



## 2017 Special Issue

## Evaluating deep learning architectures for Speech Emotion Recognition

Haytham M. Fayek<sup>a,\*</sup>, Margaret Lech<sup>a</sup>, Lawrence Cavedon<sup>b</sup><sup>a</sup> School of Engineering, RMIT University, Melbourne VIC 3001, Australia<sup>b</sup> School of Science, RMIT University, Melbourne VIC 3001, Australia

## ARTICLE INFO

## Article history:

Available online xxxx

## Keywords:

Affective computing  
Deep learning  
Emotion recognition  
Neural networks  
Speech recognition

## ABSTRACT

Speech Emotion Recognition (SER) can be regarded as a static or dynamic classification problem, which makes SER an excellent test bed for investigating and comparing various deep learning architectures. We describe a frame-based formulation to SER that relies on minimal speech processing and end-to-end deep learning to model intra-utterance dynamics. We use the proposed SER system to empirically explore feed-forward and recurrent neural network architectures and their variants. Experiments conducted illuminate the advantages and limitations of these architectures in paralinguistic speech recognition and emotion recognition in particular. As a result of our exploration, we report state-of-the-art results on the IEMOCAP database for speaker-independent SER and present quantitative and qualitative assessments of the models' performances.

© 2017 Elsevier Ltd. All rights reserved.

## 1. Introduction

In recent years, deep learning in neural networks has achieved tremendous success in various domains that led to multiple deep learning architectures emerging as effective models across numerous tasks. Feed-forward architectures such as Deep Neural Networks (DNNs) and Convolutional Neural Networks (ConvNets) have been particularly successful in image and video processing as well as speech recognition, while recurrent architectures such as Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) RNNs have been effective in speech recognition and natural language processing (LeCun, Bengio, & Hinton, 2015; Schmidhuber, 2015). These architectures process and model information in different ways and have their own advantages and limitations. For instance, ConvNets are able to deal with high-dimensional inputs and learn features that are invariant to small variations and distortions (Krizhevsky, Sutskever, & Hinton, 2012), whereas LSTM-RNNs are able to deal with variable length inputs and model sequential data with long range context (Graves, 2008).

In this paper, we investigate the application of end-to-end deep learning to Speech Emotion Recognition (SER) and critically explore how each of these architectures can be employed in this task.

SER can be regarded as a static or dynamic classification problem, which has motivated two popular formulations in the literature to the task (Verweridis & Kotropoulos, 2006): *turn-based processing* (also known as static modeling), which aims to recognize emotions from a complete utterance; or *frame-based processing* (also known as dynamic modeling), which aims to recognize emotions at the frame level. In either formulation, SER can be employed in stand-alone applications; e.g. emotion monitoring, or integrated into other systems for emotional awareness; e.g. integrating SER into Automatic Speech Recognition (ASR) to improve its capability in dealing with emotional speech (Cowie et al., 2001; Fayek, Lech, & Cavedon, 2016b; Fernandez, 2004). Frame-based processing is more robust since it does not rely on segmenting the input speech into utterances and can model intra-utterance emotion dynamics (Arias, Busso, & Yoma, 2013; Fayek, Lech, & Cavedon, 2015). However, empirical comparisons between frame-based processing and turn-based processing in prior work have demonstrated the superiority of the latter (Schuller, Vlasenko, Eyben, Rigoll, & Wendenmuth, 2009; Vlasenko, Schuller, Wendenmuth, & Rigoll, 2007).

Whether performing turn-based processing or frame-based processing, most of the research effort in the last decade has been devoted to selecting an optimal set of features (Schuller et al., 2010). Despite the effort, little success has been achieved in realizing such a set of features that performs consistently over different conditions and multiple data sets (Eyben, Scherer et al., 2015). Thus, researchers have resorted to brute-force high-dimensional

\* Corresponding author.

E-mail addresses: [haytham.fayek@ieee.org](mailto:haytham.fayek@ieee.org) (H.M. Fayek), [margaret.lech@rmit.edu.au](mailto:margaret.lech@rmit.edu.au) (M. Lech), [lawrence.cavedon@rmit.edu.au](mailto:lawrence.cavedon@rmit.edu.au) (L. Cavedon).

<http://dx.doi.org/10.1016/j.neunet.2017.02.013>

0893-6080/© 2017 Elsevier Ltd. All rights reserved.

Please cite this article in press as: Fayeek, H. M., et al., Evaluating deep learning architectures for Speech Emotion Recognition. *Neural Networks* (2017), <http://dx.doi.org/10.1016/j.neunet.2017.02.013>

Download English Version:

<https://daneshyari.com/en/article/4946668>

Download Persian Version:

<https://daneshyari.com/article/4946668>

[Daneshyari.com](https://daneshyari.com)