

AUGMENTING AUTOMATIC SPEECH RECOGNITION MODELS WITH DISFLUENCY DETECTION

Robin Amann, Zhaolin Li, Barbara Bruno, Jan Niehues

Karlsruhe Institute of Technology, Germany

ABSTRACT

Speech disfluency commonly occurs in conversational and spontaneous speech. However, standard Automatic Speech Recognition (ASR) models struggle to accurately recognize these disfluencies because they are typically trained on fluent transcripts. Current research mainly focuses on detecting disfluencies within transcripts, overlooking their exact location and duration in the speech. Additionally, previous work often requires model fine-tuning and addresses limited types of disfluencies.

In this work, we present an inference-only approach to augment any ASR model with the ability to detect open-set disfluencies. We first demonstrate that ASR models have difficulty transcribing speech disfluencies. Next, this work proposes a modified Connectionist Temporal Classification(CTC)-based forced alignment algorithm from [1] to predict word-level timestamps while effectively capturing disfluent speech. Additionally, we develop a model to classify alignment gaps between timestamps as either containing disfluent speech or silence. This model achieves an accuracy of 81.62% and an F1-score of 80.07%. We test the augmentation pipeline of alignment gap detection and classification on a disfluent dataset. Our results show that we captured 74.13% of the words that were initially missed by the transcription, demonstrating the potential of this pipeline for downstream tasks.

Index Terms— Automatic speech recognition, speech disfluency, forced alignment

1. INTRODUCTION

Speech disfluency refers to interruptions in the flow of speech, such as repetitions, interjections, and revisions. It is a natural part of conversational and spontaneous speech but can be particularly pronounced and frequent in certain speech disorders, such as stuttering [2, 3, 4]. Analyzing speech disfluency can aid in diagnosing speech disorders. It can also help in understanding language proficiency that can be applied, for example, in interviews and children’s education.

Manually annotating speech disfluency for analysis is costly, and Automatic Speech Recognition (ASR) can support the annotation process. The ASR systems transcribe the speech into readable text, which can then be passed to

the evaluator or automatic evaluation pipelines for analysis. However, ASR models show performance degradation in disfluent speech, because the models are developed to generate fluent transcripts to enhance readability [2, 5].

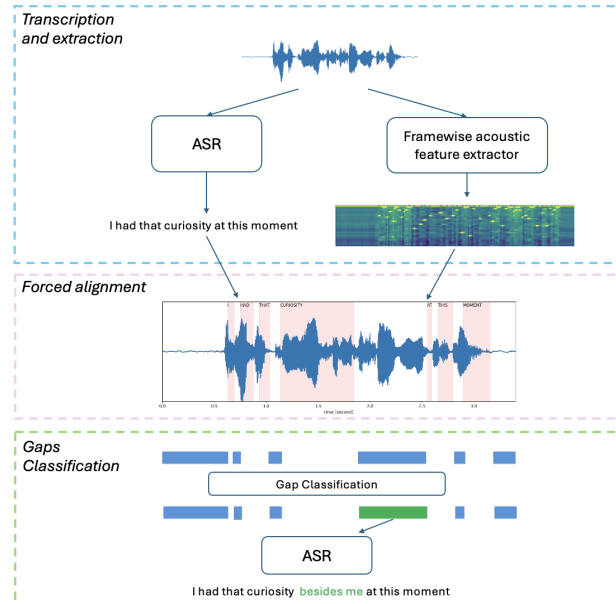


Fig. 1. The pipeline to augment ASR models with disfluency detection with a follow-up re-transcription for example of application.

To detect speech disfluency with ASR models, one popular approach is post-processing the ASR predictions as a sequence labelling problem [6, 7, 8, 9, 10]. Alternatively, [11, 12, 13, 14, 15, 16] focus on jointly predicting transcriptions and disfluencies with end-to-end speech recognition. In addition, [17] explores adapting the ASR foundation models, which are robust to recognize unfinished words, to detect disfluencies. However, those approaches only detect disfluency within the transcript while neglecting the location and duration of the speech disfluency, which plays an important role in the disfluency analysis, e.g. for the assessment of interlocutors’ alignment in collaborative activities [18].

Recent work focuses on detecting speech disfluencies at

the frame level to capture timing information. [19, 20] investigate fine-tuning ASR models with disfluent dataset. [4] explores forced alignment that aligns the audio signal with its corresponding transcript by decoding with Weighted Finite State Transducers (WFSTs) in alignment graph. [14, 21] hierarchically integrate transcription and detection modules. However, these works address predefined or restricted disfluency types and fail on the open-set disfluency detection for handling previously unseen types.

In this work, we propose a straightforward yet effective pipeline to augment ASR models by detecting open-set speech disfluency¹. As illustrated in Figure 1, the pipeline consists of three hierarchical steps: transcription and feature extraction with the ASR model and a frame-wise feature extractor, transcript and speech alignment with a modified Connectionist Temporal Classification(CTC)-based approach from [1], and alignment gaps classification for detecting the potential disfluencies. The contributions of this work are as follows:

- We examine one state-of-the-art ASR model Whisper [22] on speech disfluency detection. The experimental results show the model achieves 22.54 Word Error Rate (WER) points on a conversational dataset, but only 56% of speech disfluencies at the word level are correctly transcribed, and 73.77% of untranscribed words are disfluencies. The results indicate that the ASR model performs poorly on disfluent speech, which is reasonable since ASR models are designed to produce fluent transcriptions for better readability. The finding highlights the importance of augmenting ASR models with disfluency detection capabilities.
- Aiming to detect the location and duration of speech disfluency, this work proposes a modified CTC-based forced alignment approach from [1] to effectively locate and capture speech disfluency. We compare the proposed approach with the popular CTC-based alignment [1] and Whisper’s cross-attention alignment [22], and show that the proposed algorithm clearly captures more untranscribed words than the others.
- With the proposed forced alignment approach, we build an inference-only pipeline to augment ASR models with disfluency detection capability. The pipeline is flexible and can be adapted to any ASR model. In addition, the pipeline detects the alignment gaps containing disfluent speech, allowing the detection of the open-set disfluency beyond predefined types. With an alignment gap classification model, the pipeline achieves 81.62% accuracy in identifying gaps containing speech, covering 74.13% of all untranscribed words in the initial transcript.

¹<https://github.com/Robin-Amann/bachelor-thesis>

2. DISFLUENCY DETECTION

2.1. Augmentation pipeline

The inference-only pipeline to augment ASR models with open-set disfluency detection consists of three hierarchical steps (Figure 1). In the first step, the ASR system generates an estimated transcript, and a feature extractor model produces the frame-wise probability from speech. The ASR model could be the same as the frame-wise feature extractor models, such as Wav2Vec2 [23]. But it can be any ASR model as the augmentation pipeline only needs its transcript. After that, the pipeline applies a modified (CTC)-based forced alignment algorithm, that is based on [1] with the above generations. The algorithm generates word-level timing information and the signal gaps between the word timesteps are recognized as potential instances of disfluent speech. In the end, a developed classification model is applied to identify alignment gaps containing disfluent speech or only silence. The classification results can be utilized for downstream tasks like second-step transcription and identifying disfluency types.

2.2. Forced alignment

As the key step to extract timing information, three forced alignment approaches are employed: the standard CTC forced alignment, a modified CTC forced alignment, and the cross-attention approach of Whisper, one of the SOTA ASR models, whose attention value exhibits a high degree of correlation with timestamps.

2.2.1. Standard CTC Alignment

CTC-based forced alignment is a popular approach to extracting timing information used in many speech recognition packages, such as ESPnet [24], SpeechBrain [25] and Flashlight [26]. The alignment is calculated in three steps²:

1. The audio is fed into a feature extraction model that is pre-trained with CTC to generate frame-wise label probability over the whole alphabet.
2. From the probability, a trellis matrix is generated to represent the probability of labels occurring at each time frame. The trellis at point (t, j) (where $0 \leq t \leq T - 1$ and $0 \leq j \leq U - 1$) represents the maximal probability that the first $j - 1$ labels of the transcripts are aligned to the first $t - 1$ timeframes of the audio. The trellis is calculated in the log domain to avoid numerical instability.

The maximum probability that the first j labels of the transcript are aligned at the timeframe t is the maximum of 1) Stay on the label, which means the first j labels of the transcript are already aligned at time

²https://pytorch.org/audio/stable/tutorials/forced_alignment_tutorial.html

$t - 1$, and the alignment remains with the same label; 2) Switch to next label, which means that the first $j - 1$ labels of the transcript are aligned at time $t - 1$, and the alignment switches to label j at time t . The calculation is as follows:

$$trellis[j, t] = \max(trellis[j, t - 1] \cdot prob[t, \epsilon], trellis[j - 1, t - 1] \cdot prob[t, j]) \quad (1)$$

where $trellis[j, t]$ and $prob[j, t]$ indicate the trellis value and the frame-wise probability value at time t and label j , respectively.

3. In the third step, the path with the highest probability in the trellis is traced back. This path begins at position $(0, 0)$ and ends at position (T, U) to encompass the entire audio and transcript.

2.2.2. Modified CTC Alignment

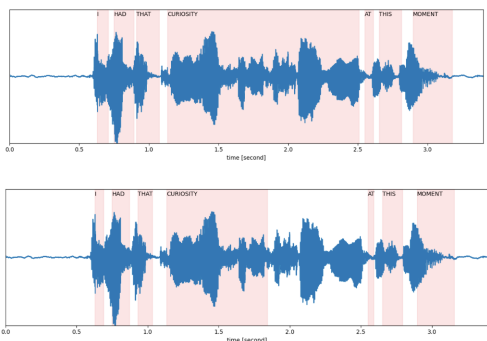


Fig. 2. Alignments of generated transcription to speech signal with standard (upper) and modified (lower) CTC forced alignments. The ASR prediction is: *I had that curiosity at the moment*, while the manual transcript contains the disfluency: *I had that curiosity beside me at the moment*.

Our preliminary experiments show that the standard CTC alignment struggles to generate correct information when the automated transcription does not include the disfluency. As Figure 2 shows, the manual transcript of this example contains the disfluency, which is removed by the ASR model for better readability. The standard CTC alignment extends the alignment around incomplete words, leading to inaccurate alignment and missing disfluency detection. This occurs because the standard CTC algorithm tends to align a word for a longer duration rather than to align silence where something is being said. In trellis generation (refer to Equation 1), the emissions for a blank token in this part of the audio are very low, as something was spoken there. As a consequence, gaps in the alignment of the transcript may not occur where they should, which is undesirable for this application.

To counteract this issue, we proposed the modified CTC alignment to enable the model to detect the alignment gaps containing the speech of untranscribed disfluency. In trellis generation, the modified Equation 2 is applied **if the current label j is a space token**. The space token infers a special label used to represent gaps or spaces between characters or tokens in the output space. With modification, the modified probability of staying on this label is the maximum of the probability acquired through the emissions and a predefined probability c . This modification incentivises the algorithm to extend silence between words to some extent.

$$trellis[j, t] = \max(trellis[j, t - 1] \cdot \max(prob[t, \epsilon], c), trellis[j - 1, t - 1] \cdot prob[t, j]) \quad (2)$$

As for the previous examples shown in Figure 2, the modified CTC alignment correctly aligns the speech signals to the ASR prediction while capturing gaps that correspond to untranscribed disfluencies. However, it is also important to note that the modified alignment for, e.g., “that” has also become shorter. This is because this modification encourages the algorithm to keep words short since an alignment containing much silence gives a better score. Setting the probability c too high may result in words being too short overall. Therefore, we experiment with different predefined probabilities to choose the value carefully.

2.2.3. Cross-attention Alignment

Whisper is trained in such a way that there exists a correlation between the cross-attention weights and the audio input. Consequently, the cross-attention weights and the output transcript allow for the calculation of alignment to the input audio through Dynamic Time Warping (DTW) [27]. As implemented in the HuggingFace library³, token-level timesteps are calculated using the encoder-decoder cross-attentions and DTW to map each output token to a position in the input audio.

2.3. Alignment Comparison Metric

Comparing the alignment approaches requires an automatic metric, which is not available. This work proposes a metric that evaluates the alignment between the manual and automatic transcripts by considering the position and length of aligned words. The correctly transcribed words are extracted using Levenshtein Alignment, and we denote $(s1, e1)$ and $(s2, e2)$ as the manual timing annotation and automatically aligned timing information.

After forced alignment, timing information of the same word from manual annotation $A1$ and automated transcription $A2$. For each pair of words $(w1, w2)$, the length, position and

³https://github.com/huggingface/transformers/blob/main/src/transformers/models/whisper/generation_whisper.py

combined scores are calculated with Equations 3, 4 and 5, respectively:

$$m_p(w_1, w_2) = \frac{1}{\left| \frac{p_1 - p_2}{l_1} \right| + 1} \quad (3)$$

$$m_l(w_1, w_2) = \frac{1}{\left| \frac{l_1 - l_2}{l_1} \right| + 1} \quad (4)$$

$$m(w_1, w_2) = \frac{1}{\left| \frac{p_1 - p_2}{l_1} \right| + 1} \cdot \frac{1}{\left| \frac{l_1 - l_2}{l_1} \right| + 1} \quad (5)$$

Where w indicate a word with starting (s) and ending (e) time, and p and q indicate the position and length of the word calculated with $p = \frac{s+e}{2}$ and $l = \frac{e-s}{2}$. The final scores for the entire alignments are then computed by averaging the scores of individual words. The three scores range from 0 to 1, and a higher value indicates a better alignment.

2.4. Gap classification

The forced alignment detects the alignment gaps, while the gaps can contain speech or only silence. Accordingly, we propose a classification step to identify alignment gaps where disfluent speech may occur.

Since the timing information of the disfluent speech is not available in the dataset, we define an alignment gap containing disfluent speech as one that covers at least one word. We define the coverage as the duration of the word has more than 50% overlapping with the duration of the alignment gap. The 50% overlap criterion strikes a balance: Considering only words completely within the gap would result in many gaps being deemed empty, despite there being speech intuitively present. Conversely, if a gap is considered non-empty as soon as a word is even partially within it, the transcription model may find it challenging to transcribe this word in the subsequent step.

Classification with all gaps is inefficient and might involve too small gaps due to alignment inaccuracy. Therefore, this work performs classification on gaps that exceed a minimum threshold length. The chosen minimum gap size should not be too small, as minimal inconsistencies in the alignment would become too noticeable. This would result in previously transcribed speech being transcribed again. On the other hand, the minimal gap should not be too large, as this could lead to overlooking too many gaps where untranscribed speech may be present. Based on preliminary experimental results, a gap size of 0.3 is selected as optimal.

3. DATASET

3.1. Dataset Preparation

As this work aims to detect the location and duration of speech disfluencies, the dataset must be composed of spontaneous speech with word-level timing annotation. Besides,

the dataset must be large enough to act as training data for the classification model. Therefore, we use the Switchboard-1 Release 2⁴ and Treebank 3⁵ datasets. The Switchboard dataset consists of approximately 260 hours of telephone conversations with word-level timing information, and the Treebank 3 dataset adds the corresponding word-level disfluency annotation to the transcripts of the Switchboard dataset.

3.2. Segmenting Audio

This dataset consists of recordings that are several minutes long. The long recordings hinder forced alignment performance, and the segmentation pre-processing is applied.

First, the audio file is segmented at all points where there is more than 5 seconds of silence. These points provide good places to divide the transcript without interrupting the speaker’s flow of speech. However, there are still segments that remain several minutes long.

In the next step, as long as the segment is longer than 30 seconds, the largest gap between two words in the middle of the transcript is sought, ensuring that it is at least 10 seconds away from the beginning and end of the segment. This ensures that the resulting segments are between 10 and 30 seconds long. Of course, it is also possible that shorter segments are created when splitting by 5-second pauses.

4. EXPERIMENTS AND RESULTS

4.1. Experimental setups

This work experiments with the SOTA ASR model Whisper⁶ to augment disfluency detection ability. As for frame-wise feature extractor, this work selects Wav2Vec2⁷ which is fine-tuned on English ASR. The ASR model can be the same as the feature extractor model in pipeline design, but we select Whisper as the ASR model because this pipeline supports augmenting any ASR model, and Whispers is a stronger ASR model in terms of accuracy and robustness for this dataset. We believe our approach is effective for ASR models providing more accurate predictions than Whisper. Besides, we acknowledge that our approach may suffer performance degradation with ASR models yielding less accurate predictions.

4.2. Are ASR models good at disfluency recognition?

For the initial transcription of the audio, we employ the Whisper model on the dataset and achieve 22.54 WER points. To assess the percentage of fluent and not fluent speech transcribed, the manual and the automatic transcript are aligned using the operations obtained during the calculation of the

⁴<https://catalog ldc.upenn.edu/LDC97S62>

⁵<https://catalog ldc.upenn.edu/LDC99T42>

⁶<https://huggingface.co/openai/whisper-large-v3>

⁷<https://huggingface.co/facebook/wav2vec2-base-960h>

WER. Table 1 illustrates the number of words from the manual transcript that were correctly transcribed, incorrectly transcribed and not transcribed, each annotated as fluent or not fluent speech.

	Correctly Transcribed	Incorrectly Transcribed	Untranscribed
Fluent	895,474	45,823	31,931
Disfluent	136,718	15,242	89,799

Table 1. Fluent and disfluent transcribed words.

It can be observed that approximately 74% of all untranscribed words are labelled as speech disfluencies. Additionally, it is evident that 37% of all speech disfluencies are not transcribed, 6% are transcribed incorrectly and only 56% are transcribed correctly, confirming the initial assumption that Whisper does not fully transcribe speech disfluencies.

4.3. What parameter to use for modified CTC algorithm?

This work proposes a modified CTC alignment algorithm to overcome the problems of the standard CTC alignment algorithm on disfluency recognition (Section 2.2.2). The modified algorithm involves a fixed probability c as to incentive gap recognition in alignment. The probability c experimented with probability values -5, -4, -3, -2, -1, -0.5, -0.1, -0.01 on a par with the log probability scale in frame-wise acoustic probability. Setting as -0.01 is essentially 0 and the value above 0 makes no sense as for comparison with log probability. A higher value of probability c indicates a higher chance to stay with the space token, corresponding to more and longer alignment gaps.

This work counts the number of words that are covered by the alignment gaps to evaluate alignment performance. We evaluate the alignment approaches on all words in the manual transcription to show the general alignment performance, and we also evaluate them on only the words next to the untranscribed words to show the performance specifically on disfluent speech. Note that the untranscribed words are determined with Levenshtein Alignment, same as Section 2.3. As Figure 3 shows, significantly more untranscribed words are reachable with a default probability of -0.01 than -5.

Besides, we evaluate using the proposed alignment scores. As shown in Table 2, the combined alignment scores for all words show no significant difference, but the score for words around the untranscribed words improves clearly with decreasing the probability value of c , which is consistent to Figure 3. Therefore, the -0.001 is chosen for the modified CTC algorithm in later experiments.

4.4. Which forced alignment algorithm works better?

The alignments calculated by the standard CTC alignment, the modified CTC algorithm and cross-attention are compared

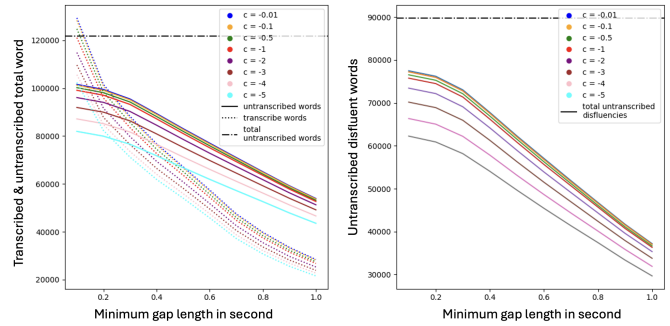


Fig. 3. The number of words that are covered by modified alignments with different probability c value. The left is for all words and the right is for the words next to the untranscribed words in the manual transcript.

Probability value	All words	Words around untranscribed words
-0.01	0.5893	0.5628
-0.1	0.5896	0.5627
-0.5	0.5906	0.5627
-1	0.5917	0.5625
-2	0.5932	0.5600
-3	0.5937	0.5550
-4	0.5937	0.5484
-5	0.5932	0.5407

Table 2. Experimental results on modified CTC alignment with different predefined probability using the combined alignment score.

with the proposed evaluation metric. Same as Section 4.3, we evaluate the performance for all words as well as only for the words around the untranscribed speech. But here we calculate the proposed alignment metrics of the position, length and the combined scores.

As Table 3 shows, the modified CTC algorithm outperforms the others in all aspects, except the length for all words where the standard CTC is slightly better. For the alignment performance in the presence of untranscribed speech, the modified CTC algorithm shows a clear performance gain.

Besides evaluating the metric scores, this work also counts how many untranscribed words are covered. Figure 4 shows for each presented alignment method the number of untranscribed words and already transcribed words within a gap for various minimum gap sizes. As can be seen, the modified CTC algorithm recognized many more untranscribed words than other algorithms. Specifically, with a total amount of 121,738, modified CTC, standard CTC and cross-attention cover 81.69%, 46.10% and 12.02% untranscribed words, respectively. Therefore, the modified CTC alignment is chosen for further analysis.

	Cross Attention	CTC	Modified CTC
All words			
Position	0.5941	0.7465	0.7702
Length	0.6855	0.7755	0.7612
Combined	0.4359	0.5880	0.5893
Words around untranscribed words			
Position	0.5138	0.7177	0.7619
Length	0.6051	0.7376	0.7555
Combined	0.3508	0.5493	0.5802

Table 3. Comparison of forced alignment algorithms with evaluation of alignment metric. *Position* and *Length* indicate the individual score, and *Combined* indicates the score considering position and length (Section 2.3).

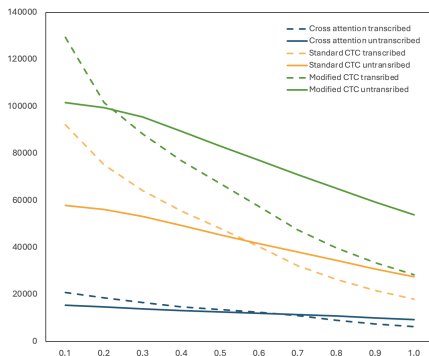


Fig. 4. Untranscribed and transcribed words covered by three forced alignment approaches.

4.5. How to build disfluency classification model?

The forced alignment brings gaps between the words of transcription, and the next question comes as to how to classify the gaps as containing speech or empty. This work proposes to build a classification model and train it with a dataset tailored for the pipeline with the gaps.

With the combined dataset, we build the datasets consisting of extracted gaps. The gap is labelled as “gap contains speech” if at least one word from the manual transcript falls into this gap. Otherwise, they were labelled as “gap is empty”. We use the modified CTC algorithm for alignment. After shuffling, we select 80% of the gaps for training data, and the rest for test data (Table 4).

	Total	Containing Speech	Empty
Training	220,344	101,207	119,137
Test	40,980	19,651	21,329

Table 4. Statistics of the gaps classification dataset

With the above dataset, this work builds a classification model by fine-tuning a wav2vec2 model in conjunction with a classification head. The classification head consists of a lin-

ear layer projecting the output of Wav2Vec2 onto the predefined classes: an empty gap and a gap with speech. With the evaluation of the test split, the classification model achieved an accuracy of 81.62%, a precision of 86.14%, a recall of 74.80%, and an F1-score of 80.07%.

4.6. How effective is the disfluency detection pipeline?

Knowing the proportion of all gaps successfully classified does not provide information about the proportion of untranscribed words successfully classified. Therefore, we count the number of transcribed and untranscribed words that are classified and covered in the gaps, classified but not covered in gaps and not classified by the gap classification model.

As Table 5 shows, 15,478 out of a total of 20,880 untranscribed words are covered by the detected alignment gaps, leading to a detection rate of 74.13%. However, 13,202 out of 168,256 already transcribed words are also labelled as predicted in the gaps, counting to a false detection rate of 8.6%. The false detection rate indicates the risk of double transcription if a follow-up re-transcription was carried on.

	Transcribed	Untranscribed
Classified and covered	13,202	15,478
Classified but uncovered	153,975	3,952
Not classified	1,079	1,450

Table 5. Pipeline performance evaluation on the type of words covered by the gaps.

5. CONCLUSION

In this work, we propose an inference-only pipeline to augment any ASR model with open-set disfluency detection. We reveal the current ASR models struggle to transcribe speech disfluency. To tackle this issue, we propose a modified CTC forced alignment algorithm to recognize the location and duration of speech disfluencies. We show the effectiveness of this approach by comparing it with popular forced alignment approaches in disfluency recognition. Additionally, we build a pipeline for disfluency detection and show that the approach captures 74.13% of the words that are not transcribed by the initial transcription.

However, the disfluency detection performance is dependent on the ASR model performance. This is because that transcribed disfluencies will not be identified as disfluencies, as they are already aligned with the transcribed words.

6. ACKNOWLEDGMENTS

This work was funded by the Baden-Württemberg Ministry of Science, Research and Art (MWK), via the state digitalisation strategy digital@bw.

7. REFERENCES

- [1] Ludwig Kürzinger, Dominik Winkelbauer, Lujun Li, Tobias Watzel, and Gerhard Rigoll, “CTC-segmentation of large corpora for german end-to-end speech recognition,” in *International Conference on Speech and Computer*. Springer, 2020, pp. 267–278.
- [2] Shaomei Wu, “The world is designed for fluent people: Benefits and challenges of videoconferencing technologies for people who stutter,” in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 2023, pp. 1–17.
- [3] Ai Leen Choo, Daphne Greenberg, Hongli Li, and Amani Talwar, “Rate of stuttering and factors associated with speech fluency characteristics in adult struggling readers,” *Journal of Learning Disabilities*, vol. 56, no. 1, pp. 7–24, 2023.
- [4] Theodoros Kouzelis, Georgios Paraskevopoulos, Athanasios Katsamanis, and Vassilis Katsouros, “Weakly-supervised forced alignment of disfluent speech using phoneme-level modeling,” in *Proc. INTERSPEECH 2023*, 2023, pp. 1563–1567.
- [5] Colin Lea, Zifang Huang, Jaya Narain, Lauren Tooley, Dianna Yee, Dung Tien Tran, Panayiotis Georgiou, Jeffrey P Bigham, and Leah Findlater, “From user perceptions to technical improvement: Enabling people who stutter to better use speech recognition,” in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, New York, NY, USA, 2023, CHI ’23, Association for Computing Machinery.
- [6] Angelica Chen, Vicky Zayats, Daniel Walker, and Dirk Padfield, “Teaching BERT to wait: Balancing accuracy and latency for streaming disfluency detection,” in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Seattle, United States, July 2022, pp. 827–838, Association for Computational Linguistics.
- [7] Johann C. Rocholl, Vicky Zayats, Daniel D. Walker, Noah B. Murad, Aaron Schneider, and Daniel J. Liebling, “Disfluency Detection with Unlabeled Data and Small BERT Models,” in *Proc. Interspeech 2021*, 2021, pp. 766–770.
- [8] Morteza Rohanian and Julian Hough, “Best of both worlds: Making high accuracy non-incremental transformer-based disfluency detection incremental,” in *Annual Meeting of the Association for Computational Linguistics*, 2021.
- [9] Motoi Omachi, Yuya Fujita, Shinji Watanabe, and Tianzi Wang, “Non-autoregressive end-to-end automatic speech recognition incorporating downstream natural language processing,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 6772–6776.
- [10] Peter Mihajlik, Yan Meng, Mate S Kadar, Julian Linke, Barbara Schuppler, and Katalin Mády, “On disfluency and non-lexical sound labeling for end-to-end automatic speech recognition,” in *Interspeech 2024*, 2024, pp. 1270–1274.
- [11] Tedd Kourkounakis, Amirhossein Hajavi, and Ali Etemad, “FLUENTNet: End-to-end detection of stuttered speech disfluencies with deep learning,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2986–2999, 2021.
- [12] Hayato Futami, Emiru Tsunoo, Kentaro Shibata, Yosuke Kashiwagi, Takao Okuda, Siddhant Arora, and Shinji Watanabe, “Streaming joint speech recognition and disfluency detection,” in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [13] Koharu Horii, Meiko Fukuda, Kengo Ohta, Ryota Nishimura, Atsunori Ogawa, and Norihide Kitaoka, “End-to-end spontaneous speech recognition using disfluency labeling,” in *Interspeech*, 2022, pp. 4108–4112.
- [14] Jiachen Lian, Carly Feng, Naasir Farooqi, Steve Li, Anshul Kashyap, Cheol Jun Cho, Peter Wu, Robbie Netzorg, Tingle Li, and Gopala Krishna Anumanchipalli, “Unconstrained dysfluency modeling for dysfluent speech transcription and detection,” in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2023, pp. 1–8.
- [15] Lavanya Venkatasubramaniam, Vishal Sunder, and Eric Fosler-Lussier, “End-to-end word-level disfluency detection and classification in children’s reading assessment,” in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [16] Dena Mujtaba, Nihar R. Mahapatra, Megan Arney, J. Scott Yaruss, Caryn Herring, and Jia Bin, “Inclusive asr for disfluent speech: Cascaded large-scale self-supervised learning with targeted fine-tuning and data augmentation,” in *Interspeech 2024*, 2024, pp. 1275–1279.
- [17] Rao Ma, Mengjie Qian, Mark Gales, and Katherine M Knill, “Adapting an ASR foundation model for spoken language assessment,” in *9th Workshop on Speech and Language Technology in Education (SLaTE)*. Aug. 2023, slate_2023, ISCA.

- [18] Utku Norman, Tanvi Dinkar, Barbara Bruno, and Chloé Clavel, “Studying alignment in a collaborative learning activity via automatic methods: The link between what we say and do,” *Dialogue & Discourse*, vol. 13, no. 2, pp. 1–48, 2022.
- [19] John Harvill, Mark Hasegawa-Johnson, and Chang D. Yoo, “Frame-Level Stutter Detection,” in *Proc. Interspeech 2022*, 2022, pp. 2843–2847.
- [20] Olabanji Shonibare, Xiaosu Tong, and Venkatesh Ravichandran, “Enhancing ASR for stuttered speech with limited data using detect and pass,” 2022.
- [21] Jiachen Lian and Gopala Anumanchipalli, “Towards hierarchical spoken language disfluency modeling,” in *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, Yvette Graham and Matthew Purver, Eds., St. Julian’s, Malta, Mar. 2024, pp. 539–551, Association for Computational Linguistics.
- [22] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever, “Robust speech recognition via large-scale weak supervision,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 28492–28518.
- [23] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in neural information processing systems*, vol. 33, pp. 12449–12460, 2020.
- [24] Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson Enrique Yalta Soplín, Jahn Heymann, Matthew Wiesner, Nanxin Chen, Adithya Renduchintala, and Tsubasa Ochiai, “ESPnet: End-to-end speech processing toolkit,” 2018.
- [25] Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, Loren Lugosch, Cem Subakan, Nauman Dawalatabad, Abdelwahab Heba, Jianyuan Zhong, Ju-Chieh Chou, Sung-Lin Yeh, Szu-Wei Fu, Chien-Feng Liao, Elena Rastorgueva, François Grondin, William Aris, Hwidong Na, Yan Gao, Renato De Mori, and Yoshua Bengio, “SpeechBrain: A general-purpose speech toolkit,” 2021.
- [26] Jacob Kahn, Vineel Pratap, Tatiana Likhomanenko, Qiantong Xu, Awni Hannun, Jeff Cai, Paden Tomasello, Ann Lee, Edouard Grave, Gilad Avidov, Benoit Steiner, Vitaliy Liptchinsky, Gabriel Synnaeve, and Ronan Collobert, “Flashlight: Enabling innovation in tools for machine learning,” 2022.
- [27] Toni Giorgino, “Computing and visualizing dynamic time warping alignments in r: The dtw package,” *Journal of Statistical Software*, vol. 31, no. 7, pp. 1–24, 2009.