



DeepFry: Identifying Vocal Fry Using Deep Neural Networks

Bronya R. Chernyak^{*1}, Talia Ben Simon^{*2}, Yael Segal^{*1}, Jeremy Steffman³, Eleanor Chodroff⁴,
Jennifer S. Cole³, Joseph Keshet¹

¹Faculty of Electrical and Computer Engineering, Technion–Israel Institute of Technology, Israel

²Department of Computer Science, Bar-Ilan University, Israel

³Department of Linguistics, Northwestern University, USA

⁴Department of Language and Linguistic Science, University of York, UK

chernroni@gmail.com

Abstract

Vocal fry or creaky voice refers to a voice quality characterized by irregular glottal opening and low pitch. It occurs in diverse languages and is prevalent in American English, where it is used not only to mark phrase finality, but also sociolinguistic factors and affect. Due to its irregular periodicity, creaky voice challenges automatic speech processing and recognition systems, particularly for languages where creak is frequently used.

This paper proposes a deep learning model to detect creaky voice in fluent speech. The model is composed of an encoder and a classifier trained together. The encoder takes the raw waveform and learns a representation using a convolutional neural network. The classifier is implemented as a multi-headed fully-connected network trained to detect creaky voice, voicing, and pitch, where the last two are used to refine creak prediction. The model is trained and tested on speech of American English speakers, annotated for creak by trained phoneticians.

We evaluated the performance of our system using two encoders: one is tailored for the task, and the other is based on a state-of-the-art unsupervised representation. Results suggest our best-performing system has improved recall and F1 scores compared to previous methods on unseen data.

Index Terms: creaky voice, vocal fry, convolutional neural networks, self-supervised speech representation

1. Introduction

Vocal fry, also known as Creaky Voice, is a type of phonation employed across a wide variety of languages with different linguistic and extralinguistic functions. A typical creaky voice has a low rate of vocal fold vibration (pitch), irregular pitch, and constricted glottis, characterized by a small peak glottal opening, long closed phase, and low glottal airflow [1]. Keating *et al.* [1] outline several kinds of creaky voice, each of which manifests a slightly different set of acoustic properties.

Cross-linguistically, creaky phonation plays many phonological roles [2]. It can serve as an utterance-final marker, signal phonemic contrasts with other voice qualities, or be an additional acoustic cue to enhance different contrasts, such as tone (as in Mandarin or Cantonese) [3, 4]. It is used as a variant of glottal stop in many languages. Creaky phonation also plays a role in social interaction: it can indicate the end of a conversational turn (Finnish) [5], indicate an irritation (Vietnamese) [6], and be a marker to establish identities [2].

Creaky voice challenges automatic speech processing algorithms due to its irregular periodicity. As a result, algorithms for

pitch tracking, spectral analysis, speaker verification, and automatic speech recognition might fail to operate in their full capacity [7, 8]. This is also a challenge in automatic processing of corpus phonetic and phonological analyses [9, 10, 11, 12]. We believe the community will benefit from an open-source tool for automatic creak detection.

There have been several studies on algorithms for the detection of creaky voice. Early methods to detect creaky voice were based on ad-hoc signal processing techniques. Vishnubhotla and Espy-Wilson [13] proposed a set of rules on the AMDF measure of periodicity. Ishi *et al.* [14] suggested to represent the speech by a pulse-synchronized analysis and then use a comparison of intra-frame periodicity and inter-pulse similarity.

Kane, Drugman, Gobl *et al.* [7, 15, 16, 17, 18] proposed several algorithms which are all based on specially designed acoustic features from the excitation and residual of the linear prediction filtering analysis of the speech and other acoustic features. This line of work used decision trees, fuzzy-input fuzzy-output support vector machine (F²SVM) algorithm, and shallow artificial neural networks for the task. Recent work has proposed more advanced learning algorithms, but have been designed for a unique and small data set. Tavi *et al.* [19] suggested an exploratory creak recognizer based on a convolutional neural network (CNN), which is generated specifically for emergency calls. Villegas *et al.* [20] used recurrent neural networks to detect creaky of single words in Burmese.

In contrast to previous work, we propose a deep learning model to detect creaky voice directly from an unprocessed speech signal. The model is built from an encoder and a multi-headed classifier trained together. The encoder is a CNN that takes the raw waveform and learns a representation of the signal. The classifier is a fully-connected network that gets as input the representation and outputs a detection score for creaky voice along with two additional auxiliary tasks: voicing and pitch. These additional predictions steer the overall network toward a better solution in detecting creaky voice [21].

We use two types of encoders. The first encoder is a specially designed encoder with a large time-span processing window proposed in [22]. It has a larger receptive field than the standard receptive field used in speech processing, and is able to “see” several pitch periods, even for very low pitch values. Therefore, we believe that this encoder can contribute to the task of creaky voice detection. This is particularly useful compared to methods that operate on a processed windowed signal (such as MFCC, STFT, and others) and might lose pitch information crucial for identifying creak.

The second encoder is based on a state-of-the-art self-supervised representation of the speech called HuBERT [23]

^{*}These authors contributed equally to this work. B. R. Chernyak is the corresponding author.

that yields state-of-the-art results for downstream tasks, such as automatic speech recognition [24, 25]. The latter encoder was pre-trained on 960 hours of read speech. The same multi-headed classifier was used in both models.

Trained phoneticians annotated two parallel corpora of connected speech each with 32 American English speakers (Datasets 1 and 2) along with a subset of speech from 14 American English speakers in the ALLSSTAR corpus [26]. The models were trained and developed on Dataset 1, and evaluated on Dataset 2 and the ALLSSTAR corpus.

Results suggest that both models outperform a heuristics-based baseline model and have better F1 and recall values than Kane *et al.* [17]. In addition, the model based on the pre-trained HuBERT representation was marginally better than the specially-designed encoder on unseen data. Our code and pre-trained models are publicly available here: <https://github.com/bronichern/DeepFry>.

2. Method

The input to all models is the raw waveform. Formally, we denote a speech waveform of T samples by $\bar{x} = (x_1, \dots, x_T)$, where $x_t \in \mathcal{X}$ for $t \in [1, T]$ and $\mathcal{X} \subseteq \mathbb{R}$. Our setting is designed to allow different input duration hence T is not fixed. We assume that there is a sequence of K multi-labels, $\bar{y} = (y_1, \dots, y_K)$, where each multi-label y_k is from the set $\mathcal{Y} = \{\text{creaky, not-creaky}, \text{voiced, unvoiced}, \text{pitch, no-pitch}\}$.

Our models are composed of two neural networks. The *encoder* $g: \mathcal{X} \rightarrow \mathcal{Z}^K$ is a function from the domain of \mathcal{X} to the embedding space $\mathcal{Z} \subseteq \mathbb{R}^N$. Specifically, the encoder generates a sequence of representational vectors $\bar{z} = (z_1, \dots, z_K)$, where $z_k \in \mathcal{Z}$, for all $1 \leq k \leq K$, such that $z_k \in \mathcal{Z}$ is the acoustic embedding of the k -th frame. The embeddings are then processed by a classifier that outputs a sequence of K predictions. The *classifier* is a function $f: \mathcal{Z}^K \rightarrow \mathcal{Y}^K$ from the domain of features vectors to the domain of target objects. Figure 1 depicts an encoder and the multi-headed classifier. The encoder g is a fully convolutional neural network and is based on the framework proposed in [22] for pitch estimation. This encoder gets as input a variable duration signal and outputs a sequence of embedding vectors every 5 msec (due to hardware limitations, the number of embedding vectors was restricted to less than 100 per input). It has a larger receptive field than the standard receptive field used in speech processing, and therefore can handle several pitch periods. We denoted this encoder-classifier combination as *DeepFry*.

In our work we also use another encoder which is based on the HuBERT model [23]. This model learns a representation of the speech in a self-supervised manner. It is a large CNN model trained to distinguish a series of subsequent samples from random future samples when some representations are masked, similar to the BERT model [27]. The rationale behind this concept is that subsequent samples are more likely to belong to the same phonetic class than random future samples. Here we used the HuBERT model that was pre-trained on 960 hours of read speech (LibriSpeech). The model operates on a 20 ms frame-rate window. We denoted this encoder-classifier combination as *HubertFry*.

Our classifier is based on the concept of multi-task learning (MTL). MTL optimizes multiple tasks simultaneously, under the assumption that the information shared by they task will help boost the performance of the model on the task of interest. Specifically our classifier f is aimed at detecting creaky voice (binary), voice/unvoiced (binary) and pitch (binary). The

voiced/unvoiced was annotated by expert phoneticians, and the pitch was extracted using [28].

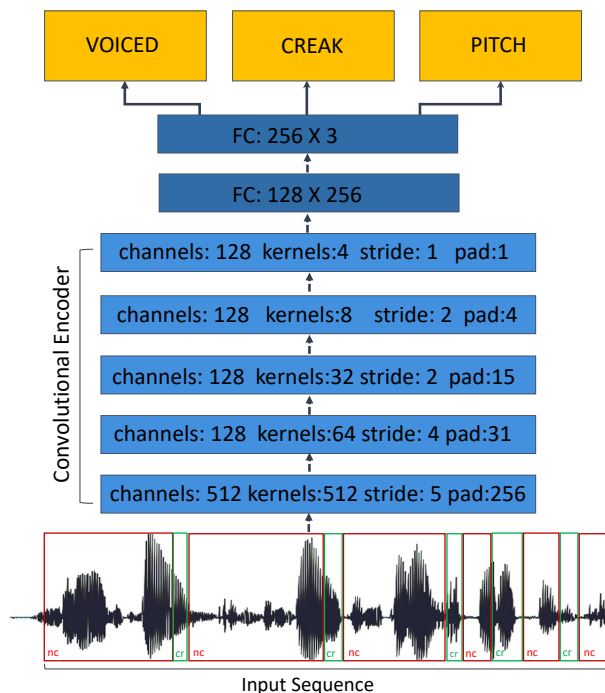


Figure 1: *Model architecture.* The encoder has 5 CNN layers, and the classifier has 2 fully-connected layers with 3 classification heads: creaky, voiced, pitch.

3. Datasets

Three datasets of American English connected speech were annotated for creaky voice. The first two datasets investigated the relationship between information structure and the realization of nuclear and prenuclear pitch accents in American English. The third dataset is the American English subset of the ALLSSTAR corpus (Archive of L1 and L2 Scripted and Spontaneous Transcripts and Recordings)[26]. Throughout the paper, the datasets are denoted as *Nuclear*, *Prenuclear* and *ALLSSTAR*, respectively.

The *Nuclear* and *Prenuclear* datasets contain speech from 32 native speakers of American English (*Nuclear*: 16 female, 16 male; *Prenuclear*: 11 female, 21 male), and the *ALLSSTAR* dataset contained speech from 14 native speakers of American English (7 female, 7 male). All participants were university-aged students.

In the *Nuclear* and *Prenuclear* datasets, participants read aloud a series of three-sentence stories. Only the final ‘target’ sentence was used in the present study of creak. Each participant read all 20 target sentences (within a unique story) in a randomized order in four separate blocks. Only the first two blocks were analyzed. Block 1 was read in a neutral speaking style and block 2 in a lively speaking style. The *ALLSSTAR* dataset contained readings of the first ten Hearing in Noise Test sentences in the ALLSSTAR HINT1 subset. The recordings in all three datasets were made at Northwestern University in a soundproof booth with a sampling rate of 22.05 kHz; all recordings were resampled to 16 kHz for the experiments. Sentences which contained a speech disfluency were removed from the

analysis. As a result, 1195 sentences were available for analysis in the *Nuclear* dataset, 1200 sentences in the *Preuclear* dataset, and 140 sentences in the *ALLSTAR* dataset.

The *Nuclear* dataset was divided into train (20 participants - about 24 minutes), validation (6 participants - about 6 minutes), and test (6 participants - about 8 minutes) folds where each fold contained an equal number of male and female participants. The *Preuclear* and *ALLSTAR* datasets were used for evaluation.

3.1. Annotation

To annotate the data, the transcripts were submitted to the Montreal Forced Aligner (MFA) for word- and phone-level alignments using the default American English pronunciation dictionary and acoustic model [29]. From these boundaries, we annotated the target sentences using two voice quality labels: modal or creaky voice, which are reasonably well-defined on sonorant segments. An interval was labeled as modal voice if a visible pitch track was present in Praat over a sonorant region, and one of the following two scenarios held true: 1) the interval had audible modal voice quality (conveyed audible pitch, had a similar ‘smooth’ voice quality as other prototypical modal regions), or 2) there was visible evidence of modal voicing from glottal pulses in the spectrogram or periodicity in the waveform. An interval was labeled as creak if no visible pitch track was present over a sonorant region in Praat and the above criteria for modality were not met. The utterance boundaries were refined during annotation and resubmitted to the MFA for a more precise phone-level alignment. We note, that in our labeling process, only unvoiced phonemes could not be tagged as creak.

4. Experiments

In this section, we present the results of our method. We begin by describing the hyper-parameters used to train our model. Then, we describe the measures used to evaluate our method and the adjustment we made to consider the different frame-rate of each method. Following that, we present ablation experiments of our model. We conclude this section by comparing our models to other methods on various unseen datasets and our test set. We trained *DeepFry* for 14 epochs, with an Adam optimizer, a learning rate of 0.001, a dropout of 0.1 and batch size 16. *HubertFry* was trained for six epochs, with an Adam optimizer, a learning rate of 0.001, and a batch size of 16. Both models were trained on the *Nuclear* dataset, and for both of them, we stipulated that the classification of a given frame was creak only if the frame was also predicted as voiced at inference time.

4.1. Evaluation method

We evaluated all models using Precision (P), Recall (R) and F1 measures. Precision is defined as the number of creaky frames correctly classified by the algorithm divided by the total number of frames classified as creaky by the algorithm. Recall is defined as the number of creaky frames correctly classified by the algorithm divided by the total number of ground truth creaky frames. The F1 measure is defined as $F1 = 2PR/(P + R)$.

Furthermore, our models, and the models we compare to, output predictions at a different frame-rate. Thus, for a fair comparison, all models classification was evaluated for a 20 ms frame-rate/window.

4.2. Baseline model

Our baseline model (referred to as the ‘Praat baseline’) used heuristics to detect the onset or offset of creaky voice for a given interval of speech [28]. The heuristic operated on the assumption that modal voice could be detected by trackable numeric pitch values, and creaky voice would be associated with an undefined pitch value in Praat. The script either moved *start-to-end* or *end-to-start* through Praat-extracted pitch values (time step = 0.15 s; pitch range = 50 to 350 Hz). When moving start-to-end, the script assumed modal voice was present until it encountered a sequence of undefined pitch values in frame n and frame $n + 2$. When moving end-to-start, the script assumed creaky voice was present until it encountered a sequence of numeric pitch values in frame n and frame $n - 2$. For intervals where the speech interval was not fully sonorant, creaky voice was frequently confused with a voiceless segment.

For the *Nuclear* and *Preuclear* datasets, the start-to-end mode was used for the first three regions of the target sentence, namely the subject noun phrase, the verb and following determiner, and the object noun phrase. The final phrase and the full sentences in the *ALLSTAR* dataset were analyzed using the end-to-start mode.

Table 1: *The contribution of the classification heads for final predictions. Precision (P), recall (R), and F1 scores on the Nuclear test set. Predictions are evaluated at a 20 ms frame-rate.*

Models	Classification heads			P	R	F1
	Creak	Voice	Pitch			
<i>DeepFry</i>	✓			68.35	69.64	68.99
	✓	✓		66.01	73.42	69.52
	✓		✓	67.30	69.67	68.46
	✓	✓	✓	67.43	76.27	71.57
<i>HubertFry</i>	✓			68.65	64.97	66.76
	✓	✓		69.29	64.52	66.82
	✓		✓	68.07	68.66	68.36
	✓	✓	✓	68.74	69.19	68.96

4.3. Ablation

We begin the experiments by analyzing the contribution of each additional task to the final detection of creak. To do so, we trained *DeepFry* and *HubertFry* with the following settings: (i) only with creak head (ii) creak head and voice head (iii) creak head and pitch head (iv) with creak, voice, and pitch heads.

Results are shown in Table 1, where each row represents a different experiment setting. Interestingly, it can be seen that each model benefits differently from each auxiliary task. By inspecting the F1 score, we can see that while *DeepFry* benefits more from the voice detection head than the pitch detection head, *HubertFry* presents higher performance with the pitch head. Finally, although training with each task individually does not necessarily significantly affect the performance, training with the combination of auxiliary tasks leads to increased performance.

4.4. Comparison to previous works

We now turn to compare our methods to previous works. We compared *DeepFry* and *HubertFry* to the baseline model presented in Section 4.2 on the test set of the *Nuclear* dataset. We

Table 2: Comparison to other methods. Precision (P), Recall (R), and F1 scores on different phonetic subsets of the test set of the Nuclear dataset. Sonorants include vowels, glides, liquids, and nasals. Metrics are reported as percentages. Predictions are evaluated at a 20 ms frame-rate.

Models	Vowels			Sonorants			All		
	P	R	F1	P	R	F1	P	R	F1
Praat baseline	59.66	59.21	59.43	60.07	60.60	60.33	32.93	60.60	42.67
Kane <i>et al.</i> [17]	86.63	30.14	44.72	87.53	28.53	43.03	75.98	28.53	41.48
<i>DeepFry</i>	81.76	75.71	78.62	82.20	76.27	79.12	67.43	76.27	71.57
<i>HubertFry</i>	78.99	68.33	73.28	79.69	69.19	74.07	68.74	69.19	68.96

Table 3: Evaluation on unseen data. Precision (P), recall (R), and F1 scores on different phonetic subsets of the Prenuclear and ALLSSTAR datasets. Sonorants include vowels, glides, liquids, and nasals. Metrics are reported as percentages. Predictions are evaluated at a 20 ms frame-rate.

Dataset	Models	Vowels			Sonorants			All		
		P	R	F1	P	R	F1	P	R	F1
Prenuclear	Praat baseline	44.01	51.99	47.67	46.59	54.89	50.40	25.52	54.89	34.84
	Kane <i>et al.</i> [17]	81.14	47.58	59.99	80.32	43.07	56.07	66.51	43.07	52.28
	<i>DeepFry</i>	74.65	77.50	76.05	74.31	76.51	75.39	61.16	76.51	67.98
	<i>HubertFry</i>	78.97	77.05	78.00	79.83	78.24	79.02	69.53	78.24	73.63
ALLSSTAR	Praat baseline	85.34	38.90	53.44	84.89	38.47	52.95	26.60	38.47	31.45
	Kane <i>et al.</i> [17]	97.15	36.50	53.06	97.39	34.02	50.42	90.03	34.02	49.38
	<i>DeepFry</i>	94.06	55.23	69.59	93.27	52.23	66.96	77.20	52.23	62.30
	<i>HubertFry</i>	75.86	71.37	73.55	76.94	70.33	73.49	66.74	70.33	68.49

also compared to the model was proposed by Kane *et al.* [17]¹. We found no other implementations or a detailed enough description of the work mentioned in Section 1.

Table 2 presents the comparison for different subsets of phonemes, to gain a deeper understanding of the results. The first column, *Vowels*, shows results only on the vowels phonemes. The second column, *Sonorants*, contains the results on the subset which includes vowels, glides, liquids, and nasals phonemes. The last column, shows results on all the phonemes. Note that the results in column *All* for *DeepFry* and *HubertFry* are the same results as in Table 1.

DeepFry received the best recall and F1 scores, while Kane *et al.* [17] received the best precision scores. Kane *et al.*'s low recall suggests that the high precision is due to miss detection bias. In addition, Praat baseline outperformed Kane *et al.* [17] on the *Vowels* and *Sonorants* subsets. Still, Praat baseline's precision score dramatically decreased when evaluated on all phonemes. That demonstrates this model's confusion between voiceless and creaky frames.

4.4.1. Testing on unseen data

Finally, we tested our models on the *Prenuclear* and *ALLSSTAR* datasets, which differ from the *Nuclear* dataset we trained on. Both the *Nuclear* and the *Prenuclear* datasets have the same lexical content. However, the *Prenuclear* and *ALLSSTAR* datasets have a different set of speakers and phonetic realizations. Furthermore, the *ALLSSTAR* dataset also has a different lexical content.

Results are shown in Table 3. It can be seen that in terms of F1 score, *HubertFry* outperformed the other methods across all datasets. However, Kane *et al.*'s precision is the highest preci-

sion on the *Vowels* and *Sonorants*, which is consistent with the results reported in Table 2. Additionally, compared to Table 2, Praat baseline algorithm recall and F1 scores drop. It is evident that *HubertFry* achieved higher F1 and recall scores than *DeepFry*, still the results of *DeepFry* are relatively higher than the other methods we compared to. Nevertheless, we note that *HubertFry* was pre-trained on 960 hours, which is much more data than the 24 minutes of the *Nuclear* dataset we trained *DeepFry* on. Furthermore, *HubertFry* has 95M parameters, compared to *DeepFry* which has less than 5M parameters. Therefore, we believe that training on more data, will improve the results of *DeepFry*.

5. Conclusions

In this work, we present a system composed of an encoder and a classifier for creak detection. We investigated two types of encoders that work on the raw wave: (i) *DeepFry* which is a small network with large receptive field and (ii) *HubertFry* which is a state-of-the-art encoder that was pre-trained on a lot of data and has a lot of parameters. Results suggest, that both of our implementations, achieve higher recall and F1 scores than the methods we compared to. Furthermore, it is evident that while *DeepFry* had better results on the test set of *Nuclear* dataset, *HubertFry* had higher recall and F1 scores on datasets from different distribution. We believe this is due to the vast amount *HubertFry* was trained on. Thus, it remains for future work to investigate the effect of training with more data on *DeepFry*.

6. Acknowledgements

Y. Segal is sponsored by the Ministry of Science & Technology, Israel.

¹We use the open-source framework - Voice Analysis Toolkit https://github.com/jckane/Voice_Analysis_Toolkit

7. References

- [1] P. A. Keating, M. Garellek, and J. Kreiman, "Acoustic properties of different kinds of creaky voice." in *ICPhS*, vol. 2015, no. 1, 2015, pp. 2–7.
- [2] L. Davidson, "The versatility of creaky phonation: Segmental, prosodic, and sociolinguistic uses in the world's languages," *Wiley Interdisciplinary Reviews: Cognitive Science*, vol. 12, no. 3, p. e1547, 2021.
- [3] K. M. Yu, "Laryngealization and features for chinese tonal recognition," in *Eleventh Annual Conference of the International Speech Communication Association*, 2010.
- [4] J. Kuang, "Creaky voice as a function of tonal categories and prosodic boundaries." in *INTERSPEECH*, 2017, pp. 3216–3220.
- [5] R. Ogden, "Turn transition, creak and glottal stop in finnish talk-in-interaction," *Journal of the International Phonetic Association*, vol. 31, no. 1, pp. 139–152, 2001.
- [6] H. Mixdorff, N. H. Bach, H. Fujisaki, and M. C. Luong, "Quantitative analysis and synthesis of syllabic tones in vietnamese," in *Eighth European Conference on Speech Communication and Technology*, 2003.
- [7] A. Cullen, J. Kane, T. Drugman, and N. Harte, "Creaky voice and the classification of affect," *Proceedings of WASSS, Grenoble, France*, 2013.
- [8] J. Keshet, "Automatic speech recognition: A primer for speech-language pathology researchers," *International journal of speech-language pathology*, vol. 20, no. 6, pp. 599–609, 2018.
- [9] M. Goldrick, J. Keshet, E. Gustafson, J. Heller, and J. Needle, "Automatic analysis of slips of the tongue: Insights into the cognitive architecture of speech production," *Cognition*, vol. 149, pp. 31–39, 2016.
- [10] E. Chodroff and C. Wilson, "Structure in talker-specific phonetic realization: Covariation of stop consonant *vot* in american english," *Journal of Phonetics*, vol. 61, pp. 30–47, 2017.
- [11] E. Chodroff, "Corpus phonetics tutorial," *arXiv preprint arXiv:1811.05553*, 2018.
- [12] K. C. Hall, J. S. Mackie, and R. Y.-H. Lo, "Phonological corpus-tools: Software for doing phonological analysis on transcribed corpora," *International Journal of Corpus Linguistics*, vol. 24, no. 4, pp. 522–535, 2019.
- [13] S. Vishnubhotla and C. Y. Espy-Wilson, "Automatic detection of irregular phonation in continuous speech," in *Ninth International Conference on Spoken Language Processing*, 2006.
- [14] C. T. Ishi, K.-I. Sakakibara, H. Ishiguro, and N. Hagita, "A method for automatic detection of vocal fry," *IEEE transactions on audio, speech, and language processing*, vol. 16, no. 1, pp. 47–56, 2007.
- [15] T. Drugman, J. Kane, and C. Gobl, "Resonator-based creaky voice detection," in *Thirteenth Annual Conference of the International Speech Communication Association*, 2012.
- [16] S. Scherer, J. Kane, C. Gobl, and F. Schwenker, "Investigating fuzzy-input fuzzy-output support vector machines for robust voice quality classification," *Computer Speech & Language*, vol. 27, no. 1, pp. 263–287, 2013.
- [17] J. Kane, T. Drugman, and C. Gobl, "Improved automatic detection of creak," *Computer Speech & Language*, vol. 27, no. 4, pp. 1028–1047, 2013.
- [18] T. Drugman, J. Kane, and C. Gobl, "Data-driven detection and analysis of the patterns of creaky voice," *Computer Speech & Language*, vol. 28, no. 5, pp. 1233–1253, 2014.
- [19] L. Tavi, T. Alumäe, and S. Werner, "Recognition of creaky voice from emergency calls," in *INTERSPEECH*, 2019, pp. 1990–1994.
- [20] J. Villegas, K. Markov, J. Perkins, and S. J. Lee, "Prediction of creaky speech by recurrent neural networks using psychoacoustic roughness," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 2, pp. 355–366, 2019.
- [21] Y. Shrem, M. Goldrick, and J. Keshet, "Dr. *vot*: Measuring positive and negative voice onset time in the wild," in *Proceedings of Interspeech*, 2019.
- [22] Y. Segal, M. Arama-Chayoth, and J. Keshet, "Pitch estimation by multiple octave decoders," *IEEE Signal Processing Letters*, vol. 28, pp. 1610–1614, 2021.
- [23] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [24] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [25] W.-N. Hsu, Y.-H. H. Tsai, B. Bolte, R. Salakhutdinov, and A. Mohamed, "Hubert: How much can a bad teacher benefit asr pre-training?" in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 6533–6537.
- [26] A. Bradlow, L. Ackerman, L. Burchfield, L. Hesterberg, J. Luque, and K. Mok, "Allstar: Archive of l1 and l2 scripted and spontaneous transcripts and recordings," in *Proceedings of the International Congress on Phonetic Sciences*, 2010, pp. 356–359.
- [27] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *NAACL-HLT*, 2019.
- [28] P. Boersma *et al.*, "Accurate short-term analysis of the fundamental frequency and the harmonics-to-noise ratio of a sampled sound," in *Proceedings of the institute of phonetic sciences*, vol. 17, no. 1193. Citeseer, 1993, pp. 97–110.
- [29] M. McAuliffe, M. Socolof, S. Mihuc, M. Wagner, and M. Sonderegger, "Montreal forced aligner: Trainable text-speech alignment using kaldif." in *Interspeech*, vol. 2017, 2017, pp. 498–502.