '16)*. Association for Computing Machinery, New York, NY, USA, 1057–1060. <https://doi.org/10.1145/2911451.2914729>
- [13] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*. 369–376.
- [14] Derek Greene, Derek O’Callaghan, and Pádraig Cunningham. 2014. How Many Topics? Stability Analysis for Topic Models. *CoRR abs/1404.4606* (2014). arXiv:1404.4606 <http://arxiv.org/abs/1404.4606>
- [15] Thomas Griffiths and Mark Steyvers. 2003. Prediction and Semantic Association. In *Advances in Neural Information Processing Systems*, S. Becker, S. Thrun, and K. Obermayer (Eds.), Vol. 15. MIT Press, Vancouver, British Columbia, Canada. <https://proceedings.neurips.cc/paper/2002/file/cb8da6767461f2812ae4290eac7cbc42-Paper.pdf>
- [16] Awni Y Hannun, Andrew L. Maas, Daniel Jurafsky, and Andrew Y Ng. 2014. First-pass large vocabulary continuous speech recognition using bi-directional recurrent dnns. *arXiv preprint arXiv:1408.2873* (2014).
- [17] Thomas Hofmann. 1999. Probabilistic Latent Semantic Analysis. In *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence* (Stockholm, Sweden) (UAI’99). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 289–296.
- [18] Rosie Jones, Hamed Zamani, Markus Schedl, Ching-Wei Chen, Sravana Reddy, Ann Clifton, Jussi Karlgren, Helia Hashemi, Aasish Pappu, Zahra Nazari, Longqi Yang, Oguz Semerci, Hugues Bouchard, and Ben Carterette. 2021. Current Challenges and Future Directions in Podcast Information Access. arXiv:2106.09227 [cs.IR]
- [19] Jacob Kahn, Morgane Rivière, Weiyi Zheng, Evgeny Kharitonov, Qiantong Xu, Pierre-Emmanuel Mazaré, Julien Karadayi, Vitaliy Liptchinsky, Ronan Collobert, Christian Fuegen, et al. 2020. Libri-light: A benchmark for asr with limited or no supervision. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 7669–7673.
- [20] Daniel Lee and H. Seung. 1999. Learning the Parts of Objects by Non-Negative Matrix Factorization. *Nature* 401 (11 1999), 788–91. <https://doi.org/10.1038/44565>
- [21] David Mimno, Hanna Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. 2011. Optimizing Semantic Coherence in Topic Models. *EMNLP 2011 - Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*, 262–272.

- [22] Lasse Lohilahti Molgaard, Kasper Winther Jorgensen, and Lars Kai Hansen. 2007. Castsearch-context based spoken document retrieval. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, Vol. 4. IEEE, IV–93.
- [23] Ramesh M. Nallapati, Susan Dittmore, John D. Lafferty, and Kin Ung. 2007. Multiscale Topic Tomography. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (San Jose, California, USA) (KDD '07)*. Association for Computing Machinery, New York, NY, USA, 520–529. <https://doi.org/10.1145/1281192.1281249>
- [24] David Newman, Jey Han Lau, Karl Grieser, and Timothy Baldwin. 2010. Automatic Evaluation of Topic Coherence. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics (Los Angeles, California) (HLT '10)*. Association for Computational Linguistics, USA, 100–108.
- [25] David Newman, Youn Noh, Edmund Talley, Sarvnaz Karimi, and Timothy Baldwin. 2010. Evaluating Topic Models for Digital Libraries. In *Proceedings of the 10th Annual Joint Conference on Digital Libraries (Gold Coast, Queensland, Australia) (JCDL '10)*. Association for Computing Machinery, New York, NY, USA, 215–224. <https://doi.org/10.1145/1816123.1816156>
- [26] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A Fast, Extensible Toolkit for Sequence Modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.
- [27] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 5206–5210.
- [28] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le. 2019. SpecAugment: A simple data augmentation method for automatic speech recognition. *arXiv preprint arXiv:1904.08779* (2019).
- [29] Daniel S Park, Yu Zhang, Ye Jia, Wei Han, Chung-Cheng Chiu, Bo Li, Yonghui Wu, and Quoc V Le. 2020. Improved noisy student training for automatic speech recognition. *arXiv preprint arXiv:2005.09629* (2020).
- [30] Radim Řehůřek and Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. ELRA, Valletta, Malta, 45–50. <http://is.muni.cz/publication/884893/en>.
- [31] Zhi-Yong Shen, Jun Sun, and Yi-Dong Shen. 2008. Collective Latent Dirichlet Allocation. In *2008 Eighth IEEE International Conference on Data Mining*, 1019–1024. <https://doi.org/10.1109/ICDM.2008.75>
- [32] Jing Su, Derek Greene, and Oisín Boydell. 2016. Topic Stability over Noisy Sources. In *Proceedings of the 2nd Workshop on Noisy User-generated Text (WNUT)*. The COLING 2016 Organizing Committee, Osaka, Japan, 85–93. <https://aclanthology.org/W16-3913>
- [33] Hanna M Wallach, Iain Murray, Ruslan Salakhutdinov, and David Mimno. 2009. Evaluation methods for topic models. In *Proceedings of the 26th annual international conference on machine learning*. 1105–1112.
- [34] Xuerui Wang and Andrew McCallum. 2006. Topics over Time: A Non-Markov Continuous-Time Model of Topical Trends. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Philadelphia, PA, USA) (KDD '06)*. Association for Computing Machinery, New York, NY, USA, 424–433. <https://doi.org/10.1145/1150402.1150450>
- [35] Longqi Yang, Yu Wang, Drew Dunne, Michael Sobolev, Mor Naaman, and Deborah Estrin. 2019. More Than Just Words: Modeling Non-Textual Characteristics of Podcasts. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining (Melbourne VIC, Australia) (WSDM '19)*. Association for Computing Machinery, New York, NY, USA, 276–284. <https://doi.org/10.1145/3289600.3290993>
- [36] Yu Zhang, James Qin, Daniel S Park, Wei Han, Chung-Cheng Chiu, Ruoming Pang, Quoc V Le, and Yonghui Wu. 2020. Pushing the limits of semi-supervised learning for automatic speech recognition. *arXiv preprint arXiv:2010.10504* (2020).

A ADDITIONAL FIGURES

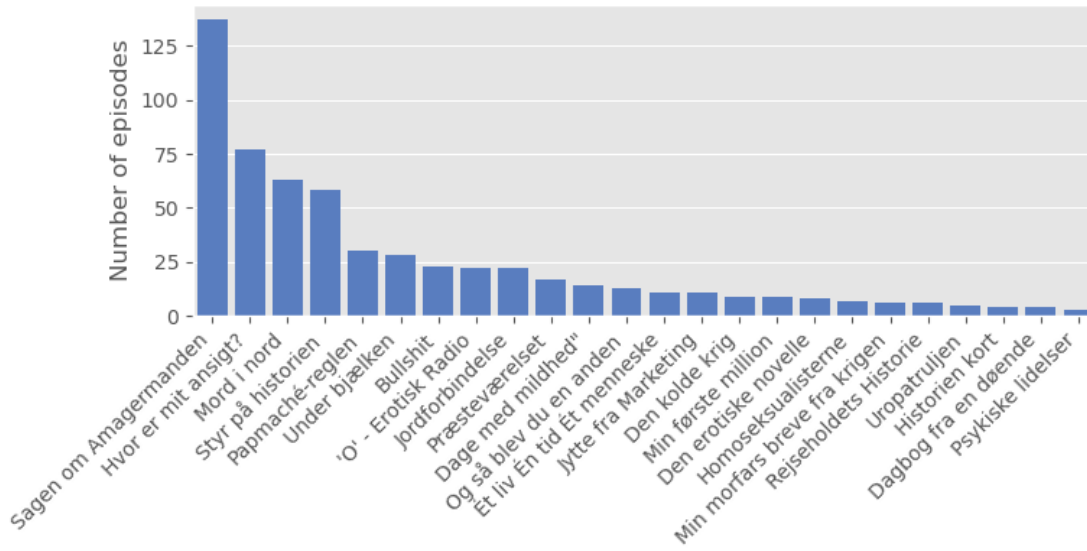


Fig. 2. Distribution of episodes per podcast.

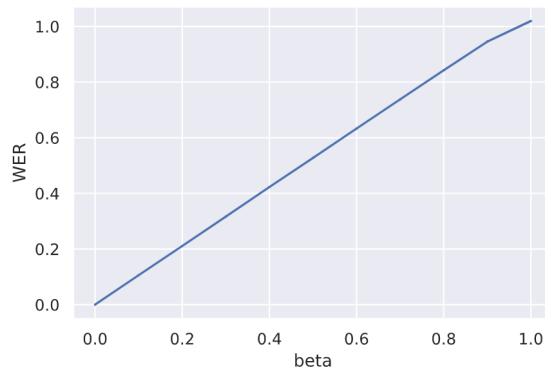


Fig. 3. Word Error Rate as function of β