56

57

58

Improved Speech Emotion Recognition using Transfer Learning and Spectrogram Augmentation

Sarala Padi* sarala.padi@nist.gov ITL, NIST Gaithersburg, MD, USA

Ram D. Sriram ram.sriram@nist.gov ITL, NIST Gaithersburg, MD, USA Seyed Omid Sadjadi* omid.sadjadi@nist.gov ITL, NIST Gaithersburg, MD, USA 61

65

66

67 68

69

70

71

72 73

74

75

80

81

82

83

84

86

87

88

89

90

93

94

95

96

97

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

Dinesh Manocha dmanocha@umd.edu University of Maryland College Park, MD, USA

ABSTRACT

Automatic speech emotion recognition (SER) is a challenging task that plays a crucial role in natural human-computer interaction. One of the main challenges in SER is data scarcity, i.e., insufficient amounts of carefully labeled data to build and fully explore complex deep learning models for emotion classification. This paper aims to address this challenge using a transfer learning strategy combined with spectrogram augmentation. Specifically, we propose a transfer learning approach that leverages a pre-trained residual network (ResNet) model including a statistics pooling layer from speaker recognition trained using large amounts of speaker-labeled data. The statistics pooling layer enables the model to efficiently process variable-length input, thereby eliminating the need for sequence truncation which is commonly used in SER systems. In addition, we adopt a spectrogram augmentation technique to generate additional training data samples by applying random time-frequency masks to log-mel spectrograms to mitigate overfitting and improve the generalization of emotion recognition models. We evaluate the effectiveness of our proposed approach on the interactive emotional dyadic motion capture (IEMOCAP) dataset. Experimental results indicate that the transfer learning and spectrogram augmentation approaches improve the SER performance, and when combined achieve state-of-the-art results.

CCS CONCEPTS

• Computing methodologies; • Artificial intelligence; • Machine learning;

KEYWORDS

Attentive pooling, IEMOCAP, ResNet, spectrogram augmentation, speech emotion recognition (SER), transfer learning

ACM acknowledges that this contribution was authored or co-authored by an employee, contractor, or affiliate of the United States government. As such, the United States government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for government purposes only.

ICMI '21, October 18–22, 2021, Montréal, QC, Canada

© 2021 Association for Computing Machinery. ACM ISBN 978-1-4503-8481-0/21/10...\$15.00 https://doi.org/10.1145/3462244.3481003

ACM Reference Format:

Sarala Padi, Seyed Omid Sadjadi, Ram D. Sriram, and Dinesh Manocha. 2021. Improved Speech Emotion Recognition using Transfer Learning and Spectrogram Augmentation. In *Proceedings of the 2021 International Conference on Multimodal Interaction (ICMI '21), October 18–22, 2021, Montréal, QC, Canada.* ACM, New York, NY, USA, 8 pages. https://doi.org/10.1145/3462244.3481003

1 INTRODUCTION

Automatic emotion recognition plays a key role in human-computer interaction where it can enrich the next-generation AI with emotional intelligence by grasping the emotion from voice and words [32, 37]. The motivation behind developing algorithms to analyze emotions is to design computer interfaces that mimic and embed realistic emotions in synthetically generated responses [6]. Furthermore, research studies have shown that emotions play a critical role in the decision-making process for humans [6]. Hence, there is a growing demand to develop automatic systems that understand and recognize human emotions.

Humans express emotions in several ways, and speech is considered the most effective communication method to express feelings. For speech emotion recognition (SER), traditionally, machine learning (ML) models were developed using hand-crafted and engineered features such as mel-frequency cepstral coefficients (MFCC), Chroma-based features, pitch, energy, entropy, and zero-crossing rate [16, 21, 43], to mention a few. However, the performance of such ML models depends on the type and diversity of the features used. Although it remains unclear which features correlate most with various emotions, the research is still ongoing to explore additional features and new algorithms to model the dynamics of feature streams representing human emotions. On the other hand, the recent advancements in deep learning, along with the available computational capabilities, have enabled the research community to build end-to-end systems for SER. A big advantage of such systems is that they can directly learn the features from spectrograms or raw waveforms [12, 23, 36, 41, 45], thereby obviating the need for extracting a large set of hand-crafted features [13]. Recent studies have proposed the use of convolutional neural network (CNN) models combined with long short-term memory (LSTM) built on spectrograms and raw waveforms, showing improved SER performance [19, 23, 24, 26, 35, 36, 46]. However, building such complex systems requires large amounts of labeled training data. Also, the insufficient labeled training data can potentially make the models

1

^{*}Both authors contributed equally to this research.

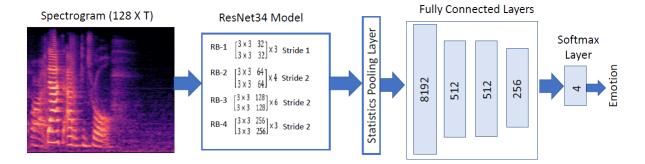


Figure 1: Block diagram of the proposed SER system. T denotes the number of frames.

overfit to specific data conditions and domains, resulting in poor generalization on unseen data.

This paper addresses the insufficient data problem using a transfer learning approach combined with a spectrogram augmentation strategy. We re-purpose a residual network (ResNet) model [17] developed for speaker recognition using large amounts of speaker-labeled data and use it as a feature descriptor for SER. The model includes a statistics pooling layer that enables processing of variable length segments without a need for truncation. Also, we increase the training data size by generating more data samples using spectrogram augmentation [30]. We evaluate the effectiveness of our proposed system on the interactive emotional dyadic motion capture (IEMOCAP) dataset [3].

2 RELATED WORK

Recently, neural network based modeling approaches along with different variations of attention mechanism (e.g., plain [18], local [24], and self [40]) have shown promise for SER. Among them, techniques such as bidirectional LSTMs (BLSTM) [10, 18, 24, 31, 42] and time-delay neural networks (TDNN) [44], which can effectively model relatively long contexts compared to their DNN counterparts, have been successfully applied for SER on the IEMOCAP. Nevertheless, as discussed previously, the lack of large amounts of carefully labeled data to build complex models for emotion classification remains a main challenge in SER [1]. To address this, two approaches are commonly used: data augmentation and transfer learning.

Data augmentation methods generate additional training data by perturbing, corrupting, mimicking, and masking the original data samples to enable the development of complex ML models. For example, [4, 29, 35] applied signal-based transformations such as speed perturbation, time-stretch, pitch shift, as well as added noise to original speech waveforms. One disadvantage of these approaches is that they require signal-level modifications, thereby increasing the computational complexity and storage requirements of the subsequent front-end processing. They can also lead to model overfitting due to potentially similar samples in the training set, while random balance can potentially remove useful information [4]. For example, in [9] a vocal tract length perturbation (VTLP) approach was explored for data augmentation along with a CNN model, and was found to result in a lower accuracy compared to a baseline model due to overfitting issues.

Since generative adversarial network (GAN) based models have demonstrated remarkable success in computer vision, several studies have recently incorporated this idea to address the data scarcity problem and to generate additional data samples for SER [4, 8]. For instance, [4] addressed the data imbalance using signal-based transformations and GAN based models for generating high-resolution spectrograms to train a VGG19 model for an emotion classification task and showed that GAN-generated spectrograms outperformed signal-based transformations. However, GAN generated features are strongly dependent on the data used during training and may not generalize to other datasets. Another challenge with GAN-based augmentation is that it is difficult to train and optimize.

Another effective way to address challenges related to data scarcity is transfer learning [2, 11, 14, 28, 49]. Transfer learning can leverage the information and knowledge learned from one related task and domain to another. Several recent studies have proposed transfer learning methods to improve SER performance and have shown these methods to outperform prior methods in recognizing emotions even for unseen scenarios, individuals, and conditions [15]. It has been shown that transfer learning can increase feature learning abilities, and that the transferred knowledge can further enhance the SER classification accuracy [7, 13, 22, 39]. To further improve the SER performance, transfer learned features have been used in combination with deep belief networks (DBN) [22], recurrent neural networks (RNN) [13], CNN [39], temporal convolutional network (TCN) [49], and sparse autoencoder [7]. However, transfer learning methods have not been fully explored and analyzed for emotion recognition. Particularly, it is unclear whether and how ML models trained for other data-rich speech applications such as speaker recognition would perform for SER.

3 PROPOSED SYSTEM

Figure 4 shows the block diagram of the proposed system for speech emotion recognition. We use an end-to-end system with a ResNet34 model [17] to perform emotion classification. ResNet models, originally developed for computer vision applications [17], have recently gained interest for speech applications such as speaker recognition [48]. The residual blocks introduced in ResNet models allow us to train much deeper models that are otherwise difficult, if not impossible, to train due to vanishing and exploding gradient problems. The ResNet models also allow the higher layers to learn the identity function so that higher-level features perform equally well on

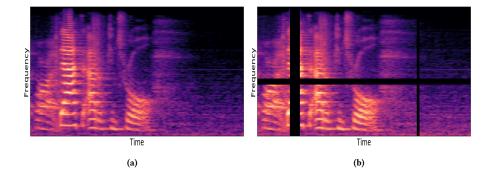


Figure 2: (a) original spectrogram, and (b) spectrogram modified using (multiple) masking blocks of consecutive time steps (vertical masks) and mel frequency channels (horizontal masks). The black horizontal and vertical stripes indicate the masked portions of the spectrogram.

unseen data compared to the lower layers of the model. In our proposed system, the convolutional layers in the model learn feature representations (feature maps) and reduce the spectral variations into compact representations, while the fully connected (FC) layers take the contextual features and generate predictions for emotion classification.

3.1 Input data

Although SER systems traditionally used a large set of low-level time- and frequency-domain features to capture and represent the various emotions in speech, in recent years many state-of-the-art SER systems use complex neural network models that learn directly from spectrograms, or even raw waveforms. Accordingly, in this study, we build and explore a ResNet based system using log-mel spectrograms as input features. We extract high-resolution spectrograms to enable the model to not only learn the spectral envelope structure, but also the coarse harmonic structure for the various emotions.

3.2 Transfer learning

As noted previously, transfer learning is a ML method where a model initially developed for one task or domain is re-purposed, partly or entirely, for a different but related task/domain. It has recently gained interest for SER [11]. In this study, we re-purpose a model initially developed for speaker recognition to serve as a feature descriptor for SER. More specifically, we first train a ResNet34 model on large amounts of speaker-labeled audio data. Then, we replace the FC layers of the pre-trained model with new randomly initialized FC layers. Finally, we re-train the new FC layers for an SER task on the IEMOCAP dataset.

3.3 Statistics pooling

As shown in Figure 4, the proposed system employs a statistics pooling layer [38] that aggregates the frame-level information over time and reduces the sequence of frames to a single vector by concatenating the mean and standard deviation computed over frames. Accordingly, the convolutional layers in the ResNet model work at the frame-level, while the FC layers work at the segment-level. This enables the system to efficiently model variable-length

Table 1: Parameter settings for the conservative and aggressive augmentation policies. Here, N_f and N_t denote the number of frequency and time masks applied.

Augmentation Policy	F	W	p	N_f	N_t
None	0	0	_	_	_
Conservative	15	50	0.2	2	2
Aggressive	27	70	0.2	2	2

sequences of frames, thereby eliminating the need for truncating the sequence of frames to a pre-specified length to match that of the segments used during training. The sequence-truncation approach, which is commonly adopted in neural network based SER systems, can have a deleterious impact on SER performance as potentially informative frames are dropped out from the input. It is worth noting here that the statistics pooling can be viewed as an attention mechanism with equal weights for all frames, which also appends second order statistics (i.e., standard deviation) to capture long-term temporal variability over the duration of segments.

3.4 Spectrogram augmentation

Currently, the majority of the features and methods for SER are adapted from speech recognition, speaker recognition, or speech synthesis fields [20]. There has been recent success in applying a computationally efficient data augmentation strategy, termed spectrogram augmentation, for speech recognition tasks [30]. The spectrogram augmentation technique generates additional training data samples by applying random time-frequency masks to spectrograms to mitigate the overfitting issue and improve the generalization of speech recognition models. Motivated by promising results seen with the spectrogram augmentation in the speech recognition field, we augment the training data using spectro-temporally modified versions of the original spectrograms (see Figure 2). Because the time-frequency masks are applied directly to spectrograms, the augmentation can be conveniently applied on-the-fly, eliminating the necessity to create and store new data files as commonly done in many augmentation approaches for speech applications.

411

412

413

414

415

418

419

420

421

422

424

425

426

427

431

432

433

434

435

437

438

439

440

441

444

445

446

447

448

449

451

452

453

454

455

457

458

459

460

461

463

464

Similar to the approach taken in [30], we consider two policies to systematically apply spectrogram augmentation for SER, namely conservative and aggressive. The frequency masking is applied over f consecutive frequency channels in the range $[f_0, f_0 + f]$, where f is sampled from a uniform distribution [0, F] and f_0 is sampled from [0, v - f]. Here, F and v denote the maximum width of frequency masks and the total number of frequency channels, respectively. The time masking, on the other hand, is applied over t consecutive frames in the range $[t_0, t_0 + t)$, where t is selected from a uniform distribution [0, W] and t_0 is sampled from [0, T - t]. Similarly, W and T denote the maximum width of time masks and the number of time frames, respectively. An upper bound is also applied on the width of the time masks such that $W = \min(W, pT)$, i.e., the width of a mask cannot be longer than p times the number of time frames. This is to ensure sufficient speech content after masking, in particular for shorter segments. Table 1 summarizes the various parameters for the two spectrogram augmentation policies used in this paper.

4 EXPERIMENTS

4.1 Dataset

349

350

351

352

353

354

355

356

357

358

361

362

363

364

365 366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

388

389

390

391

392

393

394

395

396

397

400

401

402

403

404

405

406

We evaluate the effectiveness of the proposed SER system on the IEMOCAP dataset [3], which contains improvised and scripted multimodal dyadic conversations between actors of opposite gender. It consists of 12 hours of speech data from 10 subjects, presegmented into short cuts that were judged by three annotators to generate emotion labels. It includes nine categorical emotions and 3-dimensional labels. In our experiments, we only consider the speech segments for which at least two annotators agree on the emotion label. In an attempt to replicate the experimental protocols used in a number of prior studies, we conduct three experiments on the full dataset (i.e., the combined improvised and scripted portions): Exp 1, using four categorical emotions: "angry", "happy", "neutral", "sad"; Exp 2, using the same categories as in Exp 1, but replacing the "happy" category with "excited"; Exp 3, by merging the "happy' and "excited" categories from Exp 1 and Exp 2. The total number of examples used for Exp 1 is 4490 and the number of examples per category is 1103, 595, 1708, and 1084, respectively. The number of examples in the "excited" category is 1041, making the total number of examples in the merged category (i.e., Exp 3) 1636. Table 2 summarizes the data statistics in the IEMOCAP dataset for the three experimental setups considered in this study.

The IEMOCAP dataset comprises five sessions, and the speakers in the sessions are non-overlapping. Therefore, there are 10 speakers in the dataset, i.e., 5 female and 5 male speakers. To conduct the experiments in a speaker-independent fashion, we use a leave-one-session-out (LOSO) cross-validation strategy, which results in 5 different train-test splits/folds. For each fold, we use the data from 4 sessions for training and the remaining one session for model evaluation. Since the dataset is multi-label and imbalanced, in addition to the overall accuracy, termed weighted accuracy (WA), we report the average recall over the different emotion categories, termed unweighted accuracy (UA), to present our findings. Additionally, to understand and visualize the performance of the proposed system within and across the various emotion categories, we compute and report confusion matrices for the three experiments. Note that for

Table 2: Data statistics for the various emotion classes in the IEMOCAP for the three experimental setups considered in this study. Both the improvised and scripted portions of the IEMOCAP dataset are used in our experiments.

Experiment	Emotion	#segments
Exp 1	Angry	1103
	Нарру	595
	Neutral	1708
	Sad	1084
	Total	4490
Exp 2	Angry	1103
	Excited	1041
	Neutral	1708
	Sad	1084
	Total	4936
Exp 3	Angry	1103
	Excited+Happy	1636
	Neutral	1708
	Sad	1084
	Total	5531

each experiment, we compute the average of performance metrics over the five training-test splits as the final result.

4.2 Setup and configuration

For speech parameterization, high resolution 128-dimensional logmel spectrograms are extracted from 25 ms frames at a 100 Hz frame rate (i.e., every 10 ms). For feature normalization, a segment level mean and variance normalization is applied¹. Note that this is not ideal as typically the normalization is applied at the recording/conversation level. We have found that normalizing the segments using statistics computed at the conversation level significantly improves the SER performance on the IEMOCAP. Nevertheless, this violates the independence assumption for the speech segments, hence it is not considered in this study. The front-end processing, including feature extraction and feature normalization, is performed using the NIST speaker and language recognition evaluation (SLRE) [33, 34] toolkit. While training the model, we select T-frame chunks using random offsets over original speech segments where T is randomly sampled from the set $\{150, 200, 250, 300\}$ for each batch. For speech segments shorter than T frames, signal padding is applied. On the other hand, while evaluating the model, we feed the entire duration of the test segments because the statistics pooling layer enables the model to consume variable-length inputs.

As noted previously, the proposed end-to-end SER system uses a pre-trained ResNet34 model built on a speaker recognition task. We train the ResNet34 model on millions of speech samples from more than 7000 speakers available in the VoxCeleb corpus [25]. To build the speaker recognition model, we apply the same frontend processing described above to extract high-resolution log-mel

¹No voice activity detection (VAD) is applied prior to feature normalization because it was found to be detrimental to SER performance on the IEMOCAP. We hypothesize that this is because the silence gaps in within and between utterances might be relevant in terms of speakers' emotional state.

E-Excited, H+E: Happy and Excited merged. Blanks (-) indicate unreported values.

485

486

491

492

493

494

495

496

497

498

499

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

465

466

467

469

470

471

472

473

523

524

525

527

528

529

530

531

535

536

537

538

580

Table 3: Performance comparison of our proposed approach with prior methods that use the LOSO strategy for experiments on the full IEMOCAP dataset (i.e., both the improvised and scripted portions). Abbreviations: A-Angry, H- Happy, N-Neutral,

Experiment (emotion classes)	Approach	UA [%]	WA [%]
	BLSTM+attention [18]	49.96	59.33
Exp 1 (A, H, S, N)	Transformer+self-attention [40]		59.43
	BLSTM+local attention [24]	58.8	63.5
	BLSTM+attention [31]	59.6	62.5
	Proposed	61.61	66.02
Exp 2 (A, E, S, N)	CTC-BLSTM [5]	54	-
	BLSTM+attention [42]	55.65	-
	Transformer+self-attention [40]	64.79	64.33
	Proposed	65.56	65.62
Exp 3 (A, H+E, S, N)	BLSTM+transfer learnig [13]	51.86	50.47
	VGG19+GAN augmentation [4]	54.6	-
	CNN+attention+multi-task learning [26]		56.10
	BLSTM+self-attention [10]		55.7
	CNN+attention+multi-task/transfer learning [27]	59.54	-
	ResTDNN+self-attention [44]	61.32	60.64
	Proposed	64.14	63.61

spectrograms from VoxCeleb data. We conduct experiments using models with and without transfer learning and spectrogram augmentation. For each original speech segment, we generate and augment two spectro-temporally modified versions according to the augmentation policies defined in Table 1. This is applied for both speaker and emotion recognition systems during training. To study the impact of the statistics pooling layer, we also evaluate these models with and without this layer. For all the experiments, we use a categorical cross-entropy loss as the objective function to train the models. The number of channels in the first block of the ResNet model is set to 32. The model is trained using Pytorch² and the stochastic gradient descent (SGD) optimizer with momentum (0.9), an initial learning rate of 10^{-2} , and a batch size of 32. The learning rate remains constant for the first 8 epochs, after which it is halved every other epoch. We use parametric rectified linear unit (PReLU) activation functions in all layers (except for the output), and utilize a layer-wise batch normalization to accelerate the training process and improve the generalization properties of the

RESULTS

Table 3 presents the performance comparison of our proposed system with several prior approaches for the three experimental setups (i.e., Exp 1, 2, and 3) described in Section 4. The results are obtained using the combined system that utilizes the ResNet model with the statistics pooling layer trained using the transfer learning and spectrogram augmentation approaches described in Section 3. All studies referenced in the table adopt the LOSO strategy to conduct experiments on both the improvised and scripted portions of the

To visualize the performance of the proposed system within and across the different emotion categories, confusion matrices for the three experimental setups are shown in Figure 3. It is observed from Figure 3(a) that the system confuses the "happy" class (H) with the "neutral" class (N) quite often, while performing the best on the "angry" class (A). This is consistent with observations reported in other studies on IEMOCAP [26, 47]. Our informal listening experiments confirm that the "happy" and "neutral" classes are indeed confusable emotion pairs in the IEMOCAP dataset. The system performance balance is improved in Figure 3(b) where we replace the less pronounced "happy" category with the "excited" category (E). Combining the "happy" and "excited" categories in Exp 3 further improves the performance balance across the various emotions, at the expense of increasing the confusion between the "angry" (A) and "excited" plus "happy" (H+E) categories.

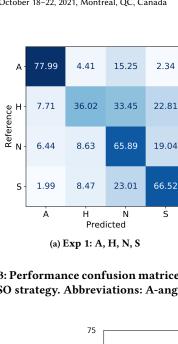
To investigate and quantify the contribution of the various system components proposed in this study for improved SER, we further conduct ablation experiments to measure the system performance with and without the transfer learning, the spectrogram augmentation, and the statistics pooling layer. For these ablation

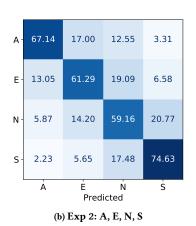
 $^3\mathrm{There}$ are other related studies in the literature that only use the improvised portion

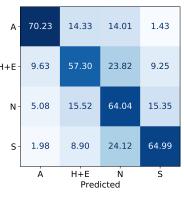
of the IEMOCAP dataset [9, 36, 49]. On the other hand, in our experiments, we use

IEMOCAP dataset³. It can be seen from the table that the proposed system consistently provides competitive performance across the three experiments, achieving state-of-the-art results. In the case of Exp 2, the proposed system outperforms a system that uses 384 engineered features [40], while for the other two experiments, our proposed system outperforms systems that use a large set of engineered features (e.g., [24] and [31]).

both the improvised and scripted portions of the IEMOCAP, which is approximately twice the size of the improvised portion alone. Because the experimental setups and the amount of data used for model training and evaluation in those studies are different than ours, we have not included them in Table 3 for comparison. The SER performance on the improvised portion is known to be better than that on the full dataset (e.g., see [13, 26, 31, 40, 42]).







(c) Exp 3: A, H+E, N, S

Figure 3: Performance confusion matrices of the proposed SER system for the three experiments conducted in this study using the LOSO strategy. Abbreviations: A-angry, E-excited, H-happy, N-neutral, S-Sad, and H+E- Happy and Excited merged.

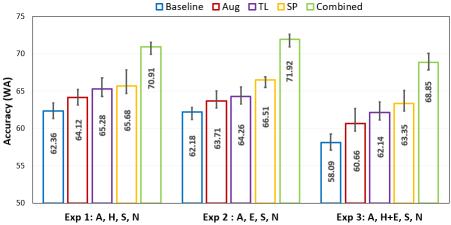


Figure 4: Performance (WA) of the proposed approach with and without transfer learning (TL), spectrogram augmentation (Aug), and statistics pooling (SP). All results are obtained using a 5-fold cross-validation. The height of each bar represents the average accuracy computed over 5 runs, while the error bars denote the standard deviations over the 5 runs. Abbreviations: A-Angry, H-Happy, N-Neutral, S-Sad, E-Excited, H+E: Happy and Excited merged

experiments, we employ a 5-fold cross-validation (CV) strategy, where we use 80% of the data for training and 20% for testing the system. This process is repeated 5 times to reduce possible partition-dependencies. Figure 4 shows the average overall classification accuracy (WA) computed across 5 folds (or 5 runs). The height of the bars represents the average accuracy, and the error bars denote the standard deviations computed over the 5 runs. It is observed that the proposed components, both individually and in combination, consistently provide performance gains across the three experimental setups (i.e., Exp 1, 2 and 3). The statistics pooling approach seems to have the greatest impact on performance, followed by the transfer learning and spectrogram augmentation methods. Furthermore, the model that combines all the system components not only consistently achieves the best performance, but

also relatively smaller variation across the 5 runs as evidenced by the error bars.

CONCLUSIONS

In this paper, we explored a transfer learning approach along with a spectrogram augmentation strategy to improve the SER performance. Specifically, we re-purposed a pre-trained ResNet model from speaker recognition that was trained using large amounts of speaker-labeled data. The convolutional layers of the ResNet model were used to extract features from high-resolution log-mel spectrograms. In addition, we adopted a spectrogram augmentation technique to generate additional training data samples by applying random time-frequency masks to log-mel spectrograms to mitigate overfitting and improve the generalization of emotion recognition

756

757

759

760

761

762

763

766

767

768

769

770

772

773

774

775

779

780

781

782

783

785

786

787

788

789

790

792

793

794

795

796

797

799

800

801

802

806

807

808

809

810

811

812

models. We evaluated the proposed system using three different experimental settings and compared the performance against that of several prior studies. The proposed system consistently provided competitive performance across the three experimental setups, achieving state-of-the-art results on two settings. The state-of-theart results were achieved without the use of engineered features. It was also shown that incorporating the statistics pooling layer to accommodate variable-length audio segments improved the emotion recognition performance. Results from this study suggest that, for practical applications, simplified front-ends with only spectrograms can be as effective for SER, and that models trained for data-rich speech applications such as speaker recognition can be re-purposed using transfer learning to improve the SER performance under data scarcity constraints. In the future, to further enhance the emotion recognition accuracy, we will extend our work along these lines by exploring more data augmentation methods, incorporating other transfer learning paradigms, and evaluating the proposed system across different datasets.

7 ACKNOWLEDGEMENT

Experiments and analyses were performed, in part, on the NIST Enki HPC cluster.

8 DISCLAIMER

697

698

699

700

701

702

703

704

705

708

709

710

711

712

713

714

715

716

718

719

721

723

724

725

726

727

728

729

730

731

732

734

735

736

737

738

739

740

741

742

743

744

745

747

748

749

750

751

752

753

754

The views and conclusions presented in this paper are those of the authors and should not be interpreted as the official findings, either expressed or implied, of NIST or the U.S. Government.

REFERENCES

- Samuel Albanie, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. 2018.
 Emotion recognition in speech using cross-modal transfer in the wild. In *Proc. ACM ICM*. 292–301.
- [2] George Boateng and Tobias Kowatsch. 2020. Speech emotion recognition among elderly individuals using multimodal fusion and transfer learning. In Proc. ACM ICMI. 12–16.
- [3] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. *Language resources and evaluation* 42, 4 (2008), 335.
- [4] Aggelina Chatziagapi, Georgios Paraskevopoulos, Dimitris Sgouropoulos, Georgios Pantazopoulos, Malvina Nikandrou, Theodoros Giannakopoulos, Athanasios Katsamanis, Alexandros Potamianos, and Shrikanth Narayanan. 2019. Data augmentation using GANs for speech emotion recognition. In Proc. INTERSPEECH. 171–175
- [5] Vladimir Chernykh and Pavel Prikhodko. 2017. Emotion recognition from speech with recurrent neural networks. arXiv preprint arXiv:1701.08071 (2017).
- [6] Roddy Cowie, Ellen Douglas-Cowie, Nicolas Tsapatsoulis, George Votsis, Stefanos Kollias, Winfried Fellenz, and John G Taylor. 2001. Emotion recognition in humancomputer interaction. *IEEE Signal Processing Magazine* 18, 1 (2001), 32–80.
- [7] Jun Deng, Zixing Zhang, Erik Marchi, and Björn Schuller. 2013. Sparse autoencoder-based feature transfer learning for speech emotion recognition. In Proc. Humaine Association Conference on Affective Computing and Intelligent Interaction. 511–516.
- [8] Sefik Emre Eskimez, Dimitrios Dimitriadis, Robert Gmyr, and Kenichi Kumanati. 2020. GAN-based data generation for speech emotion recognition. Proc. INTERSPEECH (2020), 3446–3450.
- [9] Caroline Etienne, Guillaume Fidanza, Andrei Petrovskii, Laurence Devillers, and Benoit Schmauch. 2018. CNN+LSTM architecture for speech emotion recognition with data augmentation. arXiv preprint arXiv:1802.05630 (2018).
- [10] Han Feng Feng, Sei Uno, and Tatsuya Kawahara. 2020. End-to-End speech emotion recognition combined with acoustic-to-word ASR model. In Proc. IN-TERSPEECH 501-505.
- [11] Kexin Feng and Theodora Chaspari. 2020. A review of generalizable transfer learning in automatic emotion recognition. Frontiers in Computer Science 2, 9 (2020).

- [12] Mengna Gao, Jing Dong, Dongsheng Zhou, Qiang Zhang, and Deyun Yang. 2019. End-to-end speech emotion recognition based on one-dimensional convolutional neural network. In Proc. ACM ICIAI. 78–82.
- [13] Sayan Ghosh, Eugene Laksana, Louis-Philippe Morency, and Stefan Scherer. 2016. Representation Learning for Speech Emotion Recognition. In Proc. INTERSPEECH. 3603–3607.
- [14] John Gideon, Soheil Khorram, Zakaria Aldeneh, Dimitrios Dimitriadis, and Emily Mower Provost. 2017. Progressive neural networks for transfer learning in emotion recognition. In Proc. INTERSPEECH.
- [15] John Gideon, Melvin McInnis, and Emily Mower Provost. 2019. Improving cross-corpus speech emotion recognition with adversarial discriminative domain generalization (ADDoG). IEEE Trans. Affective Computing (2019).
- [16] Kun Han, Dong Yu, and Ivan Tashev. 2014. Speech emotion recognition using deep neural network and extreme learning machine. In Proc. INTERSPEECH. 223–227.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition. 770–778.
- [18] Che-Wei Huang and Shri Narayanan. 2016. Attention assisted discovery of sub-utterance structure in speech emotion recognition. In Proc. INTERSPEECH. 1387–1391
- [19] Gil Keren and Björn Schuller. 2016. Convolutional RNN: an enhanced model for extracting features from sequential data. In Proc. IEEE IJCNN. 3412–3419.
- [20] Shashidhar G Koolagudi and K Sreenivasa Rao. 2012. Emotion recognition from speech: a review. International Journal of Speech Technology 15, 2 (2012), 99–117.
- [21] Oh-Wook Kwon, Kwokleung Chan, Jiucang Hao, and Te-Won Lee. 2003. Emotion recognition by speech signals. In Proc. INTERSPEECH. 125–128.
- [22] Siddique Latif, Rajib Rana, Shahzad Younis, Junaid Qadir, and Julien Epps. 2018. Transfer learning for improving speech emotion classification accuracy. In Proc. INTERSPEECH. 257–261.
- [23] Xi Ma, Zhiyong Wu, Jia Jia, Mingxing Xu, Helen Meng, and Lianhong Cai. 2018. Emotion recognition from variable-length speech segments using deep learning on spectrograms. In Proc. INTERSPEECH. 3683–3687.
- [24] Seyedmahdad Mirsamadi, Emad Barsoum, and Cha Zhang. 2017. Automatic speech emotion recognition using recurrent neural networks with local attention. In Proc. IEEE ICASSP. 2227–2231.
- [25] Arsha Nagrani, Joon Son Chung, Weidi Xie, and Andrew Zisserman. 2020. Voxceleb: Large-scale speaker verification in the wild. Computer Speech & Language 60 (2020), 1–15.
- [26] Michael Neumann and Ngoc Thang Vu. 2017. Attentive convolutional neural network based speech emotion recognition: A study on the impact of input features, signal length, and acted speech. In Proc. INTERSPEECH. 1263–1267.
- [27] Michael Neumann and Ngoc Thang Vu. 2019. Improving Speech Emotion Recognition with Unsupervised Representation Learning on Unlabeled Speech. In Proc. IEEE ICASSP. 7390–7394.
- [28] Sandra Ottl, Shahin Amiriparian, Maurice Gerczuk, Vincent Karas, and Björn Schuller. 2020. Group-level speech emotion recognition utilising deep spectrum features. In Proc. ACM ICMI. 821–826.
- [29] Raghavendra Pappagari, Tianzi Wang, Jesus Villalba, Nanxin Chen, and Najim Dehak. 2020. X-Vectors meet emotions: A study on dependencies between emotion and speaker recognition. In *Proc. IEEE ICASSP*. 7169–7173.
- [30] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. 2019. SpecAugment: A simple data augmentation method for automatic speech recognition. In Proc. INTERSPEECH. 2613–2617.
- [31] Gaetan Ramet, Philip N Garner, Michael Baeriswyl, and Alexandros Lazaridis. 2018. Context-aware attention mechanism for speech emotion recognition. In Proc. IEEE SLT Workshop. 126–131.
- [32] Harper Richard, R Tom, R Yvonne, and S Abigail. 2008. Being Human: Human-Computer Interaction in The Year 2020. Report, Microsoft Corporation.
- [33] Seyed Omid Sadjadi, Craig Greenberg, Elliot Singer, Douglas Reynolds, Lisa Mason, and Jaime Hernandez-Cordero. 2020. The 2019 NIST audio-visual speaker recognition evaluation. In Proc. Speaker Odyssey Workshop. 259–265.
- [34] Seyed Omid Sadjadi, Timothee Kheyrkhah, Audrey Tong, Craig Greenberg, Elliot Singer, Douglas Reynolds, Lisa Mason, and Jaime Hernandez-Cordero. 2018. The 2017 NIST language recognition evaluation. In Proc. Speaker Odyssey Workshop. 82–89.
- [35] Mousmita Sarma, Pegah Ghahremani, Daniel Povey, Nagendra Kumar Goel, Kandarpa Kumar Sarma, and Najim Dehak. 2018. Emotion identification from raw speech signals using DNNs. In Proc. INTERSPEECH. 3097–3101.
- [36] Aharon Satt, Shai Rozenberg, and Ron Hoory. 2017. Efficient emotion recognition from speech using deep learning on spectrograms. In Proc. INTERSPEECH. 1089– 1009.
- [37] Björn W Schuller. 2018. Speech emotion recognition: Two decades in a nutshell, benchmarks, and ongoing trends. Commun. ACM 61, 5 (2018), 90–99.
- [38] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur. 2018. X-vectors: Robust DNN embeddings for speaker recognition. In Proc. IEEE ICASSP. 5329–5333.

- [39] Peng Song, Yun Jin, Li Zhao, and Minghai Xin. 2014. Speech emotion recognition using transfer learning. IEICE Trans. Information and Systems 97, 9 (2014), 2530– 2532.
- [40] Lorenzo Tarantino, Philip N Garner, and Alexandros Lazaridis. 2019. Selfattention for speech emotion recognition. In Proc. INTERSPEECH. 2578–2582.
- [41] George Trigeorgis, Fabien Ringeval, Raymond Brueckner, Erik Marchi, Mihalis A Nicolaou, Björn Schuller, and Stefanos Zafeiriou. 2016. Adieu features? end-toend speech emotion recognition using a deep convolutional recurrent network. In Proc. IEEE ICASSP. 5200–5204.
- [42] Samarth Tripathi, Tripathi Sarthak, and Homayoon Beigi. 2018. Multi-modal emotion recognition on IEMOCAP dataset using deep learning. arXiv preprint arXiv:1804.05788 (2018).
- [43] Dimitrios Ververidis and Constantine Kotropoulos. 2006. Emotional speech recognition: Resources, features, and methods. Speech Communication 48, 9 (2006), 1162–1181.
- [44] Wen Wu, Chao Zhang, and Philip C. Woodland. 2021. Emotion recognition by fusing time synchronous and time asynchronous representations. In Proc. IEEE

- ICASSP, 6269-6273
- [45] Zixiaofan Yang and Julia Hirschberg. 2018. Predicting arousal and valence from waveforms and spectrograms using deep neural networks. In Proc. INTERSPEECH. 3092–3096.
- [46] Promod Yenigalla, Abhay Kumar, Suraj Tripathi, Chirag Singh, Sibsambhu Kar, and Jithendra Vepa. 2018. Speech emotion recognition using spectrogram & phoneme embedding. In Proc. INTERSPEECH. 3688–3692.
- [47] Seunghyun Yoon, Seokhyun Byun, and Kyomin Jung. 2018. Multimodal speech emotion recognition using audio and text. In Proc. IEEE SLT Workshop. IEEE, 112–118
- [48] Hossein Zeinali, Shuai Wang, Anna Silnova, Pavel Matějka, and Oldřich Plchot. 2019. BUT system description to VoxCeleb speaker recognition challenge 2019. arXiv preprint arXiv:1910.12592 (2019).
- [49] Ziping Zhao, Zhongtian Bao, Zixing Zhang, Nicholas Cummins, Shihuang Sun, Haishuai Wang, Jianhua Tao, and Björn W. Schuller. 2021. Self-attention transfer networks for speech emotion recognition. Virtual Reality & Intelligent Hardware 3, 1 (2021), 43–54.