

# Enhancing Masked Face Recognition Using Self Attention Mechanisms: A Comprehensive Study

Lekha Prajapati<sup>1</sup>

<sup>1</sup>Research Scholar, Department of Computer Science, Taywade College Koradi, (M.S.), India.

Girish Katkar<sup>2</sup>

<sup>2</sup>Assistant Professor, Department of Computer Science, Taywade College Koradi, (M.S.), India.

Ajay Ramteke<sup>3</sup>

<sup>3</sup>Assistant Professor, Department of Computer Science, Taywade College Koradi, (M.S.), India.

## Abstract-

The COVID-19 pandemic has highlighted the importance of wearing masks as a protective measure against viral infections. However, this has introduced significant challenges for existing face recognition systems. To address these challenges, advanced techniques such as self-attention mechanisms and CNN-based models are explored to enhance masked face recognition accuracy. In this study, we explore the Self-Attention Mechanism, originally a key component in Natural Language Processing (NLP), as a superior alternative to conventional CNNs for masked face recognition. Unlike CNNs, which rely on local receptive fields, self-attention captures long-range dependencies, focusing on unmasked facial regions. We evaluated to approach methods on the SMFRD and MFR2 datasets, comparing it with CNN-based methods. Experimental results demonstrate that self-attention enhances recognition accuracy by effectively learning global and local features, making it a robust solution for real-world masked face recognition applications.

## INTRODUCTION

Masked Face Recognition (MFR) has gained significant attention due to the widespread use of face masks, especially in security and authentication systems. Traditional face recognition models struggle with occlusions, as masks cover crucial facial features, making accurate identification a challenging task. The COVID-19 pandemic further highlighted this issue, as global health organizations such as the World Health Organization (WHO) and the Centers for Disease Control and Prevention (CDC) recommended wearing face masks in public spaces to prevent the spread of the virus. While face masks provide essential protection, they introduce challenges for conventional face recognition systems, leading to a notable decline in accuracy. This necessitates the development of robust techniques to mitigate the impact of face masks on recognition performance.

Deep learning-based approaches, particularly Convolutional Neural Networks (CNNs), have revolutionized face recognition by enabling robust feature extraction from large-scale datasets. Among these, ResNet and VGG16 have demonstrated remarkable success. ResNet employs a deep residual learning framework to effectively mitigate the vanishing gradient problem while capturing hierarchical facial features. On the other hand, VGG16 utilizes a uniform architecture with small convolutional filters to extract rich facial representations. Despite their effectiveness, these models primarily focus on local feature extraction and lack a global contextual understanding, making them less

effective for recognizing masked faces. The occlusion of lower facial features disrupts feature consistency, leading to degraded recognition accuracy. One of the major challenges in MFR is obtaining a diverse, large-scale dataset that includes various masked face types. Collecting such datasets is both time-consuming and labor-intensive. Additionally, maintaining diversity in facial expressions, mask types, and illumination conditions further complicates dataset creation. As a result, developing a low-cost and efficient data augmentation technique is crucial. Several methods have been proposed for generating synthetic masked faces. For example, MaskTheFace uses a Dlib-based face landmark detector to identify facial tilt and key facial points before overlaying a mask. Similarly, MaskedFace-Net employs a deformable mask-to-face mapping model and homographic transformation to align

masks onto facial areas. However, these methods primarily rely on affine transformations, often resulting in unrealistic mask placement and inconsistencies in pose and illumination. Generative Adversarial Networks (GANs) have also been explored for masked face augmentation. While GANs generate more realistic masked face images, they suffer from mode collapse, where generated images become overly similar despite different input conditions. Additionally, GAN-based methods are computationally expensive, making real-time augmentation challenging.

After reviewing multiple research papers on masked face recognition, I observed that many existing approaches

struggle with occlusions, limiting their ability to maintain high recognition accuracy. While traditional CNN-based methods such as ResNet and VGG16 have proven effective for face recognition, they primarily focus on local feature extraction and lack a global contextual understanding. Similarly, Transformer-based architectures provide strong global feature representation but often come with high computational costs. Some recent studies attempt to integrate attention mechanisms with CNNs, but they either fail to fully capture multi-scale facial features or do not prioritize critical facial regions effectively.

After carefully analyzing these limitations, I found that a lightweight MFR network integrating convolutional self-attention mechanisms presents a promising solution. This framework enhances both local and global feature extraction, ensuring robust recognition even with occluded facial regions. The proposed network consists of two key branches: the whole branch and the partial branch. The whole branch incorporates a Multi-scale Convolutional Self-Attention Module (MCSAM), which leverages CNNs and Transformers to capture multi-scale facial features. Additionally, a Multi-head Multi-scale Self-Attention (MMHSA) mechanism within the Transformer further strengthens global feature representation across spatial and channel dimensions. Meanwhile, the partial branch introduces an Upper Face Alignment Module (UFAM) that prioritizes upper facial features, enhancing overall recognition performance by assigning higher importance to the visible regions of the face. Given the persistent challenges in masked face recognition, this approach effectively balances feature extraction efficiency and recognition accuracy.

## **II. RELATED WORK**

### **A. Multi-scale Convolutional Self-Attention Module**

A lightweight masked face recognition (MFR) network is proposed, integrating convolutional self-attention architecture to enhance accuracy while preserving clean face recognition performance with minimal computational cost. The network consists of two main branches: the whole branch and the partial branch. The whole branch employs a Multi-scale Convolutional Self-Attention Module (MCSAM) to extract both local and global facial features by leveraging CNNs and Transformers. Within the Transformer, a Multi-head Multi-scale Self-Attention (MMHSA) mechanism is introduced to enhance global feature representation across spatial and channel dimensions. The partial branch incorporates an Upper Face Alignment Module (UFAM), which includes an alignment

module and MMHSA. By fusing multi-scale features of the upper face, this branch assigns higher importance to upper facial features, complementing the whole branch and improving overall MFR performance [1].

### **B. ConvNeXt-T Mechanisms**

This dual-branch design efficiently balances accuracy and computational efficiency, making it a robust solution for masked face recognition in real-world applications. This study employs ConvNeXt-T as the backbone for mask learning and recognition, leveraging its enhanced efficiency and accuracy over ResNet-50 by integrating Vision Transformer (ViT) training strategies. Key improvements include depthwise convolution, an Inverted Bottleneck, and optimized activation functions. Additionally, an Efficient Channel Attention (ECA) block is incorporated to refine feature extraction by prioritizing important features while maintaining channel dimensionality. Inspired by the human visual system, this attention mechanism enhances recognition performance by focusing computational resources on crucial information, improving accuracy and inference speed in deep learning applications such as image enhancement, segmentation, and object detection[2].

### **C. Convolutional Block Attention Module (CBAM)**

The occlusion discarding approach in face recognition (FR) removes occluded regions to reduce negative effects on recognition accuracy. Traditional methods use SVM classifiers or difference-based techniques to detect and discard occluded areas, but these rely on shallow features with limited discriminative power. More advanced techniques, like Pairwise Differential Siamese Network (PDSN) and BoostGAN, attempt to address these limitations but face challenges in real-world applications. Deep learning has significantly improved occlusion FR, leveraging deep feature extraction methods such as Dynamic Feature Matching (DFM) and GAN-based models. However, current approaches struggle with masked face recognition (MFR) due to the loss of key facial features. This paper integrates the attention mechanism into ResNet-50, using Convolutional Block Attention Module (CBAM) to focus on expressive features like the eyes while minimizing attention to occluded areas. CBAM outperforms existing methods by combining spatial and channel attention, improving feature representation in MFR tasks[3].

### **D. The CBAM-integrated ResNet-50 model**

Feature extraction plays a vital role in masked face recognition (MFR), as masks obscure essential facial features. This study employs a refined ResNet-50 architecture from ArcFace, modifying its third-stage layer structure to improve feature representation. The architecture is enhanced using a BN-Conv-BN-PReLU-Conv-BN residual unit, with a stride adjustment to refine feature extraction. The final 512-D face embedding is obtained through batch normalization, dropout, and fully connected layers. To further enhance feature extraction, the study integrates the Convolutional Block Attention Module (CBAM), a lightweight yet effective mechanism combining channel and spatial attention. The channel attention module assigns adaptive weights to crucial features using average and max pooling, followed by a multi-layer perceptron (MLP) and sigmoid activation. The spatial attention module highlights informative regions using a convolutional layer after pooling operations, refining spatial feature representation. By combining these mechanisms, the model effectively prioritizes expressive and discriminative facial features, such as the eyes, while minimizing attention to occluded areas like the lower face. This approach significantly improves masked face recognition accuracy, outperforming traditional methods. The CBAM-integrated ResNet-50 model demonstrates strong performance in handling occlusion challenges, making it suitable for real-world MFR applications[4].

#### **E. ResNet34-based face recognition model with ArcFace**

Masked face recognition presents challenges due to occlusion, requiring specialized techniques for accurate identification. This framework integrates Mask Transfer (MT), Attention-Aware Masked Face Recognition (AMaskNet), and Mask-Aware Similarity Matching (MS) to improve recognition performance. Mask Transfer (MT) generates synthetic masked face images by overlaying masks onto unmasked faces. A dataset, or mask gallery, is constructed with various mask types, including different colors, shapes, and textures. Only one representative image per mask type is needed, and frontal face images are preferred for better quality. The mask transfer process begins with pre-processing, where Dlib detects 68 facial landmarks and a Triangulated Irregular Network (TIN) is used to divide the face into triangular regions. The GrabCut algorithm segments mask areas. During transfer, affine transformations align mask regions with the unmasked face. Post-processing steps include alpha matting for seamless blending and histogram specification to normalize illumination differences.

AMaskNet improves masked face recognition by learning feature importance using a contribution estimator. A conventional ResNet34-based face recognition model with ArcFace loss is used to extract a 512-dimensional feature vector. The contribution estimator refines extracted features by learning spatial and channel-wise importance. The spatial estimator assigns higher weights to non-masked facial regions using a contribution matrix, while the channel estimator identifies crucial feature channels. The final refined features are obtained through matrix multiplications between contribution matrices and extracted features. Training employs ArcFace loss with an angular margin penalty to enhance intra-class compactness and inter-class separation. The Mask-Aware Similarity Matching (MS) strategy improves real-world face verification by addressing discrepancies between masked and unmasked images. Instead of extracting features only from the upper face, which may lose shape information, the masked region is transferred onto the gallery image to ensure spatial consistency. This method enhances masked face verification performance in scenarios such as security and authentication systems[5].

#### **F. Multi-scale Convolutional Self-Attention Module (MCSAM)**

Masked face recognition poses challenges due to occlusion, requiring both local and global feature extraction. Traditional convolutional networks struggle to capture long-range dependencies, which is critical for masked face recognition. To address this, the Convolutional Visual Self-Attention Network (CVSAN) integrates convolutional layers with multi-head self-attention (MHSA) to enhance feature representation. The CVSAN model consists of two main stages. The first stage, the pure convolutional stage, extracts local facial features using convolutional layers. It employs residual blocks to maintain spatial integrity while progressively reducing feature map dimensions. This stage focuses on learning detailed facial structures before integrating global dependencies. The second stage, the Conv-MHSA stage, incorporates self-attention to supplement convolution operations. MHSA captures long-range interactions by computing attention scores between different positions in the feature map. Position encoding ensures spatial awareness, allowing CVSAN to leverage both local and global information effectively. MHSA is embedded within residual blocks, balancing computational efficiency and feature expressiveness. Feature fusion combines outputs from convolutional and self-attention layers, refining representations for improved recognition performance. To optimize the network, an angular margin

loss function is used, which enhances intra-class compactness and inter-class separation by enforcing angular constraints during training. This helps CVSAN achieve better discrimination between identities.

Experiments on masked face datasets demonstrate CVSAN's effectiveness. Compared to traditional convolutional networks, it achieves superior recognition accuracy by integrating self-attention for global feature learning. The model is particularly useful for real-world applications where faces are partially occluded, such as authentication systems and security surveillance. CVSAN bridges the gap between local and global feature extraction, ensuring robust masked face recognition. A lightweight masked face recognition (MFR) network is proposed, integrating convolutional self-attention architecture to enhance accuracy while preserving clean face recognition performance with minimal computational cost. The network consists of two main branches: the whole branch and the partial branch. The whole branch employs a Multi-scale Convolutional Self-Attention Module (MCSAM) to extract both local and global facial features by leveraging CNNs and Transformers. Within the Transformer, a Multi-head Multi-scale Self-Attention (MMHSA) mechanism is introduced to enhance global feature representation across spatial and channel dimensions. The partial branch incorporates an Upper Face Alignment Module (UFAM), which includes an alignment module and MMHSA. By fusing multi-scale features of the upper face, this branch assigns higher importance to upper facial features, complementing the whole branch and improving overall MFR performance. This dual-branch design efficiently balances accuracy and computational efficiency, making it a robust solution for masked face recognition in real-world applications[6].

### III. TRADITIONAL CNN ARCHITECTURE

A Convolutional Neural Network (CNN) consists of multiple layers that work together to extract meaningful features from images and classify them effectively.

1. The **Input Layer** is responsible for receiving face images in the shape of  $224 \times 224 \times 3$  for RGB images or  $224 \times 224 \times 1$  for grayscale images. This layer normalizes pixel values, scaling them from 0-255 to a smaller range, such as 0-1 or -1 to 1, to facilitate efficient processing.

2. The first **Convolutional Layer** extracts low-level features like edges and textures. It consists of 32 filters with a  $3 \times 3$  kernel size, a stride of 1, and employs the ReLU activation function to introduce non-linearity. Batch normalization is

applied to stabilize training and accelerate convergence. The output is a set of feature maps that highlight key facial patterns.

3. To reduce computational complexity, the first **Max-Pooling Layer** downsamples the feature map while retaining important details. Using a  $2 \times 2$  pool size with a stride of 2, it outputs a condensed version of the extracted features, making the model more efficient.

4. The second **Convolutional Layer** enhances the representation of mid-level facial features such as the nose, eyes, and eyebrows. This layer includes 64 filters with a  $3 \times 3$  kernel size and ReLU activation, along with batch normalization to ensure stable weight updates. The output consists of more detailed feature maps that strengthen facial representation. Following this, the second Max-Pooling Layer further reduces the spatial dimensions while preserving critical facial features. It employs a  $2 \times 2$  pool size with a stride of 2, leading to a more compact yet effective feature map. The third Convolutional Layer is designed to extract high-level facial features, particularly focusing on masked and unmasked areas. With 128 filters and a  $3 \times 3$  kernel size, this layer captures complex facial structures such as the forehead and visible portions of the face. ReLU activation and batch normalization ensure the stability and efficiency of training. The output consists of abstracted feature maps that contribute significantly to the classification task. To refine feature selection, the third Max-Pooling Layer further reduces the size of feature maps while retaining only the most crucial details. A  $2 \times 2$  pool size with a stride of 2 ensures efficient dimensionality reduction without significant loss of information.

5. The **Flatten Layer** transforms the two-dimensional feature maps into a one-dimensional vector, making it suitable for fully connected layers. This step enables the extracted features to be passed into a classifier for identity recognition. The first Fully Connected Layer (Dense Layer) processes the extracted features to form a meaningful representation for classification. With 256 neurons and a ReLU activation function, this layer enhances the model's decision-making ability. A dropout rate of 0.5 is applied to randomly deactivate 50% of the neurons, reducing overfitting and improving generalization.

6. The final **Fully Connected Layer (Output Layer)** produces the classification results. The number of neurons in this layer corresponds to the number of identity classes. A Softmax activation function is used for multi-class classification, whereas a Sigmoid activation function is employed for binary classification. The output is a



probability distribution over the possible classes, determining the most likely identity of the masked face.

Architecture of Convolutional Neural Network showing in below figure 1.

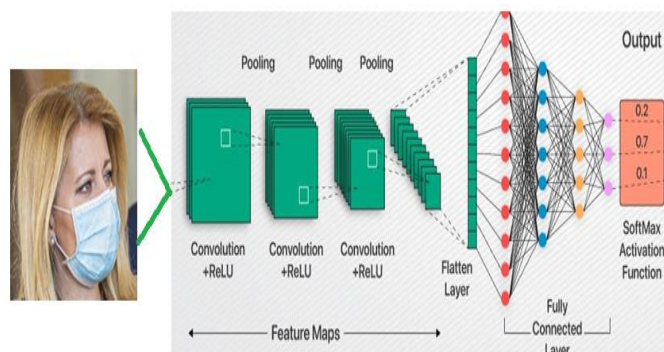


Figure 1: CNN Architecture

#### IV.RELATED DATASET

The MFR2 dataset is a small collection of real masked face images, comprising 269 images of 53 individuals, including celebrities and politicians. Each identity has an average of five images, and the dataset features diverse and unconventional mask patterns. For the masked face verification experiment, Researchers curated 800 image pairs, with 400 pairs belonging to the same identity and the remaining 400 pairs consisting of different identities.

Additionally, the LFW\_masked dataset is derived from the widely used Labeled Faces in the Wild (LFW) dataset, which is a standard benchmark for face verification. The LFW dataset contains 5,749 unique identities and a total of 13,233 face images. The LFW standard protocol, utilizing 6,000 predefined comparison pairs, where 3,000 pairs consist of images from the same identity, and the other 3,000 pairs include images of different identities.



Fig.2-MFR2 dataset on real images with mask and unmask faces

This figure 2 representing the MFR2 dataset those have two class, one is masked and another is unmasked on realistic images.

#### V.ARCHITECTURE OF THE SELF-ATTENTION MECHANISM

The architecture of the self-attention mechanism is illustrated in Figure 3, highlighting its role in enhancing masked face recognition. This mechanism allows the model to capture long-range dependencies across facial features, ensuring robust identity recognition despite occlusions.

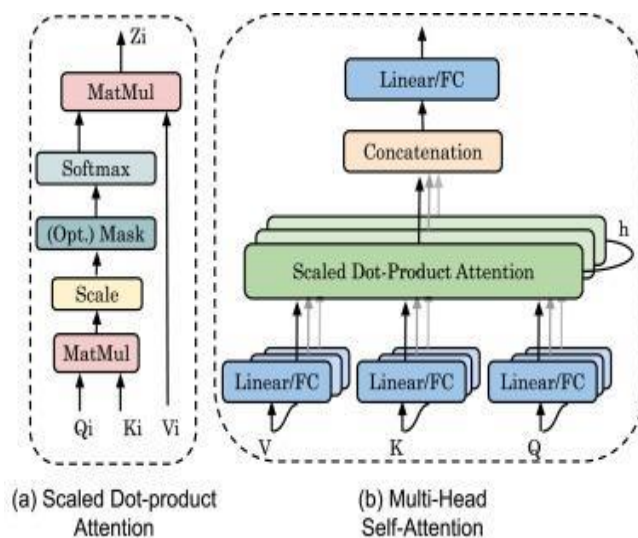


Fig 3. Architecture of the self-attention mechanism

By leveraging **Query (Q), Key (K), and Value (V) matrices**, the self-attention module selectively focuses on the most informative regions of the face, mitigating the

impact of masked areas. The integration of **Multi-Head Self-Attention (MHSA)** further refines feature extraction by attending to multiple facial attributes simultaneously.

#### **VI. MASKED FACE RECOGNITION USING SELF-ATTENTION MECHANISM**

Masked face recognition presents a significant challenge due to occlusion, where conventional CNN-based models struggle to extract identity-relevant features. To overcome this limitation, we propose a Self-Attention Mechanism that enhances both local and global feature learning, enabling the model to focus on non-occluded facial regions while maintaining robust identity recognition. Initially, feature extraction is performed using CNN, where convolutional layers capture low-level textures and mid-level facial structures such as the eyes, nose, and forehead, while max-pooling layers reduce spatial dimensions while retaining essential information. Unlike traditional CNNs that rely on local receptive fields, the Multi-Head Self-Attention (MHSA) mechanism captures long-range dependencies across the face. By computing Query (Q), Key (K), and Value (V) matrices from the CNN feature maps, the model applies Scaled Dot-Product Attention to emphasize the most informative face regions. MHSA further enhances recognition by attending to multiple facial features simultaneously, ensuring that identity-relevant details are preserved despite occlusion. Additionally, to address the issue of mask-induced occlusion, we introduce an Upper Face Alignment Module (UFAM), which prioritizes the upper facial regions—such as the eyes, forehead, and eyebrows—since these remain visible when a mask is worn. This module leverages facial landmark detection to align and extract the upper face features, ensuring the model focuses on discriminative regions unaffected by the mask.

Finally, features extracted from both the CNN and the self-attention module are fused to create a robust feature representation. These combined features are processed through fully connected layers, where a Softmax activation function classifies the identities. We compare existing method with baseline face recognition models by integrating CNN-based feature extraction with self-attention mechanisms and upper face alignment, attention mechanisms significantly enhance masked face recognition accuracy while maintaining computational efficiency, making it a viable solution for real-world applications. By integrating self-attention, the model becomes more robust to occlusions, improving recognition accuracy even when large portions of the face are covered.

#### **A. COMPREHENSIVE STUDY**

CNNs use fixed-sized kernels (e.g.,  $3 \times 3$  or  $5 \times 5$ ) that process only local features, CNNs struggle to capture global dependencies and focus too much on irrelevant regions when a mask covers a large portion of the face. Max-pooling layers reduce feature map size by selecting only the strongest activations. This operation can discard critical identity-related features, especially when masks obscure key facial regions. To overcome this issue, we will propose Self-attention mechanisms which are able to compute relationships between all facial features, dynamically focusing on visible regions, and preserving essential spatial information. Ensuring that even distant unmasked regions contribute effectively and instead of treating all regions equally, self-attention assigns higher importance to unmasked areas (e.g., eyes, forehead). Multi-Head Self-Attention (MHSA) enables the model to focus on multiple facial regions simultaneously. This ensures that different aspects of the face (e.g., shape, texture, alignment) contribute to the recognition process.

#### **VII. CONCLUSION**

Masked face recognition remains a challenging task due to occlusion, which affects traditional CNN-based models. In this study, we explored the integration of self-attention mechanisms, originally developed for natural language processing, to enhance facial feature extraction and improve recognition accuracy. Our comparative analysis demonstrated that self-attention allows the model to focus on unmasked facial regions, overcoming the limitations of CNN's local feature extraction. Experiments on the SMFRD and MFR2 datasets confirmed the effectiveness of self-attention in handling occluded faces. Future work can further optimize attention-based approaches for real-world deployment in secure authentication systems.

#### **REFERENCE**

1. Wan, Weiguo, Runlin Wen, Linghan Deng, and Yong Yang. "Masked face recognition via dual-branch convolutional self-attention network." *Applied Soft Computing* 169 (2025): 112595.
2. Wang, Yuming, Yu Li, and Hua Zou. "Masked face recognition system based on attention mechanism." *Information* 14, no. 2 (2023): 87.
3. Li, Yande, Kun Guo, Yonggang Lu, and Li Liu. "Cropping and attention based approach for masked face recognition." *Applied Intelligence* 51 (2021): 3012-3025.

4. Li, Yande, Kun Guo, Yonggang Lu, and Li Liu. "Cropping and attention based approach for masked face recognition." *Applied Intelligence* 51 (2021): 3012-3025.
5. Zhang, Meng, Rujie Liu, Daisuke Deguchi, and Hiroshi Murase. "Masked face recognition with mask transfer and self-attention under the COVID-19 pandemic." *IEEE Access* 10 (2022): 20527-20538.
6. Ge, Y., Liu, H., Du, J., Li, Z. and Wei, Y., 2023. Masked face recognition with convolutional visual self-attention network. *Neurocomputing*, 518, pp.496-506.
7. Zhang, Mengya, Yuan Zhang, and Qinghui Zhang. "Attention-Mechanism-Based Models for Unconstrained Face Recognition with Mask Occlusion." *Electronics* 12, no. 18 (2023): 3916.
8. Shahriari, Babak, Mahdi Yazdian-Dehkordi, and Ehsan Ahmadi. "Unmasking Faces: Hybrid Attention Mechanisms for Robust Masked and Unmasked Face Recognition." (2024).
9. Liu, Yang, and Wenbin Zheng. "Masked Face Recognition based on Attention Mechanism and FaceX-Zoo." In *2021 International Conference on Digital Society and Intelligent Systems (DSInS)*, pp. 107-110. IEEE, 2021.
10. Zhang, Yilin, Xiwei Peng, and Yujie Guo. "Lightweight network for masked face recognition based on improved dual attention mechanism." In *2023 IEEE International Conference on Mechatronics and Automation (ICMA)*, pp. 1621-1626. IEEE, 2023.