

SPEECH INTELLIGIBILITY CLASSIFIERS FROM 550K DISORDERED SPEECH SAMPLES

Subhashini Venugopalan¹, Jimmy Tobin¹, Samuel J. Yang¹, Katie Seaver^{1,3}, Richard J.N. Cave¹, Pan-Pan Jiang¹, Neil Zeghidour², Rus Heywood¹, Jordan Green^{1,3}, Michael P. Brenner^{1,4}

¹Google Research, USA; ²Google Research, France;
³MGH Institute of Health Professions, USA; ⁴Harvard University, USA

ABSTRACT

We developed dysarthric speech intelligibility classifiers on 551,176 disordered speech samples contributed by a diverse set of 468 speakers, with a range of self-reported speaking disorders and rated for their overall intelligibility on a five-point scale. We trained three models following different deep learning approaches and evaluated them on \sim 94K utterances from 100 speakers. We further found the models to generalize well (without further training) on the TORGO database[1] (100% accuracy), UASpeech[2] (0.93 correlation), ALS-TDI PMP[3] (0.81 AUC) datasets as well as on a dataset of realistic unprompted speech we gathered (106 dysarthric and 76 control speakers, \sim 2300 samples).

Index Terms— intelligibility, disordered speech

1. INTRODUCTION

Atypical speech can manifest from a variety of conditions. Neurological diseases such as Amyotrophic Lateral Sclerosis (ALS), Parkinson’s Disease (PD), and Cerebral Palsy (CP), are amongst the most prevalent causes of dysarthria and speech disability. Automatic assessments of speech intelligibility can help predict how well voice-based assistive technologies might aid a person with speech disorders [1]. They can be used to detect such speech e.g. in YouTube, to allow better transcriptions from specialized Automatic Speech Recognition (ASR) systems [4], or used by researchers as an objective measure to monitor decline in speech e.g., in ALS [3]. Such classifiers can also help identify variable manifestations of impaired speech, to enable automatic collection of such data at scale to teach and improve ASR systems.

Classification of speech disorders and in particular classifying dysarthric speech and speech intelligibility have been fairly well studied for different applications [5, 6]. Many works have developed machine learning models based on handcrafted acoustic features [7, 8, 9, 10]. Among deep learning methods, convolutional neural networks (CNNs) are quite popular [10, 11, 12], as are recurrent neural networks (RNNs), specifically Long-Short-Term-Memory (LSTM) [13] models [14, 15, 16] have also been used to classify dysarthric speech. Some recent works have explored transformer [17,

18] based models for non-speech classification tasks [19, 20]. However, most prior works developed models on much smaller datasets of disordered speech, with fewer utterances and speakers, and focused on a limited set of phrases or speech disorder etiologies.

We use a large dataset of 756,147 utterances contributed by 677 speakers with a range of self-reported speech disorders as part of *Project Euphonia* [4]. The speech samples are rated for their overall intelligibility on a five-point Likert scale by speech-language pathologists (SLPs). We build classification models based on different deep learning architectures including convolutional networks with learnable audio frontends [21], representations from an LSTM-based ASR encoder model [22], and representations from the self-supervised wav2vec 2.0 CNN and transformer architecture backbone [23]. The models are trained on over 550K samples to predict either a binary (typical, atypical speech) label, or the five class labels. The models achieved an accuracy of over 86% when evaluated on a test set of nearly 94K utterances from 100 speakers. To assess the flexibility and generalizability of the models, we also evaluated (inference only) on (1) the TORGO database [1] consisting of 14 speakers; (2) the UASpeech dataset [2] with 28 speakers; and (3) the ALS-TDI PMP dataset [3] with 90 speakers; and (4) on unconstrained realistic speech gathered from videos of 76 controls and 106 dysarthric speakers covering 5 etiologies. Our models showed performance competitive with the state-of-the-art (SOTA) on all datasets. We found it to perform well on speakers with ALS, PD, CP, and Ataxia. We describe our models and evaluation, and share our findings here.

2. DISORDERED SPEECH CLASSIFICATION

Our work focuses on classifying *intelligibility* of dysarthric speech. Intelligibility measures how well speech is understood by a human listener [24]. In our dataset (Sec. 3), amongst other aspects, each speaker is scored for their overall intelligibility on a five-point scale by SLPs. In this work, we consider all utterances from all speakers with ratings and develop models to predict the speech intelligibility ratings.

Tasks. We train the models on two classification tasks based on the intelligibility ratings for each utterance. First is

the **2-class MILD+** task where we predict if the speech sample is *typical* or not (i.e., disordered) by grouping *mild*, *moderate*, *severe* and *profound* into the *atypical* class. The second is the **5-class** task of predicting the 5-point SLP ratings.

2.1. SpICE: Speech Intelligibility Classifiers on Euphonia

Our speech intelligibility classification approach is partly inspired by [22]. They use CNNs, representations from an unsupervised model (TRILL [20]), and representations from an ASR-encoder model to train classifiers on a dataset of 15K utterances focusing on a narrow set of 29 short phrases from each speaker. We also use an ASR-encoder, additionally, we train a CNN with a learnable frontend [21] and representations from wav2vec 2.0 which uses a transformer backbone. Also, our work significantly scales training, using 550K+ diverse utterances to train classifiers, and extensively evaluates generalization of the models across etiologies and datasets.

ASR system encoder representations (ASR-enc). This model is identical to that in [22]. We use an LSTM encoder that models acoustic inputs in an ASR system based on an RNN transducer (RNN-T) [25] model. The specific architecture is based on He et. al. [26] trained on long-form speech [27]. As in [22], we consider the average-pooled (over time) embeddings of the encoder as the representation of a speech sample, and train linear models using logistic regression, random forest, and linear discriminant analysis on the embeddings to predict class scores.

wav2vec 2.0 representations. To compare ASR-enc with a similarly powerful model, we train linear classifiers using the self-supervised representations from the final layer of the wav2vec 2.0 model [23] publicly available on Hugging-Face [28]. This architecture consists of a multi-layer CNN that produces latent speech representations of raw audio, and uses a transformer [17] and masked language modeling [29] to build contextualized representations. We develop classifiers on representations from the the 12th (768-d) final layer.

Fully learnable convolutional classifier (LEAF + CNN) As a baseline, we train a fully learnable convolutional classifier. Unlike the CNN classifier of [22], which takes as inputs fixed mel-filterbanks, the low-level representations of our model are provided by a LEAF [21] frontend which jointly learns filtering, pooling, compression and normalization from data. This frontend feeds into a 2D CNN, based on [30], which alternates convolutions along time ((3×1) kernel) and frequency ((1×3) kernel). This is trained using cross-entropy to predict intelligibility on either 2 or 5 classes.

3. DATASETS

Euphonia-SpICE Dataset. Our training data is a subset of the Euphonia dataset [4]. We use data from 677 speakers (756,147 utterances) who were rated by SLPs using a Quality Control (QC) phrase set of 29 short phrases for each participant. SLPs listened to the QC recordings for each speaker and

Table 1: Euphonia-SpICE: Count of speakers and utterances

Intelligibility	# speakers			# utterances		
	Train	Val.	Test	Train	Val.	Test
TYPICAL	161	41	25	149,941	24,142	10,664
MILD	161	29	37	208,843	22,532	39,007
MODERATE	83	23	19	124,984	48,814	21,214
SEVERE	54	12	15	60,692	13,868	22,397
PROFOUND	9	4	4	6,716	1,691	642
OVERALL	468	109	100	551,176	111,047	93,924

assessed, among other things, the overall intelligibility of the speaker on a five-point Likert scale. The scale was mapped to 5 classes - *typical*, *mild*, *moderate*, *severe*, and *profound* (detailed in [4]). All utterances from a speaker are labeled with the same rating. While [22] only uses the QC utterances, we use the full data ($\approx 50\times$) which we call the Euphonia-SpICE dataset. The speakers were randomly split into train, val and test set in a 70:15:15 ratio. All our models were trained on the same splits. Fig. 1A shows the distribution of etiologies and Tab. 1 presents the number of speakers and utterances in each split for each label, along with the overall count.

3.1. Datasets for evaluating generalization.

We evaluate our trained models (inference only) on multiple datasets to demonstrate flexibility and generalizability of the approach to diverse disorders and data collection setups.

UASpeech [2] is a database of dysarthric speech produced by speakers with CP. Our academic collaborator obtained access to the data and evaluated our wav2vec 2.0 model on “all words” (765 utterances per speaker) from 28 consented speakers (15 dysarthric and 13 controls). We used audio from channel 5 of the 8-microphone array of recordings.

TORGO. We use a subset of the TORGO database [1], as described in [31]. In particular, the subset of control speech sentences also available as dysarthric speech, and we use only the recordings from the microphone array. This yielded 1200 utterances across 7 dysarthric speakers and 7 matched controls with intelligibility labels in [a,b,c,d,e] (‘a’ being the most intelligible and ‘e’ being least). Additionally, we had one SLP rate each speaker on the same Likert scale as in the Euphonia dataset. The SLP listened to 10 utterances for each speaker (selected randomly with no overlap in the phrases between speakers) and was asked to rate the overall intelligibility of the speaker on the five-point scale.

ALS-TDI PMP dataset [3] was collected from over 500 people living with ALS over a 4 year period. Participants recorded voice samples¹ and self-reported ALS Functional Rating Scale (ALSFRRS-R) scores (integer in [0-4]) for 12 functions one of which is speech. We use the test split from [3] consisting of 1333 recordings from 90 participants.

SpICE-V To evaluate our models on unprompted speech in realistic settings from speakers with different disorders,

¹They repeat the phrase ‘I owe you a yo-yo today’ five times

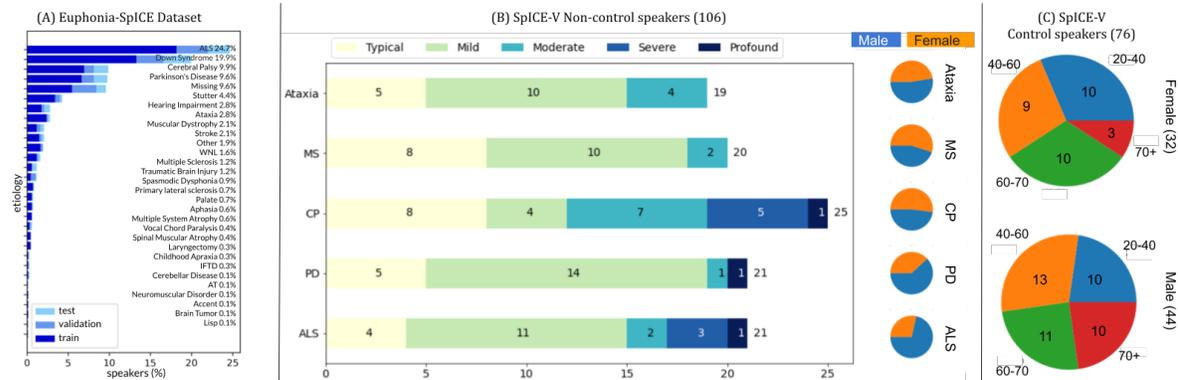


Fig. 1: [SpICE datasets] (A) Dist. of etiologies in the Euphonia-SpICE dataset. (B) SpICE-V non-control speakers split by etiology, intelligibility class and inferred gender, and (C) SpICE-V controls split by inferred age bucket and gender.

we curated our own dataset from a collection of web videos. SLPs identified ~ 20 speaker videos each for 5 etiologies: ALS, CP, PD, Ataxia, and Multiple Sclerosis (MS), accounting for balance in severity and inferred gender. They also marked time segments when the dysarthric speaker was speaking. We gathered control samples from the AudioSet [32] dataset. We watched videos labeled ‘Male speech’ and ‘Female speech’ and selected speakers to balance for inferred age² and gender. In total, we collected 106 dysarthric speaker videos containing 2221 utterances (time segments) and 76 control speaker samples (1×10 s segments each). The distribution of the data is presented in Fig. 1.

4. RESULTS

Evaluation metrics We train and evaluate our models on the Euphonia-SpICE dataset. We report utterance-level performances of the models on evaluation metrics used in [22]. Namely, **Accuracy** (Acc.), **F1 score** and **1-vs-rest AUC** (AUC) where we compute the Area Under the Receiver Operating Characteristic Curve for each class against the rest (akin to multi-label classification) and report the mean.

4.1. Euphonia-SpICE performance.

Tab. 2 presents the results of the models when trained and evaluated on the Euphonia-SpICE dataset. The ASR-enc model has the best performance on both tasks. The wav2vec 2.0 based model closely matches ASR-enc performance on the 2-class task; on the 5-class task it does slightly worse. The LEAF + CNN model which is far smaller does comparably worse, and we drop it from further evaluations.

4.2. Models generalize well on existing datasets.

TORGO. Tab. 3 presents results on the TORGO database. We compute predictions at the speaker-level, by averaging the 5-class scores of the model across all utterances and pick the argmax. We present the utterance-level accuracy (in parenthesis) on the binary classification task of whether the (5-class)

Table 2: [Euphonia-SpICE] We report the mean 1-vs-rest AUC values, F1 score, and accuracy (Acc.) at the utterance-level. Higher is better. **bold** indicates highest value.

Models	Size		2-class MILD+			5-class		
	(MB)	params	AUC	F1	Acc.	AUC	F1	Acc.
LEAF + CNN	55	8M	0.669	0.833	0.886	0.600	0.362	0.378
wav2vec 2.0	360	100M	0.742	0.857	0.863	0.652	0.416	0.423
ASR-enc	122	60M	0.761	0.861	0.862	0.714	0.422	0.432

model correctly identifies each utterance from a speaker as either *typical* or not as determined by our SLP. We observe that both the ASR-enc model and the wav2vec 2.0 based model trained on the Euphonia-SpICE dataset generalize well.

Table 3: [TORGO] Generalization (only inference) on TORGO. Per-speaker predictions and (binarized accuracy %).

Speaker #	Utts.	TORGO label	SLP label	SpICE 5-cls models		
				LEAF+CNN	wav2vec 2.0	ASR-enc
FC01	26	Control	typical	typ. (34.6)	typ. (96.2)	typ. (96.2)
FC02	122	Control	typical	typ. (68.9)	typ. (95.9)	typ. (100)
FC03	125	Control	typical	typ. (65.6)	typ. (83.2)	typ. (78.4)
MC01	118	Control	typical	typ. (55.1)	typ. (96.6)	typ. (92.4)
MC02	122	Control	typical	sev. (22.1)	typ. (94.3)	typ. (92.6)
MC03	119	Control	typical	typ. (75.6)	typ. (98.3)	typ. (98.3)
MC04	121	Control	typical	mod. (5)	typ. (98.3)	typ. (99.2)
F03	100	a	mild	typ. (63)	mild (87.0)	mild (88.0)
F04	97	a	typical	mod. (8.2)	typ. (91.8)	typ. (74.2)
M03	92	a	typical	mod. (15.2)	typ. (98.9)	typ. (100)
F01	20	d/e	moderate	mod. (85)	mod. (100)	mod. (100)
M02	92	d/e	moderate	mod. (92.4)	mild (100)	mild (100)
M04	86	d/e	severe	mod. (59.3)	sev. (100)	mod. (100)
M05	17	c	severe	typ. (41.2)	sev. (100)	mod. (100)

UASpeech contains speaker-level intelligibility ratings in the 1-100% range. We use a simple map from predicted class to intelligibility $\{0:100\%, 1: 90\%, 2: 60\%, 3: 40\%, 4: 20\%\}$ and average predictions across utterances (the mapping didn’t seem to matter as long as it was monotonic). In Tab. 4 we compare performance with prior work [33] which uses an ASR model’s error rates that requires transcription, and measures Pearson correlation between predictions and labels. Due to UASpeech access restrictions, we only evaluate on wav2vec 2.0.

²Many were public figures (e.g athletes, politicians) with wikipedia pages.

Table 4: Pearson correlation on UASpeech.

Data subset	# Speakers	SOTA [33]	wav2vec 2.0
Dysarthric; All words	15	0.94	0.91
Controls + Dysarthric; All words	28	-	0.93

Generalization to Speech ALSFRS-R prediction. Tab. 5 presents the AUC at the utterance level on predicting the Speech ALSFRS-R score (also 5 classes). We compare with [3] which uses a CNN similar to [22] trained on 3776 recordings from 389 speakers specifically to predict Speech ALSFRS-R scores, whereas we do not do any further training.

Table 5: AUC on utterance-level ALSFRS-R prediction.

# Spkr	# Utt.	SOTA [3]	ASR-enc	wav2vec 2.0
90	1333	0.86	0.82	0.81

4.3. wav2vec 2.0 generalizes well on SpICE-V

Tab. 6 presents speaker and utterance-level accuracy of the models. The performance is split based on controls (all of whom have typical speech), non-controls, and all. We separate out the group that does not include any Dysarthric speakers labeled ‘Typical’ and one which includes all Dysarthric speakers. In this more challenging dataset we can see that the self-supervised representations from wav2vec 2.0 help the model generalize better than the ASR encoder based model.

Table 6: [SpICE-V] Comparing wav2vec 2.0 and ASR-enc on speaker- and utterance-level accuracies.

Group	w. Typ. non-ctrl	Total (Atyp.) # Utts.	# Spkr	wav2vec 2.0		ASR-enc	
				spkr	Acc. (%) utt.	spkr	Acc. (%) utt.
Controls	×	76	76 (0)	76.32	76.32	96.42	96.42
Dysarthric (-Typ.)	×	1489	76 (76)	93.42	94.83	63.16	66.92
Dysarthric (all)	✓	2221	106 (76)	77.36	75.64	68.65	67.92
All (-Typ.& Dys.)	×	1565	152 (76)	84.87	93.93	78.29	68.21
All	✓	2297	182 (76)	76.92	75.66	78.57	69.47

5. DISCUSSION

Models do well on ALS, PD, CP and Ataxia. ALS and CP are the most prevalent in the evaluated datasets: TORGO, UASpeech, and ALS-TDI PMP; and our models do well on these. When we look at performance sliced by Etiology on SpICE-V (Tab. 8) and on the most prevalent 7 etiologies in Euphonia-SpICE test set (Tab. 7), we can see at the speaker level the model does well on ALS, CP, PD and Ataxia. The performance on MS is mixed, and the model has difficulty identifying speakers with MS having typical speech.

Dysarthric speakers with typical speech are harder to classify. From Tab. 6 we can observe that the models have different thresholds when predicting on dysarthric speakers with typical speech intelligibility. While the ASR-enc model identifies both controls and non-controls with typical speech as ‘Typical’ the wav2vec 2.0 model tends to identify them more often as ‘Mild’. However, when looking at Dysarthric

Table 7: [Euphonia-SpICE] Performance sliced by etiology. Both models show similar per-speaker accuracy.

Etiology	# Utts. (%)	Atyp./Total # Spkr	per-utterance AUC		Spkrs. Acc
			wav2vec 2.0	ASR-enc.	
ALS	22076 (23.7)	14 / 18	0.749	0.763	0.778
CP	14518 (15.6)	11 / 12	0.890	0.916	0.834
Down Syn.	13971 (15.0)	18 / 23	0.544	0.525	0.652
PD	13863 (14.9)	8 / 11	0.489	0.521	0.727
Hearing Imp.	8478 (9.1)	5 / 5	NA	NA	1.000
MS	6272 (6.7)	3 / 4	0.842	0.942	0.750
Musc. Dyst.	2544 (2.7)	1 / 3	0.935	0.958	0.667

Table 8: [SpICE-V] Slicing performance by etiology.

Etiology	# Utts.	# Spkr Total (Typ.)	wav2vec 2.0		ASR-enc	
			spkr	Acc. (%) utt.	spkr	Acc. (%) utt.
ALS	443	21 (4)	90.5	87.6	76.2	76.0
PD	498	21 (5)	85.7	84.9	61.9	73.0
CP	620	25 (8)	72.0	69.8	72.0	74.5
MS	352	20 (8)	55.0	57.5	60.0	48.6
Ataxia	308	19 (5)	84.2	75.6	68.4	62.1

speakers alone (Tab. 8) we see that wav2vec 2.0 performs consistently well at the speaker and utterance levels. This is explained by the significant difference in training data of the backbone models and their size.

Limitations and future work. The data in *Project Euphonia* consists of prompted English speech from participants self-identifying as having speaking disabilities. It has a male-to-female ratio of 60:37 and does not have information on race. Further there is an imbalance across many sensitive etiologies. Future work should consider typical and atypical speech that is more diverse from different demographics, minority groups, and speech in other languages and dialects. Including a fairness testing dataset would also be a valuable contribution. To use the models ‘in the wild’ it would also be necessary to include non-speech samples and unprompted speech in noisy background. The model can also be fine-tuned and calibrated on a few samples to study applicability to different etiologies.

6. CONCLUSION

In this work, we trained speech intelligibility classifiers on a large dataset of over half a million utterances from people having a range of speaking disabilities. We examined models with different backbones CNNs, LSTMs and transformers. We found our classifier to generalize well on several datasets without any additional training and does particularly well on speakers with ALS, CP, PD and Ataxia.

Acknowledgements. This study would not have been possible without the contributions and efforts of the hundreds of speakers who consented and provided their speech samples through g.co/euphonia, and members of team Euphonia for their data collection effort and feedback.

7. REFERENCES

- [1] Frank Rudzicz, Aravind Kumar Namasivayam, and Talya Wolff, “The torgo database of acoustic and articulatory speech from speakers with dysarthria,” *Language Resources and Evaluation*, 2012.
- [2] Heejin Kim, Mark Hasegawa-Johnson, Adrienne Perlman, Jon Gunderson, Thomas S Huang, Kenneth Watkin, and Simone Frame, “Dysarthric speech database for universal access research,” in *ICASSP*, 2008.
- [3] Fernando G Vieira, Subhashini Venugopalan, et al., “A machine-learning based objective measure for als disease severity,” *NPJ digital medicine*, vol. 5, no. 1, pp. 1–9, 2022.
- [4] R. L. MacDonald et al., “Disordered Speech Data Collection: Lessons Learned at 1 Million Utterances from Project Euphonia,” in *INTERSPEECH 2021*.
- [5] Andy Huang, Kyle Hall, Catherine Watson, and Seyed Reza Shahamiri, “A review of automated intelligibility assessment for dysarthric speakers,” in *2021 International Conference on Speech Technology and Human-Computer Dialogue (SpED)*.
- [6] Jangwon Kim, Naveen Kumar, Andreas Tsiartas, Ming Li, and Shrikanth S Narayanan, “Automatic intelligibility classification of sentence-level pathological speech,” *Computer speech & language*, 2015.
- [7] John HL Hansen, Liliana Gavidia-Ceballos, and James F Kaiser, “A nonlinear operator-based speech feature analysis method with application to vocal fold pathology assessment,” *IEEE Transactions on biomedical engineering*, 1998.
- [8] Ladan Baghai-Ravary and Steve W Beet, *Automatic speech signal analysis for clinical diagnosis and assessment of speech disorders*, Springer Science & Business Media, 2012.
- [9] Taha Khan, Jerker Westin, and Mark Dougherty, “Classification of speech intelligibility in parkinson’s disease,” *Biocybernetics and Biomedical Engineering*, vol. 34, 2014.
- [10] Zheli Liu et al., “GMM and CNN hybrid method for short utterance speaker recognition,” *IEEE Transactions on Industrial Informatics*, vol. 14, 2018.
- [11] Parvaneh Janbakhshi, Ina Kodrasi, and Hervé Bourlard, “Automatic dysarthric speech detection exploiting pairwise distance-based convolutional neural networks,” *ICASSP*, 2021.
- [12] Joel Shor et al., “Towards learning a universal non-semantic representation of speech,” *INTERSPEECH*, 2020.
- [13] Sepp Hochreiter and Jürgen Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, pp. 1735–80, 12 1997.
- [14] Alex Mayle, Zhiwei Mou, Razvan C Bunesco, Sadegh Mirshekarian, Li Xu, and Chang Liu, “Diagnosing dysarthria with long short-term memory networks,” in *INTERSPEECH*, 2019.
- [15] Myung Jong Kim, Beiming Cao, Kwanghoon An, and Jun Wang, “Dysarthric speech recognition using convolutional lstm neural network,” in *INTERSPEECH*, 2018.
- [16] Juliette Millet and Neil Zeghidour, “Learning to detect dysarthria from raw speech,” in *ICASSP 2019*.
- [17] Ashish Vaswani, Noam Shazeer, et al., “Attention is all you need,” *NeurIPS*, vol. 30, 2017.
- [18] Anmol Gulati, James Qin, et al., “Conformer: Convolution-augmented transformer for speech recognition,” *INTERSPEECH*, 2020.
- [19] Yu Zhang et al., “Bigssl: Exploring the frontier of large-scale semi-supervised learning for automatic speech recognition,” *IEEE Journal of Select Topics in Signal Processing*, 2021.
- [20] Joel Shor, Aren Jansen, Wei Han, Daniel Park, and Yu Zhang, “Universal paralinguistic speech representations using self-supervised conformers,” *ICASSP*, 2022.
- [21] Neil Zeghidour, Olivier Teboul, Félix de Chaumont Quitry, and Marco Tagliasacchi, “Leaf: A learnable frontend for audio classification,” *ICLR*, 2021.
- [22] Subhashini Venugopalan et al., “Comparing supervised models and learned speech representations for classifying intelligibility of disordered speech on selected phrases,” *INTERSPEECH*, 2021.
- [23] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *NeurIPS*, 2020.
- [24] Kaila L Stipancic, Yana Yunusova, James D Berry, and Jordan R Green, “Minimally detectable change and minimal clinically important difference of a decline in sentence intelligibility and speaking rate for individuals with amyotrophic lateral sclerosis,” *Journal of Speech, Language, and Hearing Research*, vol. 61, 2018.
- [25] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton, “Speech recognition with deep recurrent neural networks,” in *ICASSP*, 2013.
- [26] Yanzhang He, Tara N Sainath, Rohit Prabhavalkar, Ian McGraw, Raziq Alvarez, Ding Zhao, David Rybach, Anjali Kannan, Yonghui Wu, Ruoming Pang, et al., “Streaming end-to-end speech recognition for mobile devices,” in *ICASSP*, 2019.
- [27] A. Narayanan, R. Prabhavalkar, et al., “Recognizing long-form speech using streaming end-to-end models,” in *IEEE ASRU Workshop*, 2019.
- [28] Thomas Wolf et al., “Transformers: State-of-the-art natural language processing,” in *EMNLP: system demonstrations*, 2020, pp. 38–45.
- [29] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *NAACL*, 2019.
- [30] Marco Tagliasacchi, Beat Gfeller, Félix de Chaumont Quitry, and Dominik Roblek, “Self-supervised audio representation learning for mobile devices,” *arXiv:1905.11796*, 2019.
- [31] Chitrallekha Bhat and Helmer Strik, “Automatic assessment of sentence-level dysarthria intelligibility using blstm,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 2, pp. 322–330, 2020.
- [32] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter, “Audio set: An ontology and human-labeled dataset for audio events,” in *ICASSP*, 2017.
- [33] Ayush Tripathi, Swapnil Bhosale, and Sunil Kumar Koppurapu, “A novel approach for intelligibility assessment in dysarthric subjects,” in *ICASSP*, 2020, pp. 6779–6783.