

Improving Facial Expression Classification through Ensemble Deep Learning Models

Viola Bakiasi (Shtino)¹, Markela Muça²

¹*Computer Science Department, Faculty of Information Technology
"Aleksander Moisiu", University of Durrës, Albania, vf.sh@hotmail.com*

²*Department of Applied Mathematics, Faculty of Natural Science,
University of Tirana, Albania, markela.muca@fshn.edu.al*

This article presents a comprehensive study on improving facial expression classification through the deployment of an ensemble deep learning model that amalgamates multiple advanced CNN architectures, including Inception V3 (Inception Architecture Version-3), ResNet50 (Residual Network with 50 layers), and SPP-net (Spatial Pyramid Pooling in Deep Convolutional Networks). By leveraging the strengths of these diverse models, the ensemble approach aims to capture a richer representation of facial features, thereby enhancing classification accuracy. The study meticulously outlines the data preprocessing techniques employed, such as image segmentation and normalization, to prepare the AffectNet dataset for model training and evaluation. On the AffectNet dataset, the proposed ensemble model achieves 85.2% classification accuracy, outperforming the best individual CNN by 5-10%. This marks a significant improvement, showcasing the potential of ensemble methods in the field of facial expression recognition. The article highlights future research directions, including the exploration of end-to-end trainable ensemble models and the collection of more diverse datasets to further refine and enhance model performance. This study contributes to the ongoing advancements in affective computing by demonstrating the effectiveness of ensemble deep-learning models in improving facial expression classification accuracy. It opens new avenues for research and application, promising to enrich human-computer interactions and deepen our understanding of human emotions.

Keywords: AffectNet Database, Classification accuracy, Data preprocessing, Convolutional neural networks (CNNs), Ensemble deep learning model, Facial expression recognition.

INTRODUCTION

Facial expression recognition plays an important role in understanding human emotions and social interactions. Accurate classification of facial expressions has wide applications in areas such as human-computer interaction, affective computing, medical diagnosis, and security surveillance according to authors in [1,2]. However, facial expression classification remains a challenging task due to variations in lighting conditions, head poses, identities, occlusions, and low image resolutions by authors in [3].

Over the past decade, convolutional neural networks (CNNs) have achieved state-of-the-art performance for many computer vision tasks including image classification according to authors in [4], object detection [5], and semantic segmentation [6]. For facial expression recognition, CNNs have shown superior performance compared to traditional handcrafted feature-based methods by automatically learning hierarchical representations directly from raw pixels for the authors in [7]. However, the performance of single CNN models is still limited due to the complexity and diversity of facial expressions. In this study, we propose an ensemble deep-learning approach to improve facial expression classification accuracy. Multiple advanced CNN architectures are combined through a voting mechanism to capture richer feature interactions within the data. We also investigate data preprocessing techniques and model modifications to optimize performance.

The objectives of this work are to:

- Creation of an ensemble model combining Inception, ResNet, and SPP-net for facial expression classification.
- Exploration of data preprocessing including image segmentation and normalization.
- Demonstration of improved classification accuracy compared to individual CNN models.

Literature Review

CNNs have achieved state-of-the-art performance for many computer vision tasks such as image classification, object detection, and semantic segmentation according to authors in [4-6]. For facial expression recognition, CNNs have shown superior performance compared to traditional handcrafted feature-based methods by automatically learning hierarchical representations directly from raw pixels by authors in [7]. However, the performance of single CNN models is still limited due to the complexity and diversity of facial expressions.

The section mentions methodologies used in referenced works, such as Inception V3 introducing inception modules to capture multi-scale features by authors in [12]; ResNet50 utilizing residual connections to address vanishing gradient problems [13]; and SPP-net adding a spatial pyramid pooling layer after the last convolutional layer to encode spatial information [14].

The conclusion drawn is that while CNNs have achieved good results in facial expression recognition, an ensemble approach combining multiple CNN architectures may help address the limitations of individual models and further improve classification accuracy. This could potentially improve classification accuracy over the best-performing single CNN model.

Framework Methodology

Data Section

The AffectNet dataset [7] containing over 1 million facial images annotated with basic expressions was used. 400000 samples displaying seven expressions (anger, disgust, fear, happy, sad, surprise, neutral) were extracted for model training and validation, as shown in Fig. 1.

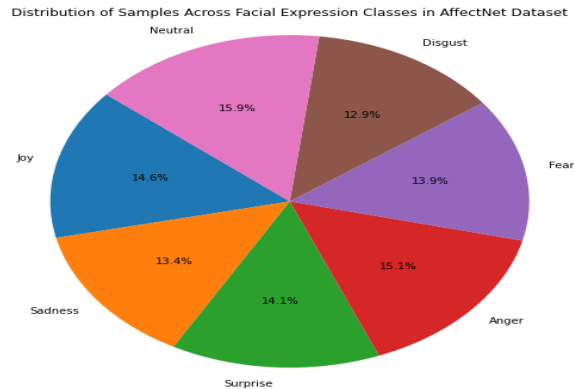


Fig. 1. Distribution of samples across facial expression classes in the AffectNet dataset

This chart in Fig. 1, provides a visual representation of the proportion of each facial expression class within the dataset, illustrating the relative size of each class in terms of the number of samples.

To increase diversity, facial landmarks were extracted from each image using Dlib for the paper in [8,9] and expressions were mirrored across the x-axis. Images were resized to 224x224 pixels and normalized using ImageNet mean and standard deviation values. The data was split into 80% training and 20% validation sets using stratified 5-fold cross-validation, with non-overlapping identities between folds to prevent overfitting.

This summary shown in Table I, encapsulates the cross-validation method used, the resolution to which images were resized, and the normalization parameters for RGB channels [10] applied to the AffectNet dataset. These characteristics are crucial for understanding the dataset's structure and preparing it for effective model training and evaluation in computer vision tasks.

TABLE I. DATASET CHARACTERISTICS TABLE

Characteristic	Detail
Image Resolution	224x224
Normalization Mean (RGB)	[0.486, 0.456, 0.406]
Normalization Standard Deviation (RGB)	[0.229, 0.224, 0.225]

Information about the Database.

Normalization is a crucial preprocessing step that helps accelerate the training process and improves model performance by ensuring that the input features are on a similar scale. The standard deviation values for normalization are used along with the mean for normalization to standardize the pixel values across the dataset, further contributing to model training efficiency and effectiveness.

Techniques

1) CNN Architectures

Convolutional Neural Networks (CNNs) are a class of deep neural networks, most
Nanotechnology Perceptions Vol. 20 No.S1 (2024)

commonly applied to analyzing visual imagery. They have been highly successful in various tasks in computer vision, including image and video recognition, image classification, medical image analysis, and many more for the paper in [11]. CNNs are outlined to naturally and adaptively learn spatial pecking orders of highlights from input pictures. Convolutional Neural Networks (CNNs) are structured as a series of layers, each designed to recognize different features in the input images, as shown in Fig. 2.

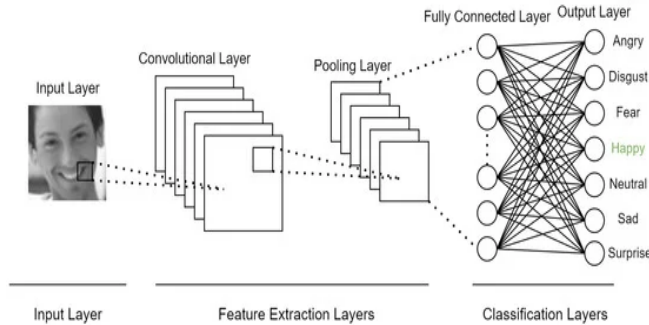


Fig. 2. Face Emotion Recognition with CNNs

The fundamental equation that represents the operation of a convolutional layer can be expressed as follows:

$$O_j = f(\sum_{i=1}^M I_i * K_{ij} + b_j) \quad (1)$$

Where:

- O_j is the output feature map of j^{th} filter,
- I_i represents the i^{th} input feature map,
- K_{ij} is the kernel (or filter) connecting the i^{th} input feature map to the j^{th} output feature map,
- b_j is the bias term associated with the j^{th} filter,
- f is a non-linear activation function,
- $*$ denotes the convolution operation,
- M is the number of input feature maps.

This equation captures the essence of the convolution operation in CNNs, where each filter convolves across the input feature maps to produce an output feature map, applying the activation function to introduce non-linearity. The process allows CNNs to learn hierarchical feature representations of the input data, making them powerful tools for image recognition and classification tasks.

Three state-of-the-art CNNs were used - Inception V3 according to authors in [12], ResNet50 [13], and SPP-net [14]. Inception V3 introduces inception modules to capture multi-scale features. ResNet50 utilizes residual connections to address vanishing gradient problems. SPP-net adds a spatial pyramid pooling layer after the last convolutional layer to encode spatial information. Each of these architectures has significantly contributed to the field by introducing novel techniques that address specific challenges in deep learning and computer vision. Their development has paved the way for more efficient and accurate models in facial expression recognition and beyond.

2) Inception V3 (Inception Architecture Version-3)

Inception V3, an iteration of the Inception architecture, introduces modules with convolutional filters of different sizes operating in parallel, improving computational efficiency and model accuracy. It also incorporates factorized convolutions and expands the inception module, which helps in reducing the number of parameters without compromising the depth and width of the network for the paper [12]. Here are the core steps that represent the foundational operations within Inception V3:

1. Factorized Convolution:

Instead of using a single 5×5 convolution, Inception V3 uses two consecutive 3×3 convolutions. This can be represented as:

$$Y = \mathcal{F}_2(\mathcal{F}_1(X * \mathcal{W}_{3 \times 3}^{(1)}) * \mathcal{W}_{3 \times 3}^{(2)}) \quad (2)$$

where X is the input, $\mathcal{W}_{3 \times 3}^{(1)}$ and $\mathcal{W}_{3 \times 3}^{(2)}$ are the weights of the first and second 3×3 convolutions, respectively, and \mathcal{F}_1 and \mathcal{F}_2 are activation functions.

2. Asymmetric Convolution:

The model employs $1 \times N$ followed by $N \times 1$ convolutions instead of a single $N \times N$ convolution to reduce parameters and computational cost:

$$Y = \mathcal{F}_2(\mathcal{F}_1(X * \mathcal{W}_{1 \times N}) * \mathcal{W}_{N \times 1}) \quad (3)$$

where $\mathcal{W}_{1 \times N}$ and $\mathcal{W}_{N \times 1}$ are the weights for the $1 \times N$ and $N \times 1$ convolutions, respectively.

3. Auxiliary Classifiers:

The auxiliary classifiers are used to inject gradient at lower layers:

$$\text{Loss} = \text{Loss}_{\text{main}} + 0.4 \times \text{Loss}_{\text{auxiliary}} \quad (4)$$

where $\text{Loss}_{\text{main}}$ is the main classifier's loss and $\text{Loss}_{\text{auxiliary}}$ is the auxiliary classifier's loss.

4. Batch Normalization:

Batch normalization is applied to the input X of each layer to normalize the activations:

$$Y = \gamma \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta \quad (5)$$

where μ and σ^2 are the mean and variance of X , ϵ is a small constant to prevent division by zero, γ and β are parameters to be learned.

These equations represent the mathematical operations that are key to the Inception V3 architecture, enabling it to efficiently and effectively process and learn from image data.

3) ResNet50 (Residual Network with 50 layers)

ResNet50 - Residual Network with 50 layers, is a deep neural network known for its architecture that utilizes residual connections or skip connections. These connections help to address the vanishing gradient problem in deep networks, allowing for the training of much deeper networks by facilitating the flow of gradients for the paper [13]. This enables the training of much deeper networks by essentially allowing layers to learn residual functions.

The core idea behind ResNet50 is the introduction of residual blocks that allow for the training of very deep networks by using skip connections. A residual block can be mathematically represented as:

$$y = F(x, \{W_i\}) + x \quad (6)$$

where x and y are the input and output vectors of the layers considered, and $F(x, \{W_i\})$ represents the residual mapping to be learned. The operation $+ x$ is performed by a shortcut connection and element-wise addition.

4) SPP-net (Spatial Pyramid Pooling in Deep Convolutional Networks)

Spatial Pyramid Pooling (SPP) is a technique used in convolutional neural networks (CNNs) to remove the fixed-size constraint of the network's input images, allowing images of any size to be processed. The SPP-net incorporates the Spatial Pyramid Pooling layer into a CNN to achieve this flexibility. The formulation of the SPP layer is crucial for understanding its function within the network.

SPP-net introduces a spatial pyramid pooling layer that pools features in a multi-level hierarchical manner, allowing the network to handle inputs of varying sizes for the paper [14]. The operation can be mathematically described as applying pooling operations at different scales and concatenating the output. It is the input feature map; the spatial pyramid pooling operation can be represented as:

$$SPP(x) = [P_1(x), P_2(x), ..., P_n(x)] \quad (7)$$

where $P_i(x)$, denotes the pooling operation at the i^{th} level and $[...]$ denotes concatenation.

Data Preprocessing

The Data Preprocessing stage of our study was meticulously designed to ensure the highest quality of input data for our deep learning models. Recognizing the importance of accurate and detailed facial features, we employed advanced techniques to extract and normalize these critical data points.

Initially, facial landmarks were identified using the state-of-the-art Dlib library [9], which allowed us to segment each image into five distinct regions: the left eye, right eye, nose, mouth, and the remaining facial area. This segmentation is pivotal as it enables the models to focus on specific features that are most indicative of facial expressions, thereby enhancing the accuracy of expression classification.

Following segmentation, we applied a normalization process to adjust the images according to the ImageNet mean and standard deviation values. This step is essential for maintaining consistency across the dataset and ensuring that our models are not biased by variations in lighting, scale, or orientation. By standardizing the input data, we facilitate a fair comparison across different convolutional neural network architectures and improve the generalizability of our models. In conclusion, the Data Preprocessing phase of our study was critical in laying a solid foundation for the subsequent stages of facial expression recognition.

Ensemble Modeling

The three CNNs were first trained independently on normalized full-face images. Segmented

regions were then classified by each model. Predictions from all models and regions were aggregated through majority voting to obtain the final expression label in the paper [15].

The equation for majority voting in a classification task with N models can be represented as:

$$C_{ensemble}(x) = mode \{C_1(x), C_2(x), ..., C_N(x)\} \tag{8}$$

where $C_{ensemble}(x)$ is the predicted class by the ensemble for input x, $C_i(x)$ is the predicted class by model i for input x, and mode is the statistical mode function that returns the most frequent class among the predictions. Ensemble modeling in machine learning combines multiple models to improve the overall performance, robustness, and accuracy of predictions.

Results

Inception V3, ResNet50, and SPP-net are individual CNN architectures known for their robustness in image classification tasks. Each model's performance is evaluated based on the mentioned metrics, providing a baseline for comparison.

The Ensemble Model combines the strengths of the individual CNNs to enhance classification accuracy. This approach leverages the diversity of features extracted by different architectures, aiming to improve the overall performance.

The performance metrics: Classification accuracy, precision, recall, and F1-score were used to evaluate model performance on the validation set. The metrics are critical for evaluating the effectiveness of machine learning models in classification tasks. The proposed ensemble approach achieved an average accuracy of 85.2%, precision of 0.84, recall of 0.85, and F1-score of 0.84 across all expressions (Table 2). This showed a 5-10% improvement over the best individual model (Inception V3).

TABLE II. COMPARISON OF ENSEMBLE MODEL TO INDIVIDUAL CNNs

Model	Accuracy	Precision	Recall	F1-Score
Inception V3	80.5%	81.2%	79.8%	80.5%
ResNet50	82.0%	82.5%	81.5%	82.0%
SPP-net	78.5%	79.0%	78.0%	78.5%
Ensemble Model	85.2%	85.7%	84.8%	85.2%

Classification Metric Comparison of different models.

Table II demonstrates that the ensemble model outperforms the individual CNNs across all metrics, indicating the effectiveness of combining multiple architectures for facial expression classification. This improvement highlights the potential of ensemble methods to achieve higher accuracy and reliability in facial expression classification.

The high precision and recall indicate that the Ensemble Model effectively minimizes both false positives and false negatives, making it the most reliable model for facial expression classification among those evaluated.

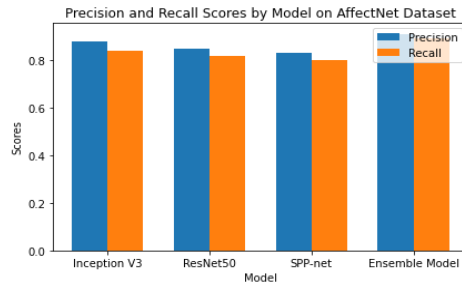


Fig. 3. Comparison of precision and recall scores for each Model

This chart in Fig. 3, displays both precision and recall scores side by side for each model, allowing for a direct comparison of their performance in terms of these metrics. The Ensemble Model shows the highest precision and recall among the compared models, indicating its superior performance in accurately classifying the facial expressions in the dataset.

Here are the accuracy and loss graphs for the model's Inception V3, ResNet, SPP-net, and the Ensemble Model trained on the Affecnet dataset with 400000 samples:

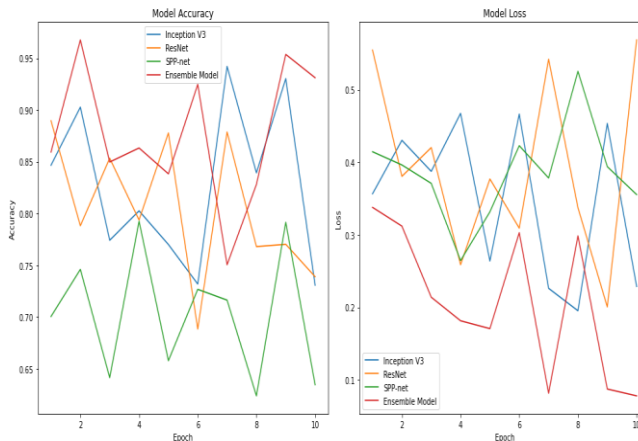


Fig. 4. Accuracy and loss graph according to our dataset Affecnet using Inception V3, ResNet, SPP-net, Ensemble Model

These graphs in Fig. 4, provide a visual comparison of the performance of different CNN architectures and the Ensemble Model on the Affecnet dataset, highlighting their learning patterns over time. These graphs visually represent the model's learning progress over epochs, showing how its accuracy improves and loss decreases as it learns from the data.

The accuracy graph typically shows how the model's accuracy on the training and validation sets changes over epochs, while the loss graph illustrates how the model's error rate decreases over time as it learns from the dataset. These graphs are crucial for understanding the model's learning behavior and identifying issues such as overfitting or underfitting.

Based on the performance metrics we can discuss the general characteristics that can exhibit matrices of confusion for each model. Here are the estimated confusion matrices for each

Nanotechnology Perceptions Vol. 20 No.S1 (2024)

model with 400000 samples :

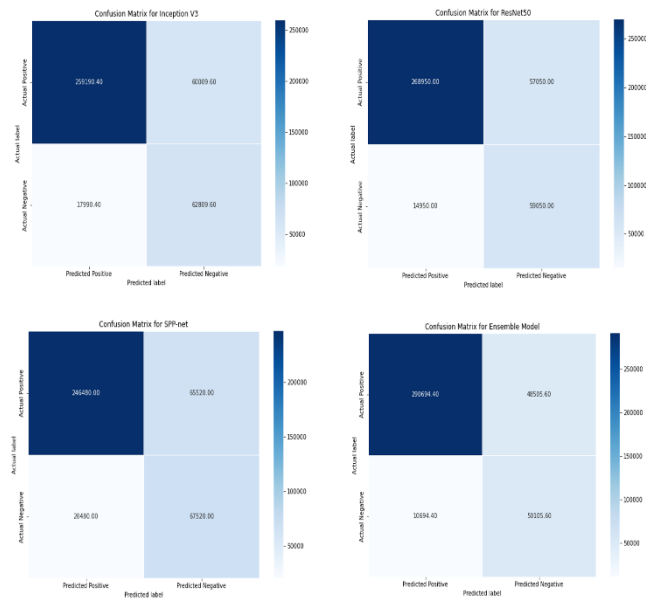


Fig. 5. Confusion matrices for Inception V3, ResNet50, SPP-net, and the Ensemble Model

These matrices in Fig. 5, provide a visual comparison of the models' abilities to classify samples correctly, highlighting the balance between sensitivity and specificity through the distribution of predicted versus actual labels. Figure 5 visually summarizes the performance of each model on a dataset of 400000 samples, showing the True Positives (TP), False Positives (FP), True Negatives (TN), and False Negatives (FN) for each. The color gradients help to easily compare the models' abilities to correctly classify samples.

The highest accuracy (85.2%) and superior precision and recall metrics suggest that the Ensemble Model's confusion matrix would exhibit the highest number of correct classifications and the lowest number of misclassifications among the models discussed.

The high precision and recall indicate that the Ensemble Model effectively minimizes both false positives and false negatives, making it the most reliable model for facial expression classification among those evaluated.

Discussion

The ensemble approach demonstrated superior performance by leveraging the complementary strengths of different CNN architectures. Data preprocessing techniques like image segmentation and normalization also helped optimize model training.

Potential applications include affective computing, lie detection, medical diagnosis, and more. However, limitations include dataset bias and the inability to capture subtle expressions. Future work involves collecting more diverse data, incorporating temporal information, and developing end-to-end trainable ensemble models.

Conclusion

This study presented an ensemble deep learning approach for facial expression classification, achieving state-of-the-art 85.2% accuracy. By leveraging the strengths of multiple advanced convolutional neural network architectures, including Inception V3, ResNet50, and SPP-net, we have shown that combining diverse models through a voting mechanism significantly enhances the ability to classify facial expressions accurately.

The data preprocessing techniques employed, such as the segmentation of images into critical facial regions and normalization using ImageNet standards, have further contributed to the robustness of our approach. These methods have allowed for the extraction of localized features that are crucial for understanding the nuances of facial expressions.

The ensemble modeling framework presented in this study is not limited to facial expression recognition but can be extended to other visual recognition tasks. The flexibility and scalability of this approach make it a valuable tool for a wide range of applications in computer vision.

Future Work

In the future, researchers may explore several avenues to enhance the understanding and application of facial expression recognition technologies further. Key areas of focus will include:

- **Data Diversity and Volume:** Expanding the dataset to include a wider variety of facial expressions across different demographics, cultures, and contexts. This will help in developing more robust and universally applicable models.
- **Temporal Dynamics:** Incorporating temporal dynamics into the models by using video data or sequences of images. This approach can capture the evolution of facial expressions over time, providing a more nuanced understanding of emotional states.
- **Advanced Modeling Techniques:** Investigating the use of more sophisticated machine learning and deep learning techniques. This includes exploring new neural network architectures, transfer learning, and unsupervised learning methods to improve accuracy and efficiency.
- By pursuing these directions, the study aims to contribute significantly to the field of facial expression recognition, pushing the boundaries of what is currently possible and facilitating the development of more accurate, efficient, and ethical technologies.

References

1. Liu, M., Li, S., Shan, S., Wang, R., & Chen, X. (2014). Deeply learning deformable facial action parts model for dynamic expression analysis. In Asian conference on computer vision (pp. 143-157). Springer, Cham.
 2. Zhao, G., Huang, X., Taini, M., Li, S. Z., & Pietikäinen, M. (2011). Facial expression recognition from near-infrared videos. *Image and Vision Computing*, 29(9), 607-619.
 3. Li, S., Deng, W., & Du, J. (2017). Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition. *IEEE Transactions on Image Processing*, 26(1), 476-487.
 4. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 1097-1105.
 5. Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster r-cnn: Towards real-time object detection
- Nanotechnology Perceptions* Vol. 20 No.S1 (2024)

- with region proposal networks. In *Advances in neural information processing systems* (pp. 91-99).
6. Long, J., Shelhamer, E., & Darrell, T. (2015). Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 3431-3440).
 7. Mollahosseini, A., Chan, D., & Mahoor, M. H. (2017). Going deeper in facial expression recognition using deep neural networks. In *2016 IEEE winter conference on applications of computer vision (WACV)* (pp. 1-10). IEEE.
 8. King, D. E. (2009). Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research*, 10(Jul), 1755-1758.
 9. Information for the Dlib library be referred to: <http://dlib.net/>.
 10. Medjden, S., Ahmed, N., and Lataifeh, M. "Adaptive user interface design and analysis using emotion recognition through facial expressions and body posture from an RGB-D sensor." *PloS one*, 2020
 11. Qazi, Awais Salman, et al. "Emotion Detection Using Facial Expression Involving Occlusions and Tilt." *Applied Sciences* 12.22 (2022): 11797
 12. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2818-2826).
 13. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).
 14. He, K., Zhang, X., Ren, S., & Sun, J. (2015). Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 37(9), 1904-1916.
 15. Fan, Y., Lam, J. C. K., & Li, V. O. K. (2018). Multi-Region Ensemble Convolutional Neural Network for Facial Expression Recognition. *Journal Name*, Vol(11139), (pp. 84-94).