

Incorporating End-to-End Speech Recognition Models for Sentiment Analysis

Egor Lakomkin¹, Mohammad Ali Zamani¹, Cornelius Weber¹, Sven Magg¹ and Stefan Wermter¹

Abstract—Previous work on emotion recognition demonstrated a synergistic effect of combining several modalities such as auditory, visual, and transcribed text to estimate the affective state of a speaker. Among these, the linguistic modality is crucial for the evaluation of an expressed emotion. However, manually transcribed spoken text cannot be given as input to a system practically. We argue that using ground-truth transcriptions during training and evaluation phases leads to a significant discrepancy in performance compared to real-world conditions, as the spoken text has to be recognized on the fly and can contain speech recognition mistakes. In this paper, we propose a method of integrating an automatic speech recognition (ASR) output with a character-level recurrent neural network for sentiment recognition. In addition, we conduct several experiments investigating sentiment recognition for human-robot interaction in a noise-realistic scenario which is challenging for the ASR systems. We quantify the improvement compared to using only the acoustic modality in sentiment recognition. We demonstrate the effectiveness of this approach on the Multimodal Corpus of Sentiment Intensity (MOSI) by achieving 73,6% accuracy in a binary sentiment classification task, exceeding previously reported results that use only acoustic input. In addition, we set a new state-of-the-art performance on the MOSI dataset (80.4% accuracy, 2% absolute improvement).

I. INTRODUCTION

Speech emotion and affect recognition are crucial aspects for a coherent human-robot interaction and have recently received growing interest in the research community [1]. The quality of human-robot interaction could be improved significantly if a robot was able to consistently evaluate the emotional state of a person and its dynamics. For instance, if a robot was able to detect that a person is speaking in an angry way, it could use this information as a sign to adjust its behavior [2].

Humans integrate information from several input modalities, such as acoustic and visual, to estimate the emotional state of the speaker [3]. Recent computational models fuse different sources of information to yield better and more robust results. For example, combining visual, linguistic and acoustic modalities resulted in state-of-the-art performance on sentiment and emotion recognition tasks and it can be observed that the linguistic modality has the biggest contribution in the overall blend [1]. However, most experiments assume that the ground-truth (manually transcribed) spoken text transcriptions are available during training and testing phases. We argue that this setup differs from the real-life

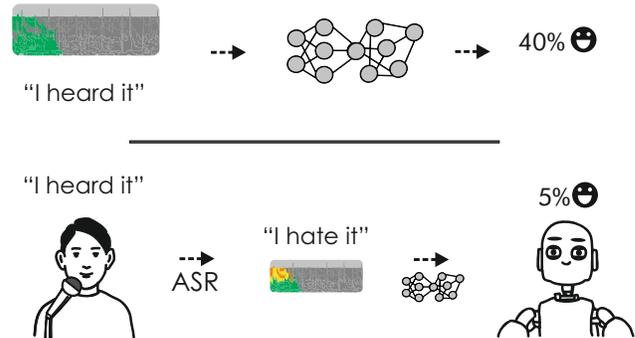


Fig. 1: Illustration of the core issue that we address in the paper. Sentiment recognition models are usually trained and evaluated using the manually transcribed text. In practice, we have to use automatic speech recognition systems to extract spoken text, which can contain errors affecting the overall performance.

condition, as we usually do not have access to the transcribed text in real time. In practice in the best case, we only extract an approximation of it: multiple hypotheses of the speech recognition system. Also, training on clean texts might result in overstating of the model’s performance and degradation during testing when we use even state-of-the-art speech recognition models.

This paper is built on our previous work [2], [4] on the application of neural emotion recognition models in the context of human-robot interaction. In this work, we conduct several experiments to address the discrepancy in the performance of a spoken sentiment recognition when the manually transcribed spoken text is not available. Our contribution is two-fold: a) We evaluate the performance of neural joint acoustic-linguistic sentiment recognition models in a human-robot interaction setup when the spoken text transcriptions are not available; b) We propose and evaluate the incorporation of recurrent character-level language model representations of spoken text for sentiment modeling to adapt to situations when the ASR output might produce word outputs with mistakes due to noise. We also analyze the models’ performance in acoustically clean conditions and when re-recorded on a robotic head, for the Multimodal Corpus of Sentiment Intensity (MOSI) dataset for sentiment identification. We compare three setups for spoken text extraction: 1) training our own end-to-end character-level neural speech recognition system, 2) using Google ASR API, 3) using ground-truth transcriptions, which we consider as

¹Knowledge Technology, Department of Informatics, University of Hamburg, Vogt-Koelln-Str. 30, 22527 Hamburg, Germany. <http://www.informatik.uni-hamburg.de/WTM/> {lakomkin, zamani, weber, magg, wermter}@informatik.uni-hamburg.de

the upper bound in our experiments.

The paper is organized as follows: section II introduces related work and section III describes our neural sentiment recognition model. Section IV outlines the methodology including the description of our neural ASR and the data used to train it. Section V introduces the conducted experiments on the original MOSI data and on the robot head recorded data in noise-realistic conditions.

II. RELATED WORK

Integration of multiple modalities like vision, auditory and linguistic with deep neural networks significantly boosted the overall performance of sentiment and emotion recognition [5]. For instance, individual modalities' representations and their paired combinations were fused by an outer tensor product [1]. The multi-attention recurrent network employs an attention method to model the integration of different modalities as has been done in multiple other works [6]. Conditioning on the context was shown beneficial on emotion and sentiment recognition [7]. As some modalities can have different contributions a certain time step, a gating mechanism which is trained with reinforcement learning, is introduced to switch on or off a particular modality's input [8].

In this work, we are focusing on the acoustic and linguistic modalities, considering situations when the speaker might not be directly observable. Jin et al. [9] used various hand-crafted acoustic and lexical features followed by late decision fusion for classification. Multimodal word-level alignment produced state-of-the-art results on the emotion and sentiment recognition tasks [10]. Automatic generation of ensemble trees with SVM classifiers as nodes was applied for affective analysis [5]. Hybrid attention mechanism was introduced to fuse acoustic and linguistic information [11]. Aldeneh et al. [12] evaluated several end-to-end approaches of pooling lexical and acoustic features extracted by recurrent neural networks, which encode each modality separately for speech valence estimation. Attention-based convolutional neural networks were proposed for acoustic-only emotion recognition [13]. Etienne et al. [14] achieved state-of-the-art results among systems using only audio modality combining convolutional and recurrent layers. However, Schuller et al. [15] compared off-the-shelf acoustic feature extractors with end-to-end approaches and demonstrated that end-to-end methods still do not consistently surpass the handcrafted representations on the paralinguistic tasks. Several ways of transfer learning to encode acoustic signals were proposed recently: tuning audio representations trained initially for other auxiliary tasks, like gender and speaker identification [16] or speech recognition [17], [18].

The main difference between our work and the previous research is that we do not assume that the ground-truth transcriptions are given as input to the model. We observe that the linguistic modality makes the biggest contribution to the overall classification [19]. In this work, we assume to have access only to the acoustic input, and spoken text is therefore extracted by a separate module.

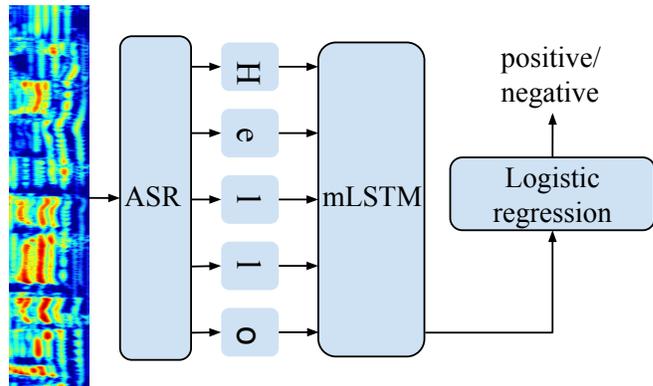


Fig. 2: Our sentiment recognition model based on the ASR character output. The multiplicative LSTM model is used to encode spoken text to a fixed-length vector with a logistic regression on top modeling the sentiment. The mLSTM model is pretrained in an unsupervised way on the Amazon reviews with a language-modelling objective.

III. SENTIMENT RECOGNITION MODEL

This section describes our sentiment recognition model. Firstly, we outline the neural spoken text representation network, followed by our proposed methods to extract transcriptions. We conclude with acoustic features description, which we used for joint acoustic-linguistic sentiment analysis.

A. Model Architecture

The essential part of our architecture is the character-level recurrent neural network for spoken text encoding. We use a single layer multiplicative LSTM [20] (mLSTM) model with 4,096 nodes and we use the hidden state corresponding to the last input character to represent the whole textual input. The mLSTM model is trained on the vast amounts of Amazon reviews in a language modelling setup [21]. This model with a linear classifier on top showed state-of-the-art performance on the Stanford Sentiment Treebank outperforming more complex architectures. This indicates that even though the mLSTM model was trained fully unsupervised, it was capable to capture the concept of sentiment just by learning to predict the next character given the context. Our main intuition to use this character-level model for spoken sentiment recognition opposes to the majority of previous work incorporating word-level processing is two-fold: a) a character-level model is capable of dealing with the spelling mistakes or out-of-vocabulary words produced by the ASR model, b) the representations learnt by the mLSTM on the reviews could be useful as well in the spoken sentiment analysis.

The overall architecture is shown in Figure 2. We feed the input text representation to an L2-regularized logistic regression for binary sentiment classification. The same architecture is used to train the acoustic model, where we use off-the-shelf acoustic feature descriptors (see section III.C) instead of the pre-trained mLSTM as a feature extractor. In addition, we fuse sentiment predictions of acoustic and lin-

guistic classifiers by computing their weighted combination. We tune hyperparameters (logistic regression regularization strength and classifiers fusion weight) on the validation data.

B. Spoken Text Extraction

In our experiments, we train our own end-to-end ASR model (see section IV). Our ASR model computes character probabilities for each timestep and the final transcription is extracted by simply taking the most probable character for each frame (greedy decoding). The only post-processing steps we perform are a) we merge together blank symbols (used by the ASR model to denote a non-speech character or a change between different characters and displayed here as '~') and character repetitions (*aaa_bbb - a.b*), b) capitalized characters are lowercased with a space imputed in front (*H_i_H_o_w_A_r_e_Y_o_u - hi how are you*). We do not apply any language model to correct the potential spelling mistakes of the model. For comparison, we extract the most probable transcription using Google Web Speech API¹.

C. Acoustic Feature Extraction

We use the *COMPARE 2016* feature set extracted by the OpenSMILE toolkit [22]. This feature set contains 6,373 features resulting from the computation of various functionals (for example mean, standard deviation, maximum value) over low-level descriptors (like mel cepstral coefficients, pitch, loudness, etc.) described in [23].

IV. SPEECH RECOGNITION MODEL

To extract spoken text from the acoustic signal we train an end-to-end neural automatic speech recognition system. In this section, we outline the ASR architecture, describe the data preprocessing and feature extraction pipeline, and datasets used for training.

A. Architecture

Our ASR model (see Figure 3) is based on several stacked Long Short Term Memory (LSTM) [24] recurrent layers [25]. The model contains five bi-directional LSTMs with batch normalization layers in between [25] processing log mel-spectrograms extracted from input audio. We use 40 mel coefficients, extracted using a Hamming window of 25ms width and 10ms stride. We stack three consecutive speech frames resulting in 120 features for each timestep. Frame stacking greatly speeds up the training and makes it more stable as the input and output are three times shorter. Recurrent layers are followed by a fully connected layer with a softmax activation on top, modelling the character probability distribution for each speech frame. Along with the standard 26 English characters, we introduce capital characters to the overall set denoting the beginning of words for the model. Overall, our model has around 61 million parameters. Connectionist Temporal Classification (CTC) [26] is used as a loss criterion to measure how the alignment produced by the network matches to the ground-truth transcription.

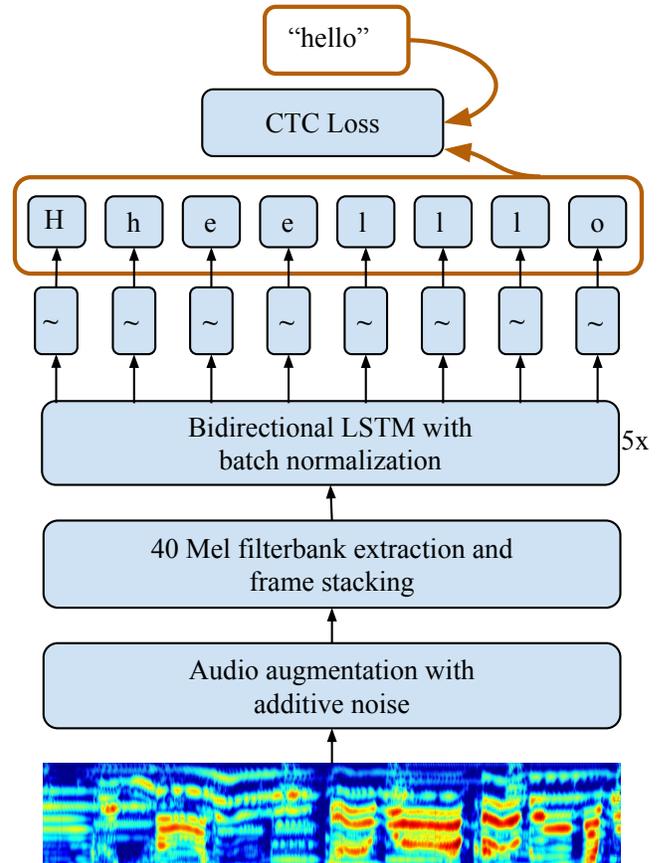


Fig. 3: Architecture of the ASR model used in this work. The model is trained end-to-end by mapping mel-spectrograms to characters using a stack of LSTM layers and CTC loss function.

The Stochastic Gradient Descent optimizer is used in all experiments with a learning rate of 0.0001 and a Nesterov momentum value of 0.9, clipping the norm of the gradient at the level of 400 with a batch size of 24. During the training, we apply learning rate annealing with a factor of 1.1. We apply the SortaGrad algorithm [25] during the first epoch by sorting utterances by their duration [27]. We select the model with the best word error rate measured on the LibriSpeech validation set (combined clean and noisy splits) to prevent model overfitting. We train the model on two GTX1080TI, and it takes around five days to train the model until convergence.

B. Data Augmentation

Previous research in end-to-end speech recognition demonstrated the importance of introducing random perturbations into the speech signal like a change of pitch, tempo, loudness, and adding noise [4], [25], [27], [28]. Since such perturbations do not alter the target label (spoken text in the case of speech recognition, or an emotion category), they can be conveniently applied with some occurrence probability during training. Data augmentation can be considered also as

¹<https://pypi.org/project/SpeechRecognition/>

a way to increase the training data size. In the case of human-robot interaction, it is crucial to have a noise-robust model due to the presence of a robot’s ego-noise or background noise.

In our experiments, we 1) change the tempo of the recording by sampling the speed factor uniformly in a range of [85, 120] percent, 2) change the loudness of the recording by sampling gain uniformly in a range of [-6, 5] dB, 3) add random background noise, where non-speech noise samples are selected from Google’s AudioSet², by sampling the noise-to-signal ratio uniformly in the range [0.1, 0.4] and mixing it with the original utterances resulting in over 530,000 samples of 10 seconds in length, and 4) perturb vocal tract length in the range [0.9, 1.1].

C. Speech Data

In our experiments, we use only freely available datasets. We concatenate five datasets to train the ASR model: LibriSpeech, TED-LIUM v3, Mozilla Common Voice, Google Speech Commands v2 and VoxForge. LibriSpeech [29] contains around 1,000 hours of English-read speech from audiobooks. TED-LIUM v3 [30] is a dataset composed of transcribed TED talks, containing 452 hours of speech and 2,351 speakers. VoxForge is an open-source collection of transcribed recordings collected using crowd-sourcing. We downloaded all English recordings³, which are around 100 hours of speech. Common Voice⁴ is a crowdsourced dataset, where utterances were collected through a web interface. Each participant was asked to pronounce a predefined text and submit it to the website. In addition, other volunteers were asked to check if the spoken text matches the actual requested one. Overall, Common Voice contains around 300 hours of validated speech data. Google Speech Commands contains 100,000 short recordings with only one word pronounced (out of 30 possible ones) by a variety of speakers. Overall, 850,000 utterances containing 1,600 hours of speech from more than 3,000 speakers are used to train the ASR model. We conduct no preprocessing other than the conversion of recordings to WAV format with single-channel 16-bit signed integer format and a sampling rate of 16,000. Utterances longer than 15 seconds are filtered out due to GPU memory constraints.

V. EXPERIMENTS AND ANALYSIS

In this section, we outline the data used in our experiments and the evaluation protocol and metrics, followed by the evaluation results and comparisons to previous work.

A. Data and Evaluation Measure

1) *CMU-MOSI*: This dataset is a multimodal sentiment intensity and subjectivity dataset consisting of 93 review videos in English with 2,199 utterance segments collected

²<https://research.google.com/audioset/>

³http://www.repository.voxforge1.org/downloads/SpeechCorpus/Trunk/Audio/Original/48kHz_16bit/

⁴<https://voice.mozilla.org/>

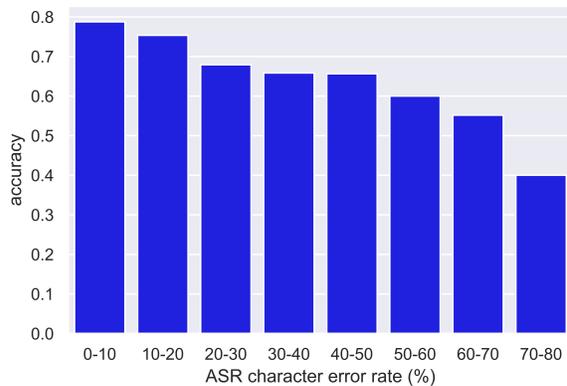


Fig. 4: Histogram plot visualizing the dependency between the character error rate of our trained ASR system and the accuracy of the sentiment recognition model. We note that even for samples with 20-30% character error rate the accuracy score is greater than 70%.

from YouTube [32]. Each segment is labelled by five individual annotators between -3 (strong negative) and +3 (strong positive). We binarize the labels (positive and negative) based on the sign of the annotations’ average to compare to the previously published methods. We use an 80%-20% training-testing speaker-independent split following the same strategy as in previous work [1], leaving 10% of training data for validation. Specifically, there are 1,279 utterances for training, 233 for validation and 686 utterances for testing. We report accuracy and macro F1-score calculated over the test set.

B. Human Robot Simulation

To simulate an acoustically close to real-life scenario, we re-recorded the CMU-MOSI corpus in our lab. The experimental setup [33], as shown in Fig. 5, consists of loudspeakers which are placed around the Soundman wooden head [34], behind the white display between 0° and 180° along the azimuth plane with the same elevation. The Soundman head is 1.6 meters away from the speakers, designing a far-field communication setup. We only use 4 speakers out of 13 speakers. In our previous work [2], [4] we use the iCub robotic head to test the robustness of our models against the robot’s ego noise. However, in this paper, we use the Soundman wooden head to focus on background noise generated by the projectors, computers, air conditioner, power sources as well as noise from airplanes frequently passing nearby, and reverberation noise. The entire recording was done in our lab.

C. Spoken Text Extraction Performance

We calculate Word Error Rate (WER) and Character Error Rate (CER) for Google ASR and our ASR model evaluated on the MOSI and MOSI-Soundman datasets. Google ASR has on average 51.3% WER and 39.7% CER on the original MOSI dataset, and 49.1% WER and 38.6% CER on the MOSI-Soundman. Our ASR model has 53.2% WER and

TABLE I: Sentiment prediction results on the CMU-MOSI test set. The best result of the model which does not use ground-truth transcriptions is highlighted in bold. Text source denotes how the spoken text was extracted: either ground-truth from the MOSI data or ASR output (Google Web Speech API or our ASR) was used. MOSI-Soundman is the MOSI dataset re-recorded in our lab emulating human-robot interaction scenario.

Model	Dataset	Modalities	Text Source	Accuracy	F-Score
Tensor fusion network, Zadeh et al. [1]	MOSI	audio + text	ground-truth	74.6%	74.5%
MARN, Zadeh et al. [6]	MOSI	text	ground-truth	77.1%	77.0%
Word-level alignment, Gu et al. [10]	MOSI	text+audio+vision	ground-truth	76.4%	76.8%
Recurrent multi-stage fusion, Liang et al. [31]	MOSI	text+audio+vision	ground-truth	78.4%	78.0%
Ours char-RNN + LogReg	MOSI	text	ground-truth	80.4%	79.8%
Ours, char-RNN + LogReg	MOSI	audio	-	54.8%	54.1%
Ours, char-RNN + LogReg	MOSI-Soundman	audio	-	53.6%	53.4%
Ours, char-RNN + LogReg	MOSI	text	our ASR	69.9%	68.7%
Ours, char-RNN + LogReg	MOSI	text	Google ASR	69.6%	69.3%
Ours, char-RNN + LogReg (2x models fused)	MOSI	text	Google ASR + our ASR	72.3%	71.9%
Ours, char-RNN + IS16 + LogReg (3x models fused)	MOSI	text + audio	Google ASR + our ASR	73.6%	73.1%
Ours, char-RNN + LogReg	MOSI-Soundman	text	our ASR	58.4%	58.2%
Ours, char-RNN + LogReg	MOSI-Soundman	text	Google ASR	67.9%	67.7%
Ours, char-RNN + IS16 + LogReg (2x models fused)	MOSI-Soundman	text + audio	Google ASR	70.1%	69.7%
Ours, char-RNN + IS16 + LogReg (3x models fused)	MOSI-Soundman	text + audio	Google ASR + our ASR	70.2%	69.8%



Fig. 5: Lab setup of the Soundman head in front of loud-speakers behind a screen. The positions of the speakers are highlighted with external light. See also [33].

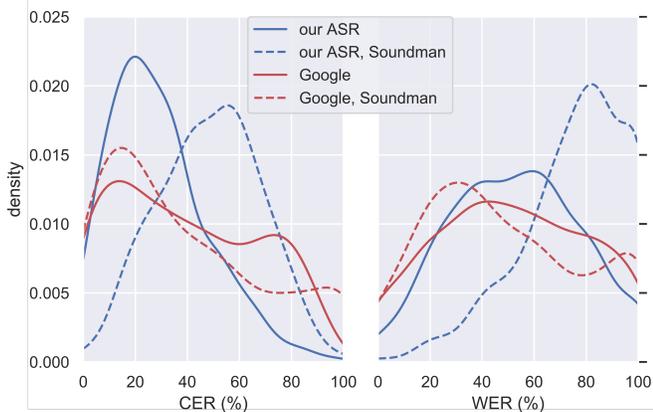


Fig. 6: CER and WER density plots for our ASR (blue) and Google ASR (red) evaluated on the MOSI (solid line) and MOSI-Soundman datasets (dashed line).

28.9% CER on the MOSI dataset, and 75.7% WER and 48.6% CER on the MOSI-Soundman. On the close-field communication (original MOSI) our model shows competitive results, while Google ASR being superior on the far-field scenario. Figure 6 shows the density plots of CER for our ASR and Google Web Speech API, demonstrating the difference in the models: the latter system has a lower WER on average and two peaks on the CER density plots, where the right one occurring at the 75%+ CER area could be explained by language model application. The language model corrects spelling mistakes, but, on the other hand, can change the word completely (for example, names) and increase CER while having better WER.

D. Experimental Results

We present the results of our experiments in Table I. We observe that the character-level mLSTM model pretrained on the Amazon reviews with a simple logistic regression can outperform all the previously published methods, which are more complex and integrate multiple modalities. One argument could be that this mLSTM model even though trained unsupervised has seen significantly more textual data than any work we compared it with. On the other hand, the pretrained word embeddings used in [1], for instance, can be an example of encoded external knowledge. Therefore, we believe that this pre-trained mLSTM model is extremely useful for the spoken sentiment recognition task. In addition, we achieve 69.9% and 69.6% accuracy score using the same model, but with the character output of our trained ASR system and Google ASR respectively. Interestingly, the fusion of these two models yields a significant gain in performance, achieving 72.3% accuracy and 73.6% accuracy by adding acoustic features. These ASR models make different and uncorrelated mistakes and this can explain the boost of combining them. Overall, the performance gap between the previous best-reported model using acoustic and linguistic modalities and our setup is around 1% accuracy and 1.3% F-score. However, in our experiments, we did not use ground-truth transcriptions, but only raw and processed

TABLE II: Examples from the CMU-MOSI dataset. For each example, we show ground-truth transcription, our ASR model output, sentiment ground-truth and model prediction, and ASR character error rate for this example.

ground-truth text	ASR transcription	Sentiment	Model output	CER
a) I hated it	I hated I	neg	pos	25%
b) its a pointless scene for the audience and the characters	it's appointmen seen from the audience tho the characters	neg	pos	22.9%
c) anyway oh you can see im still speechless this movie was just beautiful	ay ligt o you can see o stell speechlesses will be was beautiful	pos	pos	35.6%
d) who have named themselves after the places um to which they have traveled in which i think is a really nice	sen who rae named themselves after the places and to which with travelling which i think is a really nice to or	pos	pos	32.2%
e) yeah it really is good i mean	really is good i mean	pos	pos	26%
f) it was terrible	ws terrible i	neg	neg	30.7%
g) but if you are a child who grew up in that time period youre not going to enjoy this movie very much	but if youere a child who grew up in that time cerrod you're not going to enjoy this movvie very much because	neg	neg	13.9%
h) but nevertheless another really cool thing about this movie	but never was anoter really calling about this ruie	pos	neg	27.5%
i) um that being said you can tell that lot people were having fun with this	that being said you can sell that lot people were having to phumbl ois	pos	neg	20.3%

with ASR acoustic signal ones. Our results show that Google ASR is robust to the change of recording conditions and we get similar results on the MOSI and MOSI-Soundman data, while our ASR system performs significantly better on the original MOSI data. We hypothesize that our data augmentation pipeline needs to be improved further by simulating different room conditions during random training to achieve better results in noisy and reverberant conditions for far-field communication.

We provide several examples from the MOSI dataset in table II. The examples *b* and *i* demonstrate a situation when the ASR did not recognize correctly the key word, which changes the sentiment of the phrase completely. However, the examples *d* or *e* demonstrate that even with a relative high CER value of 32% the character-level model can tolerate those errors and correctly classify the sentiment.

E. Importance of ground-truth transcriptions for word-level sentiment model

As we observed a significant drop in performance when using spoken text extracted from Google ASR or our own ASR system, we performed an additional experiment using a word-level model with architecture similar to [35]. It consists of a 1-dimensional convolution network with 100 filters of sizes 2,3,4 and 5, followed by a fully connected layer with 400 units and an output node with a sigmoid activation for binary sentiment modelling. We achieved 74.7% accuracy on the MOSI test set, similar to spoken text-only results [1]. However, if we substitute test set transcriptions with the Google ASR results, we observe a drop to 56.8% accuracy, which can be a sign of significant overfitting to the specific words and text modality in general. This result is the additional testimony that it is crucial to take into account potential ASR mistakes during training to achieve robust sentiment recognition in practice.

VI. CONCLUSIONS

We addressed spoken sentiment recognition in conditions when ground-truth text is not available. Multiple previous

works demonstrated that the linguistic features dominate audio-visual input in sentiment and emotion recognition tasks. However, we note that those systems were trained and evaluated using human-transcribed text and, practically, we are not able to use it, for example, during human-robot interaction, when spoken text should be recognized in real time. We demonstrated the discrepancy in performance when ground-truth transcriptions are not present as input and ASR output is used instead on two models: character-level mLSTM with the linear classifier and word-level CNN. However, we observe significant improvements over the acoustic-only baseline on the original and re-recorded MOSI data by adding the ASR hypothesis, which shows that the linguistic modality is still crucial to achieving high-performance sentiment recognition.

In future work, we plan to investigate further ways to integrate multiple ASR hypotheses for robust sentiment and emotion recognition. The ensemble of Google Web Speech API and our ASR model show a significant boost in performance indicating the need of having a diverse set of hypotheses to make a better judgement of the affective state of a speaker. Character-level representations learned from the unsupervised language modelling task show very promising performance and we plan to research further whether a similar approach can be transferred to learning robust acoustic representations.

The demo, our pre-trained ASR model and its parameters are available at https://github.com/EgorLakomkin/icra_2019_speech

ACKNOWLEDGMENT

The authors thank Erik Strahl for his continuous support with the experimental setup, Julia Lakomkina for her help with illustrations, and Tayfun Alpay for his help in preparing the paper. This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 642667 (SECURE) and the German Research Foundation DFG under project CML (TRR 169).

REFERENCES

- [1] A. Zadeh, M. Chen, S. Poria, E. Cambria, and L.-P. Morency, "Tensor Fusion Network for Multimodal Sentiment Analysis," *Conference on Empirical Methods in Natural Language Processing*, pp. 1103–1114, 2017.
- [2] E. Lakomkin, M. A. Zamani, C. Weber, S. Magg, and S. Wermter, "EmoRL: Real-time Acoustic Emotion Classification using Deep Reinforcement Learning," in *Robotics and Automation (ICRA), 2018 IEEE International Conference on*, 2018, pp. 4445–4450.
- [3] R. W. Picard and R. W., *Affective computing*. MIT Press, 1997.
- [4] E. Lakomkin, M. A. Zamani, C. Weber, S. Magg, and S. Wermter, "On the Robustness of Speech Emotion Recognition for Human-Robot Interaction with Deep Neural Networks," *To appear in International Conference on Intelligent Robots (IROS)*, 2018. [Online]. Available: <http://arxiv.org/abs/1804.02173>
- [5] V. Rožgić, S. Ananthkrishnan, S. Saleem, R. Kumar, and R. Prasad, "Ensemble of SVM Trees for Multimodal Emotion Recognition," *Proceedings of The 2012 Asia Pacific Signal and Information Processing Association Annual Summit and Conference*, pp. 1–4, 2012. [Online]. Available: http://www.apsipa.org/proceedings_2012/papers/157.pdf
- [6] A. Zadeh, P. P. Liang, S. Poria, P. Vij, E. Cambria, and L.-P. Morency, "Multi-attention Recurrent Network for Human Communication Comprehension," 2 2018. [Online]. Available: <http://arxiv.org/abs/1802.00923>
- [7] S. Poria, E. Cambria, D. Hazarika, N. Mazumder, A. Zadeh, and L.-P. Morency, "Context-Dependent Sentiment Analysis in User-Generated Videos," *Association for Computational Linguistics*, pp. 873–883, 2017.
- [8] M. Chen, S. Wang, P. P. Liang, T. Baltrušaitis, A. Zadeh, and L.-P. Morency, "Multimodal sentiment analysis with word-level fusion and reinforcement learning," in *Proceedings of the 19th ACM International Conference on Multimodal Interaction - ICMI 2017*. New York, New York, USA: ACM Press, 2017, pp. 163–171.
- [9] Q. Jin, C. Li, S. Chen, and H. Wu, "Speech emotion recognition with acoustic and lexical features," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 4 2015, pp. 4749–4753. [Online]. Available: <http://ieeexplore.ieee.org/document/7178872/>
- [10] Y. Gu, K. Yang, S. Fu, S. Chen, X. Li, and I. Marsic, "Multimodal Affective Analysis Using Hierarchical Attention Strategy with Word-Level Alignment," *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, pp. 2225–2235, 2018. [Online]. Available: <https://aclanthology.info/papers/P18-1207/p18-1207>
- [11] —, "Hybrid Attention based Multimodal Network for Spoken Language Classification," *ACL*, pp. 2379–2390, 2018.
- [12] Z. Aldeneh, S. Khorram, D. Dimitriadis, and E. M. Provost, "Pooling acoustic and lexical features for the prediction of valence," in *Proceedings of the 19th ACM International Conference on Multimodal Interaction - ICMI 2017*. New York, New York, USA: ACM Press, 2017, pp. 68–72. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=3136755.3136760>
- [13] M. Neumann and N. T. Vu, "Attentive Convolutional Neural Network based Speech Emotion Recognition: A Study on the Impact of Input Features, Signal Length, and Acted Speech," *Interspeech*, pp. 1263–1267, 2017.
- [14] C. Etienne, G. Fidanza, A. Petrovskii, L. Devillers, and B. Schmauch, "Speech Emotion Recognition with Data Augmentation and Layer-wise Learning Rate Adjustment," 2018. [Online]. Available: <https://arxiv.org/pdf/1802.05630.pdf>
- [15] B. W. Schuller, S. Steidl, A. Batliner, P. B. Marschik, H. Baumeister, F. Dong, S. Hantke, F. B. Pokorny, E.-M. Rathner, K. D. Bartl-Pokorny, C. Einspieler, D. Zhang, A. Baird, S. Amiriparian, K. Qian, Z. Ren, M. Schmitt, P. Tzirakis, and S. Zafeiriou, "The INTERSPEECH 2018 Computational Paralinguistics Challenge: Atypical & Self-Assessed Affect, Crying and Heart Beats," in *Interspeech*, 2018.
- [16] J. Gideon, S. Khorram, Z. Aldeneh, D. Dimitriadis, E. M. Provost, and A. Arbor, "Progressive Neural Networks for Transfer Learning in Emotion Recognition," *Interspeech*, pp. 1098–1102, 2017.
- [17] H. M. Fayek, M. Lech, and L. Cavedon, "On the Correlation and Transferability of Features between Automatic Speech Recognition and Speech Emotion Recognition," *Interspeech*, pp. 3618–362, 2016.
- [18] E. Lakomkin, C. Weber, S. Magg, and S. Wermter, "Reusing Neural Speech Representations for Auditory Emotion Recognition," *Proceedings of the Eighth International Joint Conference on Natural Language Processing*, vol. 1, pp. 423–430, 2017.
- [19] N. Blanchard, D. Moreira, A. Bharati, and W. Scheirer, "Getting the subtext without the text: Scalable multimodal sentiment classification from visual and acoustic modalities," *Proceedings of Grand Challenge and Workshop on Human Multimodal Language (Challenge-HML)*, pp. 1–10, 2018. [Online]. Available: <https://aclanthology.info/papers/W18-3301/w18-3301>
- [20] B. Krause, L. Lu, I. Murray, and S. Renals, "Multiplicative LSTM for sequence modelling," 9 2016. [Online]. Available: <http://arxiv.org/abs/1609.07959>
- [21] A. Radford, R. Jozefowicz, and I. Sutskever, "Learning to Generate Reviews and Discovering Sentiment," 4 2017. [Online]. Available: <http://arxiv.org/abs/1704.01444>
- [22] F. Eyben, F. Weninger, F. Gross, and B. Schuller, "Recent developments in openSMILE, the Munich open-source multimedia feature extractor," in *Proceedings of the 21st ACM international conference on Multimedia - MM '13*. New York, New York, USA: ACM Press, 2013, pp. 835–838.
- [23] F. Weninger, F. Eyben, B. W. Schuller, M. Mortillaro, and K. R. Scherer, "On the Acoustics of Emotion in Audio: What Speech, Music, and Sound have in Common." *Frontiers in psychology*, vol. 4, p. 292, 2013.
- [24] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 11 1997. [Online]. Available: <http://www.mitpressjournals.org/doi/10.1162/neco.1997.9.8.1735>
- [25] D. Amodei, R. Anubhai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, J. Chen, M. Chrzanowski, A. Coates, and et al, "Deep Speech 2: End-to-End Speech Recognition in English and Mandarin," *Proceedings of The 33rd International Conference on Machine Learning*, vol. 48, pp. 173–182, 2016. [Online]. Available: <https://arxiv.org/pdf/1512.02595.pdf>
- [26] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification," in *Proceedings of the 23rd International Conference on Machine Learning - ICML '06*. New York, New York, USA: ACM Press, 2006, pp. 369–376.
- [27] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh, S. Sengupta, A. Coates, and A. Y. Ng, "Deep Speech: Scaling up end-to-end speech recognition," *CoRR*, vol. abs/1412.5567, 2014.
- [28] Y. Zhou, C. Xiong, and R. Socher, "Improved Regularization Techniques for End-to-End Speech Recognition," *CoRR*, vol. abs/1712.07108, 2017. [Online]. Available: <http://arxiv.org/abs/1712.07108>
- [29] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2015.
- [30] F. Hernandez, V. Nguyen, S. Ghannay, N. Tomashenko, and Y. Estève, "TED-LIUM 3: twice as much data and corpus repartition for experiments on speaker adaptation," 5 2018. [Online]. Available: <http://arxiv.org/abs/1805.04699>
- [31] P. P. Liang, Z. Liu, A. Zadeh, and L.-P. Morency, "Multimodal Language Analysis with Recurrent Multistage Fusion," 8 2018. [Online]. Available: <http://arxiv.org/abs/1808.03920>
- [32] A. Zadeh, R. Zellers, E. Pincus, and L.-P. Morency, "MOSI: Multimodal Corpus of Sentiment Intensity and Subjectivity Analysis in Online Opinion Videos," *IEEE Intelligent Systems*, 2016.
- [33] J. Bauer, J. Davila-Chacon, E. Strahl, and S. Wermter, "Smoke and mirrors Virtual realities for sensor fusion experiments in biomimetic robotics," in *2012 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI)*. IEEE, 9 2012, pp. 114–119. [Online]. Available: <http://ieeexplore.ieee.org/document/6343022/>
- [34] J. Dávila-Chacón, J. Twiefel, J. Liu, and S. Wermter, "Improving Humanoid Robot Speech Recognition with Sound Source Localisation." Springer, Cham, 2014, pp. 619–626.
- [35] Y. Kim, "Convolutional Neural Networks for Sentence Classification," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Stroudsburg, PA, USA: Association for Computational Linguistics, 2014, pp. 1746–1751.