

A Fusion of Self-Cure Network and Transfer Learning-Based Deep Convolutional Neural Network for Facial Emotion Recognition

Tien Phan

Earlham College

Richmond, Indiana, United States

tphan21@earlham.edu

ABSTRACT

Recent advancements in deep learning, especially the application of Convolutional Neural Networks (CNNs), have guided significant progress in Facial Emotion Recognition (FER). CNNs excel at automatically extracting facial features, leading to the development of robust FER systems. Deep Convolutional Neural Networks (DCNNs) have further elevated FER capabilities, outperforming traditional CNNs. However, DCNNs still face a limitation concerning their demand for high-dimensional datasets. To address this constraint, this paper proposes an ensemble model combining Self-Cure Network (SCN) with DCNN, bolstered by transfer learning.

KEYWORDS

Convolutional neural networks (CNN), deep convolutional neural networks (DCNN), self-cure network (SCN), transfer learning, facial emotion recognition (FER)

ACM Reference Format:

Tien Phan. 2024. A Fusion of Self-Cure Network and Transfer Learning-Based Deep Convolutional Neural Network for Facial Emotion Recognition. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nmnnnnn.nnnnnnn>

1 INTRODUCTION

Facial expressions play a vital role in human communication, acting as a universal language conveying emotions, intentions, and social cues. The accurate recognition and interpretation of these expressions have far-reaching implications across various domains such as psychological treatment or surveillance [8]. Consequently, Facial Expression Recognition (FER) has attracted substantial attention and investment over the years.

Among the methods used in FER, Convolutional Neural Networks (CNNs) have gained immense popularity due to their ability to automatically extract features and adapt to diverse datasets with high precision. Deep Convolutional Neural Networks (DCNNs), a specialized form of CNNs, excel at processing high-dimensional images. However, the widespread adoption of DCNNs has been limited by the substantial computational effort and extensive training data they demand.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, July 2017, Washington, DC, USA

© 2024 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nmnnnnn.nnnnnnn>

To enhance the training process of DCNNs, Akhand et al. [1] introduced a promising approach—pre-training DCNN models through transfer learning. While this method has improved the efficiency and accuracy of DCNNs, it encounters challenges when dealing with datasets containing low-resolution images or highly imbalanced cases. To overcome these obstacles, this paper will propose the fusion of the Self-Cure Network [12] and DCNN models utilizing transfer learning, offering a solution to these limitations.

The paper is structured in four sections: Section 2 provides background information and related works. Section 3 describes the design of the project. Section 4 provides the timeline of the project.

2 RELATED WORK

This section will provide the background knowledge related to the method proposed in this paper including topics like DCNN, transfer learning, and SCN.

2.1 Transfer learning based Deep Convolutional Neural Networks Model

A CNN typically comprises a series of convolutional layers, followed by pooling layers and fully connected layers. DCNNs are a kind of CNN that are characterized by their depth, consisting of multiple layers, including convolutional, pooling, and fully connected layers. These layers collaborate to autonomously learn and extract features from images. Convolutional layers capture basic features like edges and textures, while deeper layers identify complex patterns and structures [1].

Building and training a DCNN from scratch is a heavy task due to its complexity. To overcome this challenge, this project employs the approach proposed by Akhand et al. [1], which involves utilizing pre-trained DCNN models. Some popular DCNN models such as AlexNet [2], VGG, and ResNet [7] have demonstrated exceptional performance in image classification tasks on ImageNet [11], a vast dataset of labeled images. These models are pre-trained on ImageNet, where they acquire the ability to recognize a wide range of objects and features within images. This is where the concept of transfer learning comes into play.

Transfer learning is a machine learning technique that applies knowledge learned from one task to a different yet related task. In DCNN, transfer learning entails taking a pre-trained model, such as one trained on ImageNet, and adapting it to a new task, such as facial emotion recognition. This approach leverages the knowledge acquired during the pre-training phase, providing the model with a strong foundation, even when the new task has limited data available. The method proposed by Akhand et al. [1] underscores the significance of fine-tuning within the context of transfer learning.

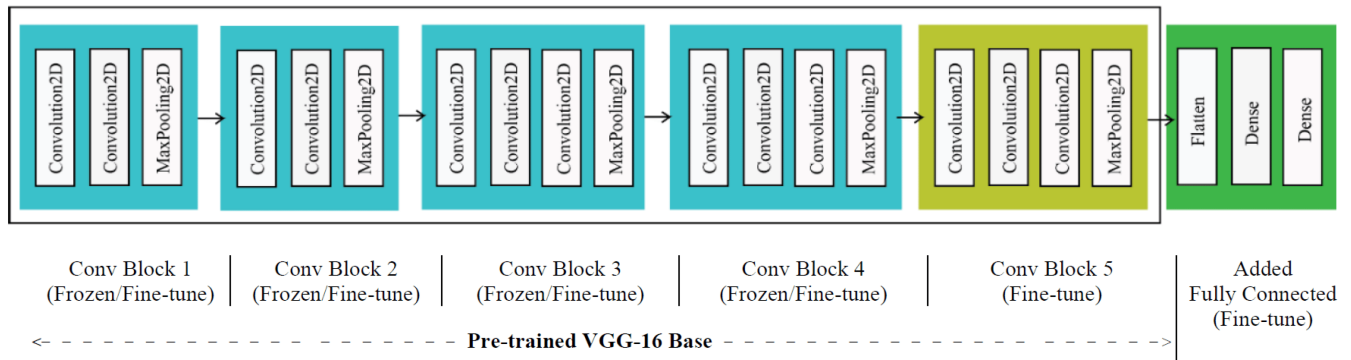


Figure 1: The architecture of DCNN with a detailed pre-trained VGG-16 model with dense layers for FER [1]

Fine-tuning is a critical step in transfer learning. It involves modifying the pre-trained DCNN by replacing or adding layers specifically tailored to the new task, such as facial emotion recognition. Typically, the existing convolutional base is retained, and new layers, often fully connected layers, are added to adapt the model to the target task. These newly added layers are initialized with random weights and are then trained using the new dataset.

In the case of transfer learning-based DCNN for FER, the pre-trained model (e.g., VGG-16 trained on ImageNet) receives modifications to suit emotion recognition by redefining its dense layers. Subsequently, fine-tuning is executed using emotion data. In this process, the last dense layers of the pre-trained model are replaced with new dense layers, which are designed to recognize facial images categorized into one of seven emotion classes, such as afraid, angry, disgusted, sad, happy, surprised, and neutral.

The complete model, composed of the pre-trained DCNN and the newly added dense layers, facilitates the fine-tuning of dense layers and selected layers of the pre-trained model using emotion data. The pre-trained VGG-16 model, as shown in Figure 1, consists of five convolutional blocks, with each block featuring two or three convolutional layers and a pooling layer.

2.2 Self-Cure Network

Wang et al. [12] introduced the Self-Cure Network (SCN), an extension of conventional CNNs. The primary objective of SCN is to address the challenges arising from the subjective nature of FER datasets obtained from the internet, which often result in inconsistent and erroneous labels. SCN comprises three essential components: self-attention importance weighting, ranking regularization, and relabeling. When given a dataset containing uncertain samples, the system first extracts deep features using a backbone network like CNNs. The self-attention importance weighting module assigns significant weights to each image, while the rank regularization module adjusts the attention weights to downplay the importance of uncertain samples. Lastly, the relabeling module modifies some of the uncertain samples within the low importance group. Figure 2 visually shows the structure of SCN: the process starts by using a backbone CNN to extract features from facial images. Then, a self-attention importance weighting module learns weights for each sample based on facial features, influencing the loss calculation.

The rank regularization module uses these weights and applies constraints using a ranking operation and a margin-based loss function. In the labeling module, reliable samples are identified by comparing predicted probabilities with given labels. Mislabeled samples are marked with red rectangles, and ambiguous ones with green dashed rectangles. It's important to note that the Self-Cure Network mainly uses re-weighting to handle uncertainties and adjusts only certain samples.

To assess the effectiveness of Self Cure Network (SCN), the study conducted experiments using four distinct datasets: RAF-DB [9], FERPlus [3], AffectNet [10], and a dataset referred to as WebEmotion, which contained data collected from the Internet. The primary focus was on addressing the challenge of incorrect or noisy annotations, which led to the creation of the WebEmotion dataset. It is a video dataset that was originally obtained from YouTube, but for the purpose of this study, it was treated as an image dataset by assigning emotion labels to individual frames. The dataset was curated by searching for videos using a set of 40 emotion-related keywords, combining them with 45 country-related keywords including countries in Asia, Europe, Africa, and America and six age-related keywords (baby, lady, woman, man, old man, old woman). The WebEmotion dataset comprises the same eight emotion classes as FERPlus [3], with each class being associated with various emotion-related keywords. For example, the "happy" class is linked to keywords such as happy, funny, ecstatic, smug, and kawaii. To establish meaningful correlations between these keywords and the collected videos, only the top 20 retrieved videos, each lasting less than four minutes, were selected for inclusion in the dataset. This curation process resulted in approximately 41,000 videos, which were further segmented into 200,000 video clips, with the requirement that a human face, detected using the Multi-task Cascaded Convolutional Neural Networks (MTCNN) [13] face detection method, must appear for at least five seconds in each clip. Moreover, they employed ResNet-18 [7] as the foundational network, which had been pre-trained on the MS-Celeb-1M face recognition dataset [6]. The extraction of facial features was performed from the final pooling layer of this network.

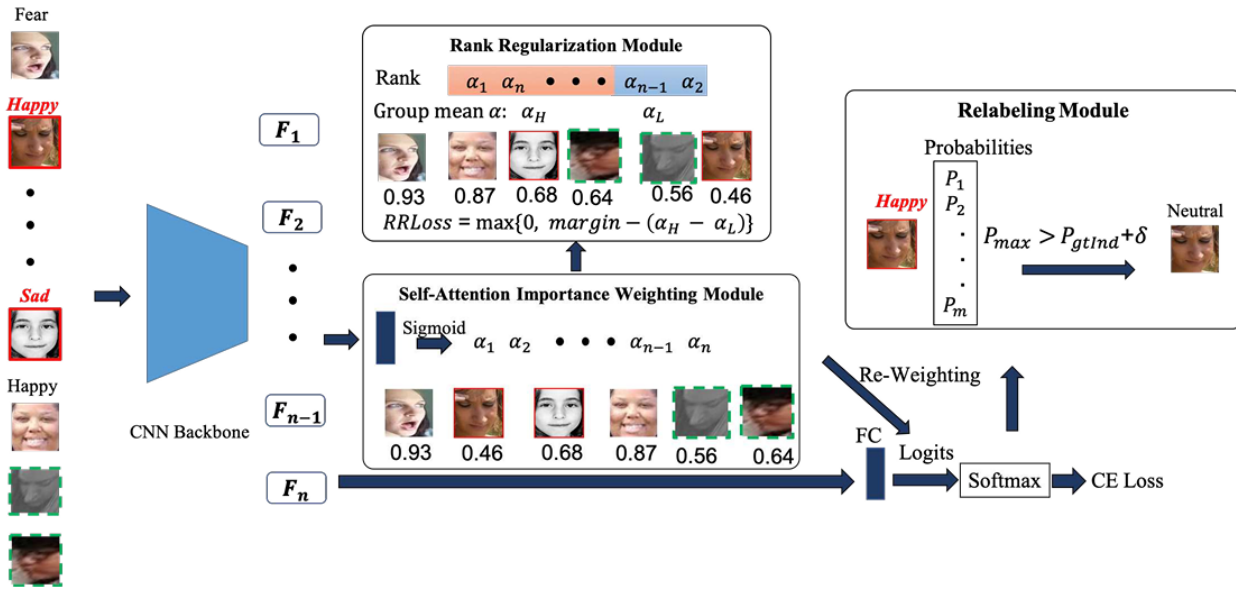


Figure 2: The architecture of SCN [8]

3 DESIGN

In this project, an ensemble approach will be implemented, combining the SCN with a transfer learning-based DCNN using a stacking strategy. The results obtained from this ensemble method will be compared with the existing state-of-the-art methods to assess its performance.

3.1 Implementation

For the implementation of this project, the integration of the Self-Cure Network (SCN) and the Deep Convolutional Neural Network (DCNN) will be achieved through a well-defined process. Both the SCN and DCNN will share the same pre-trained model, ResNet-18, which has previously been trained on the extensive MS-Celeb-1M dataset. To ensure data compatibility, image resizing will be done using PyTorch. The resized images will serve as inputs for the DCNN, which features the ResNet-18 model along with the addition of new dense layers for FER. The next step involves channeling the DCNN's output into the SCN, which will also make use of the same pre-trained ResNet-18 model. The SCN's role is to address uncertainties and improve the overall quality of predictions. At the end, the result generated by the SCN will serve as the final outcome for comparison and evaluation against existing methods and state-of-the-art solutions. The integration process ensures that the strengths of both the SCN and DCNN are utilized effectively, providing a robust solution for FER. The shared pre-trained model, ResNet-18, forms a common foundation for feature extraction.

3.2 Datasets

The project will use KDEF [4] dataset, which was developed by Karolinska Institute, Department of Clinical Neuroscience, Section of Psychology, Stockholm, Sweden. The dataset contains 4900

images with high-resolution and clear labels of standard facial expression from multiple angles.

In addition to that, the project will use the FER-2013 [5] dataset, which was derived from the Google search engine and consists of a large set of faces automatically registered and labeled with basic emotions. The images are more varied in terms of quality, angle, lighting, and occlusion.

3.3 Project Setup

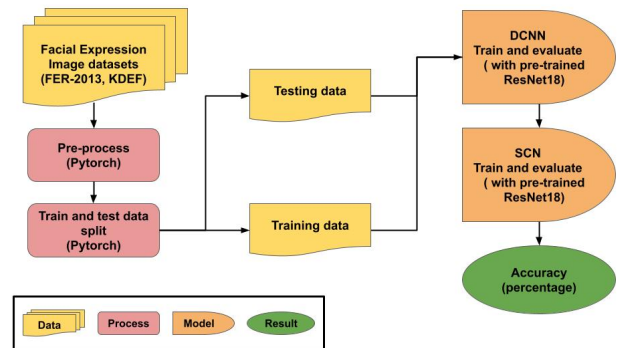


Figure 3: Data Architecture Diagram

Before training the models, I followed the data transformation based on DCNN code source [1]. For training data, the images are moderated using torchvision model from PyTorch. First, it is resized to 224x224 pixels since this is a common input size for many CNNs. After that, images are randomly flipped horizontally and converted to a torch.FloatTensor of shape (CxHxW). Images are also normalized with a mean and a standard deviation: the

normalization values for the mean are $[0.485, 0.456, 0.406]$ and for the standard deviation are $[0.229, 0.224, 0.225]$. These values are specific for model that has been pretrained on ImageNet dataset, which is ResNet18 for this project.

For testing data, images are resized to 2256x256 pixels and cropped at the center to 224x224 pixels. After that, testing images follow the same conversion and normalization conditions as training images.

After preparing the data, the next step is loading the data using DataLoader from PyTorch, which specify that each batch of data will contain 4 images and after each epoch, the images will be shuffle to prevent the model from the learning the order of the images.

3.4 Result

Currently, I finished training and evaluating the DCNN model. The best accuracy for FER-2013 is 65%. Here is some visualization output from one batch:

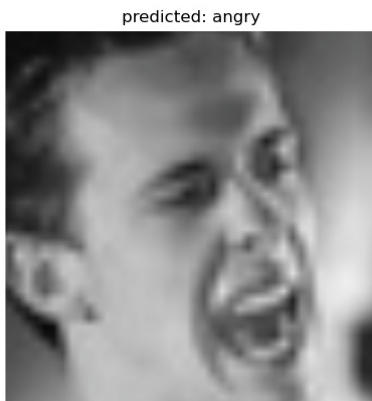


Figure 4: Output 1

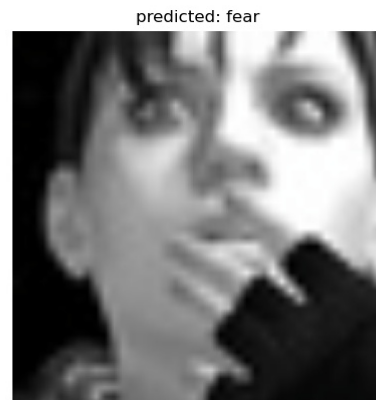


Figure 5: Output 2

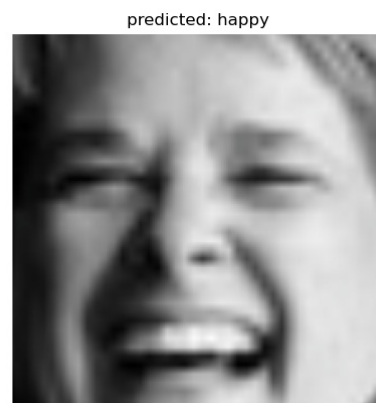


Figure 6: Output 3

3.5 Discussion

Before starting to work on the SCN model, all testing images in FER-2013 need to be fed into DCNN to create a new input for the SCN model. FER-2013 is a dataset with low quality images, which explained the accuracy from DCNN. Since SCN is supposed to improve the weakness of DCNN has with low quality images, I am excited to see how SCN will help improve the overall accuracy. However, there are still concerns in how the output of DCNN should be used in SCN: whether as training or testing data. Besides, the output of DCNN needs to be organized and format in a suitable way for SCN's input.

ACKNOWLEDGMENTS

I would like to thank Dr. David Barbella and Dr. Charlie Pecks for his support and feedback on this project.

REFERENCES

- [1] M.A.H. Akhand, Shuvendu Roy, Nazmul Siddique, Md Abdus Samad Kamal, and Tetsuya Shimamura. 2021. Facial emotion recognition using transfer learning in the deep CNN. *Electronics* 10, 9 (2021), 1036.
- [2] Grigorios Antonellis, Andreas G. Gavras, Marios Panagiotou, Bruce L Kutter, Gabriele Guerrini, Andrew C. Sander, and Patrick J Fox. 2015. Shake table test of large-scale bridge columns supported on rocking shallow foundations. *Journal of Geotechnical and Geoenvironmental Engineering* 141, 5 (2015), 04015009.
- [3] Emad Barsoum, Cha Zhang, Cristian Canton Ferrer, and Zhengyou Zhang. 2016. Training deep networks for facial expression recognition with crowd-sourced label distribution. In *Proceedings of the 18th ACM international conference on multimodal interaction*. 279–283.
- [4] Manuel G Calvo and Daniel Lundqvist. 2008. Facial expressions of emotion (KDEF): Identification under different display-duration conditions. *Behavior research methods* 40, 1 (2008), 109–115.
- [5] Ian J Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, et al. 2013. Challenges in representation learning: A report on three machine learning contests. In *Neural Information Processing: 20th International Conference, ICONIP 2013, Daegu, Korea, November 3-7, 2013. Proceedings, Part III* 20. Springer, 117–124.

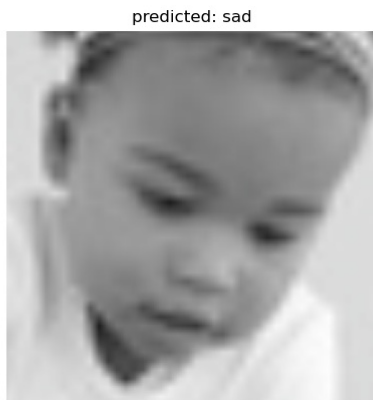


Figure 7: Output 4

- [6] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. 2016. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III 14*. Springer, 87–102.
- [7] 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.
- [8] Yunxin Huang, Fei Chen, Shaohe Lv, and Xiaodong Wang. 2019. Facial expression recognition: A survey. *Symmetry* 11, 10 (2019), 1189.
- [9] Shan Li, Weihong Deng, and JunPing Du. 2017. Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2852–2861.
- [10] Ali Mollahosseini, Behzad Hasani, and Mohammad H. Mahoor. 2017. Affectnet: A database for facial expression, valence, and arousal computing in the wild. *IEEE Transactions on Affective Computing* 10, 1 (2017), 18–31.
- [11] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [12] Kai Wang, Xiaojiang Peng, Jianfei Yang, Shijian Lu, and Yu Qiao. 2020. Suppressing uncertainties for large-scale facial expression recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 6897–6906.
- [13] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. 2016. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters* 23, 10 (2016), 1499–1503.