

On the Impact of Word Error Rate on Acoustic-Linguistic Speech Emotion Recognition: An Update for the Deep Learning Era

Shahin Amiriparian¹, Artem Sokolov^{2,3}, Ilhan Aslan², Lukas Christ¹, Maurice Gerczuk¹, Tobias Hübner¹, Dmitry Lamanov², Manuel Milling¹, Sandra Ottl¹, Ilya Poduremennykh², Evgeniy Shuranov^{2,4}, Björn W. Schuller^{1,5}

¹EIHW – Chair of Embedded Intelligence for Health Care and Wellbeing, University of Augsburg, Germany

²Huawei Technologies

³HSE University, Nizhny Novgorod, Russia

⁴ITMO University, Saint Petersburg, Russia

⁵GLAM – Group on Language, Audio, & Music, Imperial College London, UK

shahin.amiriparian@uni-a.de, sokolov.artem@huawei.com, ilhan.aslan@huawei.com, lukas.christ@informatik.uni-augsburg.de, maurice.gerczuk@informatik.uni-augsburg.de, tobias.huebner@informatik.uni-augsburg.de, lamanov.dmitry@huawei.com, manuel.milling@informatik.uni-augsburg.de, sandra.ottl@informatik.uni-augsburg.de, poduremennykh.ilya@huawei.com, evgeniy.shuranov@huawei.com, schuller@ieee.org

Abstract

Text encodings from automatic speech recognition (ASR) transcripts and audio representations have shown promise in speech emotion recognition (SER) ever since. Yet, it is challenging to explain the effect of each information stream on the SER systems. Further, more clarification is required for analysing the impact of ASR’s word error rate (WER) on linguistic emotion recognition per se and in the context of fusion with acoustic information exploitation in the age of deep ASR systems. In order to tackle the above issues we create transcripts from the original speech by applying three modern ASR systems, including an end-to-end model trained with recurrent neural network-transducer loss, a model with connectionist temporal classification loss, and a WAV2VEC framework for self-supervised learning. Afterwards, we use pre-trained textual models to extract text representations from the ASR outputs and the gold standard. For extraction and learning of acoustic speech features, we utilise OPENSIMILE, OPENXBOW, DEEPSPECTRUM, and AUDEEP. Finally, we conduct decision-level fusion on both information streams – acoustics and linguistics. Using the best development configuration, we achieve state-of-the-art unweighted average recall values of 73.6 % and 73.8 % on the speaker-independent development and test partitions of IEMOCAP, respectively.

Index Terms: emotion recognition, automatic speech recognition, computational paralinguistics

1. Introduction

As technology is becoming increasingly ubiquitous, speech input is gaining popularity as an accessible interaction modality. The rise of voice assistants, e. g., Amazon’s Alexa, exemplifies this trend. While today’s technologies may understand speech commands well, the conversation quality is still far from what we as humans experience in interpersonal communication. Emotional expressions are a key part of interpersonal communication. They are embodied in our gestures, body posture, and speech. Humans typically express and recognise emotional speech effortlessly, while, for machines, recognising emotions in speech is still a hard challenge.

In this paper, we present an update to previous research (i. e., [1]) on the trade-off between automatic speech recogni-

tion (ASR) accuracy (i. e., Word error rate (WER)) and linguistic emotion recognition, and the impact thereof on the later fusion with voice based emotion recognition. Such an update is urgently required, as I) most papers analysing the fusion of acoustics and linguistics use human transcripts and not actual ASR (e. g., [2], [3], [4]), hence, oversimplifying the problem, and II) practically no systematic investigation of the WER on linguistic speech emotion recognition (SER) exists, besides [1] – however, more than a decade since this investigation has witnessed massive improvements in ASR in the era of deep ASR approaches, and III) the modelling of linguistic information itself has changed dramatically with the advent of deep text modelling and the existence of large pre-trained according models. Hence, a re-investigation is urgently needed. To provide a comprehensive and ecologically valuable overview, we juxtapose and contrast variations of contemporary solutions for ASR and feature sets for emotion recognition from both text and voice. We highlight the best performing fusion solution, that to the best of our knowledge sets a new state-of-the-art. Further, we describe in detail the overall system that we use for our experiments. Moreover, we provide different variations of every system’s component.

2. Methodology

In this section, we introduce feature extraction methods that are well suited to process acoustic and linguistic cues. The features are used as inputs to Support Vector Machines (SVMs) and, therefore, build the basis for our SER analysis. We further introduce several ASR approaches, which will be investigated with respect to their WER and corresponding suitability in the SER context.

2.1. Audio Features

We examine four different audio feature sets. The first feature set is extracted with the OPENSIMILE toolkit using the ComParE_2016.conf configuration file [5]. It contains 6 373 static features resulting from the computation of functionals (statistics) over low-level descriptor (LLD) contours¹ [5, 6]. A full description of the feature set can be found in [7].

¹<https://github.com/audeer/opensmile>

In addition to the default Computational Paralinguistics Challenge (ComParE) feature set, we provide Bag-of-Audio-Words (BoAW) features by using OPENXBOW [8]. These have been applied successfully for, e. g., acoustic event detection [9] and speech-based emotion recognition [10]. After a quantisation based on a codebook, audio chunks are represented as histograms of acoustic LLDs. One codebook is learnt for 65 LLDs from the COMPARE feature set, and another one for 65 deltas of these LLDs. Codebook generation is done by *random sampling* from the LLDs and its deltas in the training data. Each LLD and delta is assigned to 10 audio words from the codebooks with the lowest Euclidean distance. Subsequently, both BoAW representations are concatenated. Finally, a logarithmic term frequency weighting is applied to compress the numeric range of the histograms.

The feature extraction DEEP SPECTRUM toolkit² is applied to obtain deep representations from the input audio data utilising pre-trained Convolutional Neural Networks (CNNs) [11]. DEEP SPECTRUM features have been shown to be effective, e. g., for speech processing [12, 13] and sentiment analysis [14]. First, audio signals are transformed into Mel-spectrogram plots using a Hanning window of width 32 ms and an overlap of 16 ms. From these, 128 Mel frequency bands are computed. The spectrograms are then forwarded through a pre-trained DENSENET121 [15] and the activations from the ‘avg_pool’ layer are extracted, resulting in a 1 024 dimensional feature vector.

Another feature set is obtained through unsupervised representation learning with recurrent sequence-to-sequence autoencoders, using AUDEEP³ [16, 17]. This feature set models the inherently sequential nature of audio with Recurrent Neural Networks (RNNs) within the encoder and decoder networks [16, 17]. First, Mel-scale spectrograms are extracted from the raw waveforms. In order to eliminate some background noise, power levels are clipped below four different given thresholds in these spectrograms. The number of thresholds results in four separate sets of spectrograms per data set. Subsequently, a distinct recurrent sequence-to-sequence autoencoder is trained on each of these sets of spectrograms in an unsupervised way, i. e., without any label information. The learnt representations of a spectrogram are then extracted as feature vectors for the corresponding instance. Finally, these feature vectors are concatenated to obtain the final feature vector. For the results shown in Table 1, the autoencoders’ hyperparameters are not fine tuned.

2.2. Text Features

DeepMoji, proposed by Felbo et al. [18], is a model pre-trained for emotion-related text classification tasks. It consists of two bidirectional long short-term memory (LSTM) layers, followed by an attention layer and yields a sentence encoding of length 2 304. Even though DeepMoji is pre-trained on emotional tweets only, the authors show that it also performs well for other kinds of emotional text data, e. g., reports of emotional experiences. We extract DeepMoji sentence encodings via the PyTorch implementation *TorchMoji*⁴.

Moreover, we employ several variants of Bidirectional Encoder Representations from Transformers (BERT) [19] that has set new standards for many text processing tasks in recent years. In its *base* configuration, BERT consists of 12 transformer ([20]) encoder layers. This network is pre-trained on large text data sets using two unsupervised language modelling tasks, namely

masked word prediction and next sentence prediction. Here, we employ the pre-trained BERT-base model to obtain sentence encodings. The output of the last layer’s hidden state for the special token [CLS], followed by one tanh-activated linear layer (*pooler_output*) is considered as the sentence encoding.

ALBERT (*A Lite BERT*) [21] is a popular variant of BERT. It is of the same size as the original BERT model but comes with considerably less parameters due to parameter sharing across layers and factorisation of the embedding matrix. Furthermore, the next sentence prediction task in BERT’s pre-training has been replaced by sentence order prediction, i. e., deciding whether two sentences are given in the correct order. ALBERT has been shown to outperform BERT on many tasks. Similar to our BERT baseline, we take the *pooler_output* of ALBERT in its *base* version as our sentence encoding.

Another recent variant of the BERT language model (LM) is given by ELECTRA [22], referring to an alternative method of pre-training transformer language models. In this approach, corrupted input words are detected. First, a generator model corrupts the input sentence. Then, the discriminator, i. e., the actual language model, predicts for every word whether it has been changed by the generator or not. BERT-like transformer networks pre-trained in this fashion outperform other BERT variants on several tasks. The architecture of the model is nearly identical to BERT-base. We take the embedding of the special token [CLS] as the sentence encoding. For all three BERT variants, huggingface implementations and pre-trained weights^{5 6 7} are used to extract 768 features.

2.3. Automatic Speech Recognition

To obtain text encodings from audio waveforms, we employ several pre-trained ASR systems: a system based on QuartzNet (QN), streaming Transformer Transducer (TT) and wav2vec (W2V). Transformer Transducer (TT) [23] is an end-to-end model trained with RNN-Transducer (RNN-T) loss [24]. Its encoder implementation entails Transformer blocks with multi-headed self-attention masking future context making the network suitable for stream audio processing. The label encoder of this architecture can be interpreted as a small built-in internal LM as it takes the previous predicted output label as input. The joint network combines audio and label encoder outputs and passes them to the the final softmax. Our solution is trained on LibriSpeech [25], CommonVoice [26], and Tedlium [27]. Additionally, we utilise our internal audio sets with various eastern accents for model fine-tuning. In total, about 4 000 hours of speech are used for TT training. For tokenisation, we use the Byte-Pair Encoding (BPE) [28] sentence-piece model with a vocabulary size of 4 096 items. Despite the fact that the system has an internal LM, we additionally evaluate our model in combination with an external LM based on the transformer architecture. The LM is trained on 30 gigabytes of corpora that include wiki texts and books. Furthermore, cold fusion is used to connect the outputs of the LM to the external LM with the lambda parameter set to 0.2.

QuartzNet (QN) [29] is a Connectionist Temporal Classification (CTC) [30] loss based model composed of blocks with separable convolutions and residual connections between them, with a fully connected decoder at the end. The model has fewer parameters than TT while still showing near state-of-the-art accuracy. In our experiments, first, a pre-trained configuration with

²<https://github.com/DeepSpectrum/DeepSpectrum>

³<https://github.com/auDeep/auDeep>

⁴<https://github.com/huggingface/torchMoji>

⁵<https://huggingface.co/bert-base-cased>

⁶<https://huggingface.co/albert-base-v2>

⁷<https://huggingface.co/google/electra-base-discriminator>

Table 1: *SER comparison of linguistic features on the Interactive Emotional Dyadic Motion Capture (IEMOCAP) corpus (Chance level: 25.0% UAR). Every text feature extractor is tested with different ASR systems as input as well as with the gold standard (GS). UAR: Unweighted Average Recall. QN: QuartzNet. TT: Transformer Transducer. W2V: wav2vec. LM: Language model.*

[UAR %]	GS Text		ASR QN		ASR QN-LM		ASR TT		ASR TT-LM		ASR W2V		ASR W2V-LM	
Network	Dev	Test	Dev	Test	Dev	Test	Dev	Test	Dev	Test	Dev	Test	Dev	Test
DEEPMOJI	61.7	63.0	52.2	49.6	52.3	48.6	53.8	54.0	53.8	45.2	57.4	58.0	56.3	55.2
BERT BASE	57.6	58.0	47.1	45.9	47.5	45.0	50.1	50.9	49.1	50.7	51.6	55.2	53.9	55.1
ALBERT BASE	47.9	52.7	38.1	40.3	40.8	42.4	42.2	43.1	41.4	44.3	42.7	46.3	44.3	49.2
ELECTRA BASE	56.9	56.2	44.2	43.1	46.7	43.2	48.2	45.5	47.1	45.6	52.6	49.7	53.3	49.9

15 blocks and 5 sub-blocks in each block is used. Subsequently, we fine-tune it on the dataset with British accents and recordings generated by a text-to-speech (TTS) system. For training the QN model, we employ around 2 000 hours of internal and public datasets. Unlike TT, we set up the configuration to predict graphemes. By default, the CTC loss does not consider the use of a built-in LM. We use a 4-gram statistical language model learnt on Gigaword [31] and fuse it with QN during the inference decoding.

The wav2vec (W2V) [32] system is a new framework for self-supervised training. The model can be broken down into three parts: i) a feature encoder, which represents speech waveforms in latent states that concurrently goes further to a contextualised representation part, ii) a quantisation module, iii) a contextualised representation joined with Transformers learns the relative positional information in the latent states. A quantization module represents the infinite output of a feature encoder to a discrete set via product quantization. The framework exploits the CTC loss for training. We use a pre-trained large model from the official repository⁸. First, we choose the checkpoint obtained by the training on 60 *k* hours of LibriVox, subsequently, we finetune it with additional 960 hours of LibriSpeech. For the external language model, we fuse W2V, which is the same as for QN.

3. Experiments

3.1. Dataset

The Interactive Emotional Dyadic Motion Capture (IEMOCAP) dataset [33] is an English emotion dataset that comes with both textual transcriptions and raw audio information. The dataset contains scripted and improvised dialogues between 5 female and 5 male speakers. In order to be consistent with previous research with IEMOCAP, we choose the main emotions happiness (fused with excitement), sadness, anger, and neutral. The choice of emotions results in 5 531 utterances totaling 7.0 hours of audio data. In literature, there is no agreement on the partitioning of the dataset. In our experiments, we split the IEMOCAP dataset speaker-independently into session 1 – 3 for training, session 4 for development, and session 5 for testing.

LibriSpeech [25] is a publicly available and popular dataset for speech recognition system experiments and evaluations. It includes transcriptions for around 960 hours of public domain audio books dictated by many speakers. For our experiments, we use test-clean part which is about 4.5 hours of audio with 20 male and 20 female speakers.

3.2. Automatic Speech Recognition

We evaluate each of our models on LibriSpeech test-clean with and without the external LM to compare the accuracy. We measure Word error rate (WER) and character error rate (CER). As two models are adopted for different accents, state-of-the-art results for the chosen data are not expected. On the contrary, the W2V trained on data with the same distribution as LibriSpeech and setups with this model show more precise predictions for both datasets. The evaluation on test-clean compares the accuracy of models on public and well-known data. Our measurements for IEMOCAP and LibriSpeech test-clean are demonstrated in Table 2.

Table 2: *Evaluation results of ASRs on LibriSpeech test-clean and IEMOCAP waveforms. WER: Word error rate. CER: Character error rate.*

SetUp	LibriSpeech		IEMOCAP	
	WER [%]	CER [%]	WER [%]	CER [%]
QN	15.69	7.26	41.21	25.43
QN + LM	12.98	7.35	42.14	31.21
TT	4.98	1.89	31.81	22.30
TT + LM	4.93	1.86	31.47	22.08
W2V	2.16	0.57	25.71	13.56
W2V + LM	2.24	0.63	22.23	13.37

Table 3: *SER Results of audio features on IEMOCAP. The best resulting model of DEEPSPECTRUM is a DENSENET201 with 128 Mel bins and the viridis colour map.*

[UAR %]	Dev	Test
OPENSIMILE (ComParE_2016)	57.8	58.5
OPENXBOW ($N = 2\,000$)	55.7	59.1
DEEPSPECTRUM	53.2	59.8
AUDEEP ($X = -60\text{ dB}$)	55.0	53.3

3.3. Late Fusion

We apply late fusion in form of majority voting on the basis of SVM predictions on the individual feature sets (cf. Figure 1). We evaluate all combinations of at least three feature sets from the audio, Gold Standard (GS) human transcribed spoken content, and ASR-based spoken content transcription. Table 4 summarises our results. Taking into account a high number of feature set combinations, we restrict our evaluation of the ASR-based text

⁸<https://github.com/pytorch/fairseq/tree/master/examples/wav2vec>

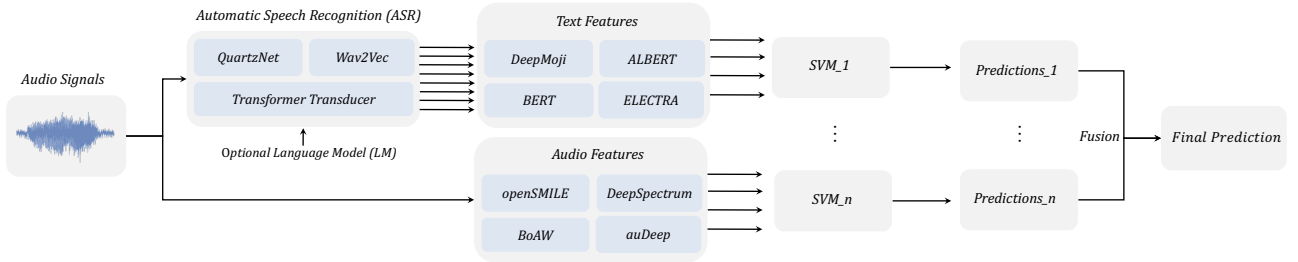


Figure 1: A general overview of our emotion recognition by multi-modal fusion of different trained Support Vector Machines (SVMs).

features to the ASR system with the best overall baseline performance which is W2V with LM. Feature set combinations, based on the GS text only, clearly outperform those, which are based on ASR-generated text. Considering the fact that this result is in line with performance of the individual feature sets (showed in Table 1), most likely, this behaviour is caused by the WER of the ASR system. However, when combining audio and text features, both ASR-based and GS-based feature sets lead to a similar performance, being higher than any individual information stream, both on average and for the best performance. Combinations of ASR-based and GS-based features show no improvement in the average Unweighted Average Recall (UAR) compared to the GS feature sets, most likely due to the similarity of the features. When considering all combinations of available feature sets the average and best performance can be further increased, which can most reasonably be explained by the vast number of combinations. Several combinations consisting of audio, GS-TEXT and ASR-TEXT features – including the combination GS-BERT, GS-DEEPMOJI, ASR-BERT, ASR-DEEPMOJI, ASR-ELECTRA, AUDEEP, DEEPSPECTRUM, OPENSMILE, BOAW – achieve the best observed performance of 73.6% UAR on the development set and a corresponding UAR of 73.8% on the test set.

Table 4: Results of the majority voting late fusion. The possible number of feature set combinations (#), as well as the mean UAR and standard deviation are reported. We further provide the UAR of the best performing feature set combination on the development set, as well as the corresponding performance of said combination on the test set.

[UAR %]	#	Mean Dev	Max Dev	Test
AUDIO	5	59.9 ± 0.1	60.8	63.3
GS-TEXT	5	60.3 ± 1.5	61.7	63.0
ASR-TEXT	5	54.6 ± 1.0	55.8	56.2
AUDIO + GS-TEXT	219	65.0 ± 3.5	71.0	69.9
AUDIO + ASR-TEXT	219	63.4 ± 3.2	69.5	70.5
ASR-TEXT + GS-TEXT	219	58.8 ± 3.6	63.3	64.8
ALL SYSTEMS	4017	66.2 ± 3.8	73.6	73.8

4. Discussion

When comparing the performance of different ASR systems in Table 2 and Table 1, a correlation between low WER values and high UAR values becomes obvious. Accordingly the Gold Standard System using human annotations, which is considered to have a much lower WER than any of the ASR systems, clearly achieves the highest UAR. A similar effect has previously been observed in [1], however, utilising a – from today’s point of view –

outdated ASR systems with a much more limited vocabulary size. Table 4 suggests that a higher number of considered feature sets leads to a higher UAR on average. This effect is known in general, however, it should be noted that the pairwise dependence of classifiers plays a considerable role in such a late fusion system [34]. A pairwise dependence of ASR-based and GS-based feature sets could therefore explain why combinations of both sets perform worse than combinations which combine either ASR-based or GS-based feature sets with audio-based feature sets. Assuming a high dependence between ASR-based and GS-based feature sets would further suggest that a well-suited weighted fusion method combining only audio-based and ASR-based features might further increase results towards the best-performing configuration introduced in 3.3, which combines two instances of BERT and DeepMojj features.

5. Conclusions

In this paper, we presented current ASR systems to create transcriptions for the linguistic SER. Without adapting the ASR systems to the target database IEMOCAP, we were able to achieve state-of-the-art results by fusing acoustic and linguistic information. We further observed that higher WERs on the ASR systems lead to higher UAR values for emotion recognition.

For future work, the number of feature sets and the respective feature set sizes can be reduced in order to increase computational efficiency. Furthermore, evaluation could be performed on more natural or in-the-wild databases. We mainly chose IEMOCAP as it is widely established and thereby suitable for comparison with state-of-the-art approaches. Moreover, IEMOCAP contains transcriptions making it easier to evaluate the impact of WER achieved by the ASR on the final emotion classification. Finally, the fusion with ASR could be implemented on the levels of the embeddings instead of the text.

6. Acknowledgements

This research was partly supported by the Affective Computing & HCI Innovation Research Lab between Huawei Technologies and University of Augsburg. We acknowledge funding from Deutsche Forschungsgemeinschaft (DFG) under grant agreement No. 421613952 (ParaStiChaD), and Zentrales Innovationsprogramm Mittelstand (ZIM) under grant agreement No. 16KN069455 (KIRun). The work of Artem Sokolov is partially supported by RSF (Russian Science Foundation) grant 20-71-10010.

7. References

- [1] F. Metze, A. Batliner, F. Eyben, T. Polzehl, B. Schuller, and S. Steidl, "Emotion recognition using imperfect speech recognition," in *Proc. Interspeech*, 01 2010, pp. 478–481.
- [2] E. Cambria, D. Hazarika, S. Poria, A. Hussain, and R. B. V. Subramanyam, "Benchmarking multimodal sentiment analysis," 2017.
- [3] J. Cho, R. Pappagari, P. Kulkarni, J. Villalba, Y. Carmiel, and N. Dehak, "Deep neural networks for emotion recognition combining audio and transcripts," *arXiv preprint arXiv:1911.00432*, 2019.
- [4] M. Chen and X. Zhao, "A multi-scale fusion framework for bimodal speech emotion recognition," in *Proc. Interspeech*, 2020, pp. 374–378.
- [5] F. Eyben, F. Wening, F. Groß, and B. Schuller, "Recent Developments in openSMILE, the Munich Open-Source Multimedia Feature Extractor," in *Proc. ACM Multimedia*, Barcelona, Spain, 2013, pp. 835–838.
- [6] B. Schuller, S. Steidl, A. Batliner, A. Vinciarelli, K. Scherer, F. Ringeval, M. Chetouani, F. Wening, F. Eyben, E. Marchi, M. Mortillaro, H. Salamin, A. Polychroniou, F. Valente, and S. Kim, "The Interspeech 2013 Computational Paralinguistics Challenge: Social Signals, Conflict, Emotion, Autism," in *Proc. Interspeech*, Lyon, France, 2013, pp. 148–152.
- [7] F. Wening, F. Eyben, B. Schuller, M. Mortillaro, and K. R. Scherer, "On the Acoustics of Emotion in Audio: What Speech, Music and Sound have in Common," *Frontiers in Emotion Science*, vol. 4, pp. 1–12, 2013.
- [8] M. Schmitt and B. W. Schuller, "openXBOW – Introducing the Passau Open-Source Crossmodal Bag-of-Words Toolkit," *Journal of Machine Learning Research*, vol. 18, pp. 1–5, 2017.
- [9] H. Lim, M. J. Kim, and H. Kim, "Robust Sound Event Classification Using LBP-HOG Based Bag-of-Audio-Words Feature Representation," in *Proc. Interspeech*, Dresden, Germany, 2015, pp. 3325–3329.
- [10] M. Schmitt, F. Ringeval, and B. Schuller, "At the Border of Acoustics and Linguistics: Bag-of-Audio-Words for the Recognition of Emotions in Speech," in *Proc. Interspeech*, San Francisco, CA, 2016, pp. 495–499.
- [11] S. Amiriparian, M. Gerczuk, S. Ottl, N. Cummins, M. Freitag, S. Pugachevskiy, and B. Schuller, "Snore sound classification using image-based deep spectrum features," in *Proc. Interspeech 2017*, Stockholm, Sweden, 2017, pp. 3512–3516.
- [12] S. Amiriparian, "Deep representation learning techniques for audio signal processing," Ph.D. dissertation, Technische Universität München, 2019.
- [13] S. Amiriparian, M. Gerczuk, S. Ottl, N. Cummins, S. Pugachevskiy, and B. Schuller, "Bag-of-deep-features: Noise-robust deep feature representations for audio analysis," in *Proc. IJCNN*, Rio de Janeiro, Brazil, 2018, pp. 2419–2425.
- [14] S. Amiriparian, N. Cummins, S. Ottl, M. Gerczuk, and B. Schuller, "Sentiment analysis using image-based deep spectrum features," in *Proc. ACIIW 2017*, San Antonio, TX, 2017, pp. 26–29.
- [15] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. CVPR*, 2017, pp. 4700–4708.
- [16] S. Amiriparian, M. Freitag, N. Cummins, and B. Schuller, "Sequence to Sequence Autoencoders for Unsupervised Representation Learning from Audio," in *Proc. DCASE 2017*, Munich, Germany, 2017, pp. 17–21.
- [17] M. Freitag, S. Amiriparian, S. Pugachevskiy, N. Cummins, and B. Schuller, "auDeep: Unsupervised Learning of Representations from Audio with Deep Recurrent Neural Networks," *Journal of Machine Learning Research*, vol. 18, pp. 1–5, 2018.
- [18] B. Felbo, A. Mislove, A. Søgaard, I. Rahwan, and S. Lehmann, "Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm," *arXiv preprint arXiv:1708.00524*, 2017.
- [19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 4171–4186.
- [20] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.
- [21] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "Albert: A lite bert for self-supervised learning of language representations," *arXiv preprint arXiv:1909.11942*, 2019.
- [22] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, "Electra: Pre-training text encoders as discriminators rather than generators," *arXiv preprint arXiv:2003.10555*, 2020.
- [23] Zhang, Q., Lu, H., Sak, H., Tripathi, A., McDermott, E., Koo, S., and Kumar, S., "Transformer transducer: A streamable speech recognition model with transformer encoders and rnn-t loss," *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7829–7833, 2020.
- [24] Graves, Alex., "Sequence transduction with recurrent neural networks," *arXiv preprint arXiv:11211.3711*, 2012.
- [25] Panayotov, V., Chen, G., Povey, D. and Khudanpur, S., "Librispeech: an asr corpus based on public domain audio books," *EEE international conference on acoustics, speech and signal processing*, pp. 5206–5210, 2015.
- [26] Ardila R, Branson M, Davis K, Henretty M, Kohler M, Meyer J, Morais R, Saunders L, Tyers FM, Weber G., "Common voice: A massively-multilingual speech corpus." *arXiv preprint arXiv:1912.06670*, 2019.
- [27] Rousseau, A., Deléglise, P. and Esteve, Y., "Ted-lium: an automatic speech recognition dedicated corpus," *LREC*, pp. 5206–5210, 2012.
- [28] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 1715–1725. [Online]. Available: <https://www.aclweb.org/anthology/P16-1162>
- [29] Kriman, Samuel, Stanislav Beliaev, Boris Ginsburg, Jocelyn Huang, Oleksii Kuchaiev, Vitaly Lavrukhin, Ryan Leary, Jason Li, and Yang Zhang, "Quartznet: Deep automatic speech recognition with 1d time-channel separable convolutions," *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6124–6128, 2020.
- [30] Graves, Alex, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber., "Sequence transduction with recurrent neural networks," in *Proceedings of the 23rd international conference on Machine learning*, pp. 369–376, 2012.
- [31] V. D. B. Napoles C, Gormley MR, "Annotated gigaword," *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction (AKBC-WEKEX)*, pp. 95–100, 2012.
- [32] Baevski, Alexei, Henry Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *34th Conference on Neural Information Processing Systems (NeurIPS 2020)*, 2020.
- [33] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. N. Chang, S. Lee, and S. S. Narayanan, "Iemocap: Interactive emotional dyadic motion capture database," *Language, resources and evaluation*, vol. 42, no. 4, p. 335, 2008.
- [34] L. Kuncheva, C. Whitaker, C. Shipp, and R. Duin, "Limits on the majority vote accuracy in classier fusion," *Formal Pattern Analysis and Applications*, vol. 6, pp. 22–31, 04 2003.