

PAtt-Lite: Lightweight Patch and Attention MobileNet for Challenging Facial Expression Recognition

Jia Le Ngwe, Kian Ming Lim, *Senior Member, IEEE*, Chin Poo Lee, *Senior Member, IEEE*, and Thian Song Ong, *Senior Member, IEEE*

Abstract—Facial Expression Recognition (FER) is a machine learning problem that deals with recognizing human facial expressions. While existing work has achieved performance improvements in recent years, FER in the wild and under challenging conditions remains a challenge. In this paper, a lightweight patch and attention network based on MobileNetV1, referred to as PAtt-Lite, is proposed to improve FER performance under challenging conditions. A truncated ImageNet-pre-trained MobileNetV1 is utilized as the backbone feature extractor of the proposed method. In place of the truncated layers is a patch extraction block that is proposed for extracting significant local facial features to enhance the representation from MobileNetV1, especially under challenging conditions. An attention classifier is also proposed to improve the learning of these patched feature maps from the extremely lightweight feature extractor. The experimental results on public benchmark databases proved the effectiveness of the proposed method. PAtt-Lite achieved state-of-the-art results on CK+, RAF-DB, FER2013, FERPlus, and the challenging conditions subsets for RAF-DB and FERPlus. The source code for the proposed method will be available at <https://github.com/JLREx/PAtt-Lite>.

Index Terms—Facial Expression Recognition, Patch Extraction, Self-Attention, MobileNetV1.

1 INTRODUCTION

FACIAL expression is a complex and fascinating aspect of nonverbal human communication that involves a range of facial muscle movements. These changes can convey a wide range of emotions and mental states, including happiness, sadness, anger, surprise, fear, and disgust. Given the importance of facial expressions in communication, it is not surprising that there has been a growing interest in automated facial expression recognition (FER) technology. FER has the potential to revolutionize a wide range of fields, from education to healthcare. For example, FER could be used in educational settings to measure the effectiveness and quality of teaching [1], [2], or in healthcare settings to assist in the analysis of the psychological condition of a patient [3], [4]. Along with the advances made in GPU technology, the enormous potential for downstream applications of FER also contributed to its increasing popularity.

The main challenges that FER poses differently from other image classification tasks are the inter-class similarities and intra-class differences in human facial expressions. Inter-class similarities refer to the subtle differences between facial expressions, which makes it difficult to highlight the small differences between facial expressions and recognize them correctly. On the other hand, intra-class differences, also known as subject variability, refer to the characteristic of FER databases that images from an expression class are made up of different subjects with different facial structures, gender, age, and race. This variability can hinder the learn-

ing performance of a solution, as the model may struggle to generalize across different subjects, leading to reduced accuracy and reliability. For example, the differences between an angry face and a disgusted face may be minimal, while the differences between two different individuals within the same expression class can be quite significant.

In addition, existing work has exposed other FER challenges on in-the-wild databases, namely the recognition of negative expressions, FER under challenging conditions, and reliance on large neural networks. The scarcity of negative expression images on the Internet has made it difficult to collect a representative database that can reflect real-world scenarios. Therefore, it can result in a class imbalance in the in-the-wild FER databases, which can cause the recognition rate of negative expressions to be lower than that of positive expressions. FER under challenging conditions refers to the recognition of facial expressions when the subjects are posed at certain angles or when the subject faces are partially occluded by other objects. The accurate recognition of these samples is important, especially since the challenging conditions are likely conditions identical to the downstream applications. Meanwhile, in the pursuit of classification performance, existing work is also slowly leaning towards large neural networks to achieve these performance improvements. However, considering the computing resources of downstream applications, FER methods should be readily available for these applications without requiring powerful resources.

In this paper, PAtt-Lite, a lightweight patch and attention network is proposed to improve the FER performance under challenging conditions. First, a truncated MobileNetV1 is employed as the backbone model. A patch extraction block is proposed for the truncated backbone model to enforce the

The authors are with the Faculty of Information Science and Technology, Multimedia University, Melaka, Malaysia.
 E-mail: 1181101714@student.mmu.edu.my, {kmlim, cplee, tsong}@mmu.edu.my
 The research in this work was supported by the Telekom Malaysia Research & Development under grant number RDTC/231075.

model to extract significant local facial features to classify facial expressions under challenging conditions accurately. It is designed to be lightweight while responsible for splitting the MobileNetV1 feature maps into 4 non-overlapping regions. A self-attention classifier is proposed for the backbone model to improve the learning of the output feature maps. With a dot product self-attention layer sandwiched between two fully connected layers, the attention classifier is able to learn the patched feature maps better than a vanilla classifier, hence enhancing the performance of the proposed PAtt-Lite under challenging conditions. Finally, to evaluate the performance of the proposed method, one lab-controlled database, i.e., CK+, and three in-the-wild databases, i.e., RAF-DB, FER2013, and FERPlus are employed as the benchmark databases of this research. Extensive experiments to determine the performance of the proposed method under challenging conditions such as occlusion and posed subjects using the challenging condition subsets introduced by [5] are also conducted.

The main contributions of this work are as follows:

- 1) A lightweight patch extraction block is proposed and added to the truncated MobileNetV1 to extract significant local facial features for accurately classifying facial expressions of occluded or posed subjects.
- 2) Attention classifier is proposed to relate the global average pooled output feature maps and detect their underlying patterns for better classification performance.
- 3) Extensive experimental results demonstrate the superiority of the proposed method over state-of-the-art methods on all benchmark databases, including the challenging condition subsets, despite its lightweight nature and significantly lesser parameters than state-of-the-art methods.

The remaining of this paper is organized as follows. Section 2 reviews related work with a focus on the application of the patch extraction block and the attention mechanism. Section 3 provides an overview of the architecture of the proposed PAtt-Lite, followed by detailed explanations of each module in the proposed solutions. Section 4 introduces the benchmark databases and details the experimental setting of the proposed method, along with an ablation analysis of the modules presented in the solutions, and a comparison of the proposed method with the state-of-the-art. Finally, the conclusion for this paper is included in Section 5.

2 RELATED WORK

2.1 Convolutional Neural Network

Convolutional Neural Network (CNN) is a class of deep learning models designed to process grid-like data such as images. It employs convolutional layers to automatically detect patterns or features through spatial hierarchies, enabling the network to learn progressively complex information, and pooling layers to reduce dimensionality and computational complexity, while maintaining important features. These networks have achieved significant results in computer vision tasks such as image classification, image segmentation, and object detection. The main advantage

of CNNs is their ability to learn complex features automatically without the need for manual feature engineering. Besides, CNNs are also highly adaptable to different input sizes while being able to handle complex patterns and data variations.

LeNet-5 [6] was the first CNN introduced in 1998 for the handwritten digit recognition problem, which has laid a foundation for modern CNN architecture. The CNN architecture has then gained its popularity with the introduction of AlexNet [7] in the 2012 ImageNet Large Scale Visual Recognition Challenge (ILSVRC). Since then, novel CNN architectures like the GoogLeNet [8], ResNet [9], Inception-v3 [10], MobileNet [11], and more have been proposed for the general image classification task.

With the advancement in GPU technology and the availability of mature deep learning libraries, existing work for FER has focused more on deep learning solutions recently. These solutions often outperform the handcrafted methods, especially in in-the-wild databases. While many deep learning methods exist for FER, CNN has remained the more popular choice for existing work. Most of the CNN-based methods, such as [5], [12], [13], [14], [15], [16], [17], attempt to improve the FER performance by exploiting local information in different ways with their additional modules.

The development of CNN architectures has brought forward many innovations, including residual connections [9], bottleneck design [8], [9], batch normalization [18] and its alternatives [19], [20], depthwise separable convolutions [11], and many more. However, the architectures that integrate some of these innovations are complex and have a higher number of training parameters, which in turn require larger databases and longer training time. Hence, the ImageNet-pre-trained MobileNetV1 [11] is selected as the baseline architecture for this research, due to its high performance despite its lightweight nature and relatively simpler architecture.

The base feature extractor is paired with the proposed patch extraction block as our attempt to improve FER performance under challenging conditions. The patch extraction block is designed to solely extract local features. This distinguishes it from the methods used in [12], [14], [16], which employed patch attention mechanisms. Specifically, the proposed patch extraction is inspired by that of the vision transformer architecture but remains different from its inspiration in terms of implementation details, which will be explained in the next subsection.

2.2 Vision Transformers

In recent years, the Transformer architecture has gained increasing attention in various deep learning tasks, particularly natural language processing. It is a type of neural network introduced by [21] that was designed for sequence-to-sequence tasks. The architecture consists of a stacked encoder and/or decoder layers that allow for efficient and scalable processing of large input data while learning even more complex patterns in the data. The vision transformer (ViT) architecture introduced by [22] is a novel approach that adapts the Transformer architecture to computer vision tasks. It divides images into smaller, non-overlapping patches and reshapes them into 1D sequences before processing them as a sequence using a Transformer model. The

success of this architecture has also attracted researchers' attention for the development of ViT alternatives such as DeiT [23] and Swin Transformer [24], [25]. Overall, ViT and its alternatives have achieved state-of-the-art performance on various tasks, including image classification, thus demonstrating the versatility and effectiveness of the Transformer architecture beyond natural language processing.

Hence, researchers also have begun to introduce the Transformer or the ViT architecture for FER in recent years [26], [27], [28], [29], [30], [31], [32], motivated by their performance achieved across different tasks. Based on the results posted in existing work, the application of vision transformers in FER is proven to be useful with ViT+SE [26] posting the state-of-the-art performance of 99.80% mean accuracy across 10 folds on the CK+ database, POSTER++ [30] being the best-performing method on the RAF-DB database with 92.21% accuracy, and POSTER [28] being the second best-performing method on the FERPlus database by achieving 91.62% accuracy. However, this performance often comes with large architectures with significantly more parameters than CNN-based methods. Nevertheless, the raw performance of the ViT architecture also attracted our attention to draw some inspiration for integrating into the MobileNetV1 backbone for better FER performance.

Although the patch extraction block is ultimately inspired by ViT, there exist some differences in terms of the implementation details. The first difference is the design and placement of the patch extraction block. The patch extraction mechanism in ViT is a single-layer convolution that is placed at the beginning of the architecture, while the patch extraction block in the proposed PAtt-Lite is a multi-layer convolution that is placed within the architecture. This placement allows the proposed method to fully utilize the pretrained weights of the backbone MobileNetV1, which were trained on ImageNet samples of size 224×224 . Secondly, ViT splits the input image of size 224×224 into 196 non-overlapping patches of size 16×16 , while PAtt-Lite splits the output feature maps from MobileNetV1 of size 14×14 into 4 non-overlapping patches of size 1×1 . The larger receptive regions of the proposed patch extraction block also help the proposed PAtt-Lite in extracting significant and high-level facial features.

The attention mechanism is intended to model the human attention mechanism, by highlighting parts of the input feature while ignoring the others, which enables better learning of the correlation between two input sequences. There are many variations of the attention mechanisms which sport different score functions, like additive attention [33] and dot-product attention [34]. The key component of the Transformer architecture is its self-attention, which is an attention mechanism that relates different positions of the same sequence. In this paper, an attention classifier inspired by the Transformer architecture is proposed to further improve the learning of output feature maps from the modified lightweight feature extractor. Specifically, the proposed attention classifier attempts to replicate the performance of vision transformers without requiring its series of self-attention blocks. Instead, a dot-product self-attention operation is integrated between the fully connected layers of the classifier. Through this design decision, the model can be kept lightweight while retaining high feature extractive and

classification performance for FER.

3 METHODOLOGY

3.1 Overview

Fig. 1 illustrates the overall architecture of the proposed PAtt-Lite. The proposed PAtt-Lite is built upon a truncated pre-trained MobileNetV1, combined with the proposed patch extraction block and attention classifier. Specifically, layers after the depthwise convolution of block 9 are truncated. The proposed patch extraction block and attention classifier are added to the truncated backbone model.

Given an image sample, the input will first go into the truncated MobileNetV1 to leverage the feature-extracting capability of the pre-trained model on the lower-level details of the image. The output feature maps are then used as input for our patch extraction block, where meaningful local features are extracted. The output feature maps from the patch extraction block are in the dimensions of $2 \times 2 \times D$, where D represents the depth of the feature maps. The attention classifier takes the feature maps that have been global average pooled as input, and outputs the probabilities of the facial expressions.

3.2 MobileNetV1

CNNs have been used in various computer vision tasks, such as object detection, image classification, and semantic segmentation, with state-of-the-art performance. However, CNNs can be computationally intensive and require large memory footprints, which makes them impractical for deployment on mobile or edge devices. MobileNetV1 [11] is a family of lightweight CNN architectures designed to be used on mobile and embedded devices. By leveraging depthwise separable convolutions, MobileNetV1 achieves a significant reduction in the number of model parameters and the number of multiplication and addition operations required for inference (Mult-Adds).

Depthwise separable convolution is a departure from the standard convolutional operation, as they split a standard convolution into two separate operations by performing depthwise convolutions followed by pointwise convolutions. Depthwise convolution is different from conventional convolution in that depthwise convolution applies a single convolutional filter for each input channel, whereas conventional convolution has filters that are as deep as its input. Meanwhile, pointwise convolution can be achieved using the standard convolutional operation by setting the kernel size to 1. Essentially, pointwise convolutions enable the mixing of input channels as conventional convolutions do.

As a lightweight CNN architecture, MobileNetV1 is selected as the base feature extractor in the proposed method as it provides easy finetuning of the pre-trained weights on our benchmark databases without overfitting the training samples. Mathematically, the feature extractive process from the truncated MobileNetV1 is formulated as follows:

$$X_{FE} = \text{MobileNetV1}(X) \quad (1)$$

where X is the original sample image and X_{FE} is the output feature maps from the backbone feature extractor.

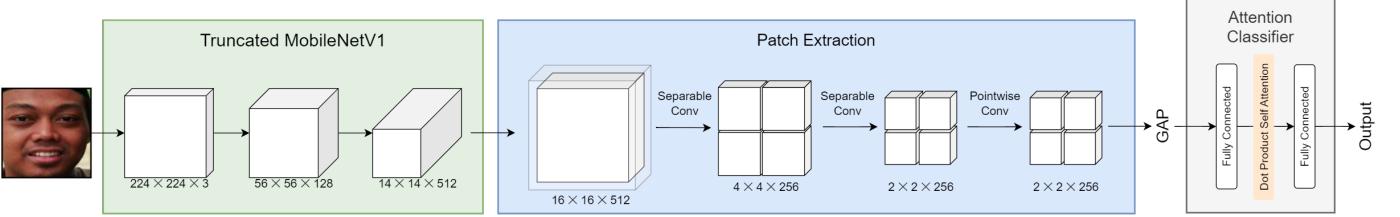


Fig. 1. Architecture of the proposed PAtt-Lite.

3.3 Patch Extraction

The key advantage of using a pre-trained CNN for transfer learning is that earlier layers have learned generic features of the training samples, such as the edges, while the later layers have learned specific features of training samples. In the context of PAtt-Lite, layers after the depthwise convolution of block 9 are skipped. The patch extraction block is added to better adapt to the FER databases than simply fine-tuning the final layers. This modification to the feature extractor also results in a shorter training period, as a higher learning rate can be used as opposed to when the pre-trained weights are being finetuned.

Our proposed patch extraction block consists of three different convolutional layers, the first two being depthwise separable convolutional layers while the last being a pointwise convolutional layer. Operating on feature maps from the MobileNetV1 that are padded to the dimension of 16×16 , the first separable convolutional layer is responsible for splitting the feature maps into four patches while learning higher-level features from its input. Subsequently, the second separable convolutional layer and the pointwise convolutional layer are responsible for learning the higher-level features from the patched feature maps, resulting in output with a dimension of 2×2 . Instead of the standard convolutional layer used in conventional CNNs, the depthwise separable convolutional layer is selected for PAtt-Lite. This design decision improves the classification performance of the proposed method on challenging subsets while reducing the number of model parameters.

3.4 Global Average Pooling

Global average pooling (GAP) is a technique that was first introduced in [35] to address the problem of overfitting in CNNs. GAP is a type of pooling operation that computes a single value for each feature map by taking the average of all the values in that map. Unlike conventional pooling operations, which reduce the spatial resolution of feature maps, GAP is normally applied at the end of a CNN architecture. The application of GAP can result in a much smaller output volume, with the output value also acting as a confidence map for each category that CNN is trained to recognize.

GAP in the proposed PAtt-Lite is responsible for averaging the patch representation from our patch extraction block, which removes the need of flattening the feature maps and feeding them into fully connected layers, hence resulting in a slight reduction in the number of parameters while further minimizing the possibility of overfitting.

Let X_{PE} be the output feature maps from the patch extraction block and \bar{X}_{PE} be the output from the GAP operation, this operation can be represented with the equation as follows:

$$\bar{X}_{PE} = \text{GAP}(X_{PE}) \quad (2)$$

3.5 Attention Classifier

An attention classifier is introduced in the proposed method for better learning of representation from the backbone MobileNetV1 and the patch extraction block. The attention classifier comprises a dot-product [34] self-attention [36] layer placed between two fully connected layers of the newly added classifier.

Dot product attention is a specific type of self-attention mechanism where the attention weights are computed as a dot product between the query vector and the key vector, divided by the square root of the dimension of the key vectors. Self-attention, also known as intra-attention in [36], is a mechanism that allows a neural network to focus on specific parts of its input during computation selectively. The idea behind self-attention is to allow the network to learn a set of attention weights that indicate how important each input element is to the output of the network. It has become a popular technique in natural language processing and computer vision tasks as it can help improve performance by selectively attending to the most relevant parts of the input.

Let Q , K , and V be the query, key, and value vectors, respectively, and $d_q = d_k$. The dot-product self-attention score can be computed as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_q}}\right)V \quad (3)$$

where d_k is the dimensionality of the key vectors. The softmax function is applied to the dot-product similarity scores to obtain a set of attention weights that sum up to 1. These weights are used to compute a weighted sum of the value vectors, resulting in the final attention output. Hence, together with the fully connected layers, the attention classifier can be represented with the following equations:

$$X_R = \text{ReLU}(\bar{X}_{PE}) \quad (4)$$

Let Q , K , and V be the query, key, and value vectors computed from the input vector X_R , the final attention output, X_A can be computed as follows:

$$X_A = \text{Attention}(Q, K, V) \quad (5)$$

$$Y = \text{softmax}(X_A) \quad (6)$$

where \tilde{X}_{PE} is the output values from GAP, X_R is the output values from the first fully connected layer with ReLU activation function, and Y represents the predicted target label as output from the final fully connected layer with softmax activation function.

4 EXPERIMENTS AND COMPARISON

4.1 Databases

Both laboratory-controlled and in-the-wild databases are used to evaluate the proposed PAtt-Lite, namely CK+, RAF-DB, FER2013, and FERPlus.

CK+ [37] is a well-known laboratory-controlled database extended from the CK database. The database consists of 593 image sequences from 123 subjects, of which 327 are labeled with one of the 7 discrete emotions: Anger, Disgust, Fear, Happy, Sadness, Surprise, and Contempt, with the first images in the sequence being the neutral expression. This research evaluates the 7 emotions CK+ with 10-fold subject-independent cross-validation to have a fair comparison with most existing work.

RAF-DB [38] is another widely used database in recent years. The database contains great variation in terms of gender, age, race, and pose of the subjects. Nearly 30,000 sample images are included in the database with crowdsourced annotations from 40 taggers. This research evaluates the basic expression subset of the database, which contains 12,271 training images and 3,068 testing images. The challenging condition test subsets of the RAF-DB database introduced by [5] are also evaluated in this research.

FER2013 [39] is introduced during the FER challenge hosted on Kaggle. It is a database collected through the Google image search API, with nearly 36,000 sample images included. The sample images are annotated with 7 basic expression labels, i.e., Angry, Disgust, Fear, Happy, Neutral, Sad, and Surprise, by 1 tagger. Compared to RAF-DB, this is a relatively more challenging database, as some of the samples are incorrectly labeled and some are without a face.

FERPlus [40] is extended from FER2013 by relabeling the original database through crowdsourcing from 10 taggers. The sample images are annotated with 8 basic expression labels, through the addition of the Contempt label. This process has corrected incorrectly labelled samples and has removed faceless samples, resulting in 35,710 sample images in the database. The challenging condition test subsets of the FERPlus database introduced by [5] are also evaluated in this research.

4.2 Implementation Details

The proposed method is implemented with the TensorFlow library on an NVIDIA TESLA P100 GPU from the Kaggle platform. A resizing operation is added to ensure that all sample images are resized to 224×224 . Random horizontal flip and random contrast are performed for data augmentation.

A two-stage training-finetuning process from [41] is employed for the training process. The pre-trained weights are frozen to solely adapt the new components to the output from MobileNetV1 during the training stage. For the finetuning process, several layers were unfrozen for finetuning

the feature extractor to the benchmark databases. To extract the best feature extractive performance from the backbone MobileNetV1, 40 layers, 59 layers, 46 layers, and 49 layers were unfrozen for CK+, RAF-DB, FER2013, and FERPlus, respectively.

We use sparse categorical cross entropy as the loss function, Adam as the optimizer, and a batch size of 8 for all experiments. For better stability of the proposed method, global gradient norm clipping is also added to the experiments. The initial learning rate is set to 1×10^{-3} for the initial training process. The learning rate is scheduled by decreasing it when the model accuracy is not improving for longer than the number of epochs that were set as patience. For the finetuning process, the learning rate is scheduled based on the inverse time decay schedule with the initial learning rate of 1×10^{-5} . The number of epochs for both the initial training and finetuning process is determined by the early stopping callback with restoration to the best weights when the training process is terminated.

4.3 Ablation Study

For the ablation study, the effectiveness of the patch extraction block and the attention classifier are evaluated by comparing them to the baseline performance of the finetuned MobileNetV1. Experiments are conducted on all benchmark databases for a proper evaluation on the effect of the proposed modules. Additionally, experimental results for the comparison between patch extraction and patch attention on in-the-wild databases are also included to justify the proposed patch extraction instead of conventional patch attention. The experimental results for this study are presented in Table 1 and Table 2.

4.3.1 Effectiveness of proposed modules

As shown in Table 1, the proposed patch extraction block is observed to have a slight decrease in performance compared to the MobileNetV1 baseline in in-the-wild databases. The performance drops in in-the-wild databases are mainly due to the difference in the number of trainable parameters between the original layers and the layers from the patch extraction block, which the newly initialized layers also resulted in the absence of pre-trained weights from the final layers. On the other hand, the MobileNetV1 baseline struggled to get a 100.00% mean accuracy on CK+ based on our experiments. However, with the small scale of CK+, this performance is immediately achievable with the addition of the patch extraction block, hence validating its effectiveness.

For the effectiveness of the proposed attention classifier, it improved the classification accuracy compared to the baseline. Moreover, the attention classifier has also significantly improved the classification accuracy on in-the-wild databases, achieving near-state-of-the-art performance on all in-the-wild databases that we benchmarked on. Specifically, the attention classifier provided an improvement of 5.83% for RAF-DB, 14.46% for FER2013, and 7.83% for FERPlus.

While the results show that the performance dropped with the patch extraction block alone, further performance improvement is achieved with both modules combined. This improvement is believed to stem from the self-attention

TABLE 1

Ablation study for the proposed method on all benchmark databases. The best result is highlighted in bold.

Attention Classifier	Patch Extraction	CK+	RAF-DB	FER2013	FERPlus
✗	✗	99.90	85.17	68.54	82.88
✗	✓	100.00	81.10	61.52	77.72
✓	✗	100.00	91.00	83.00	90.71
✓	✓	100.00	95.05	92.50	95.55

TABLE 2

Comparison between patch extraction and patch attention on in-the-wild databases. The best result is highlighted in bold.

	RAF-DB	FER2013	FERPlus
Patch Attention	94.17	89.86	92.72
Patch Extraction	95.05	92.50	95.55

layer between the fully connected layers, enabling the classifier to better adapt to the representations from the patch extraction block. The performance of the attention classifier is further boosted with the introduction of the newly initialized patch extraction block, which allowed these two modules to be trained at a higher learning rate, as opposed to the small learning rate normally used in finetuning over the pre-trained weights. Overall, the proposed patch extraction block provides a further improvement of 3.92% accuracy on RAF-DB, 2.37% on FER2013, and 4.84% on FERPlus over the MobileNetV1 with attention classifier.

4.3.2 Comparison between patch extraction and patch attention

As shown in Table 2, the conventional patch attention mechanism is not bringing any performance gain when compared to the proposed patch extraction block. Instead, experimental results have demonstrated that the proposed patch extraction block performed better in the proposed method than the patch attention across all in-the-wild databases.

While the purpose to integrate a patch attention mechanism is often to highlight significant local facial regions and hence improve the classification performance under challenging conditions, the experimental results and confusion matrices depicted in Fig. 2 - 5 have shown that the patch attention is identical to the patch extraction in terms of performance under challenging conditions. Specifically, the classification accuracy and per-class performance under challenging conditions are similar between these two designs. Both performed similarly well overall and in most classes except the Contempt class and the Disgust class, where both struggled to perform well.

4.4 Method Comparison

In this section, PAtt-Lite is compared with state-of-the-art methods on the benchmark databases. The performance comparisons are shown in Table 3 - 7. In addition, the confusion matrices of the proposed PAtt-Lite on RAF-DB, FER2013, and FERPlus are depicted in Fig. 6.

4.4.1 Results on CK+

The performance comparison of the proposed method with the state-of-the-art methods for CK+ is presented in Table 3.

TABLE 3

Comparison of the state-of-the-art results on CK+. The best result is highlighted in bold.

Methods	Accuracy
gACNN [14]	96.40
pACNN [14]	97.03
SCAN-CCI [15]	97.31
IF-GAN [42]	97.52
FDRL [43]	99.54
ViT + SE [26]	99.80
PAtt-Lite	100.00

TABLE 4

Comparison of the state-of-the-art results on RAF-DB and FERPlus. The best result is highlighted in bold.

Methods	# Params	RAF-DB	FERPlus
VTFF [31]	80.1M	-	88.81
RAN [5]	11.2M	86.90	89.16
VTFF [31]	51.8M	88.14	-
SCAN-CCI [15]	70M	89.02	89.42
ARM [44]	11.2M	90.42	-
TransFER [27]	65.2M	90.91	90.83
Facial Chirality [45]	46.2M	91.20	-
APViT [29]	65.2M	91.98	90.86
POSTER [28]	71.8M	92.05	91.62
POSTER++ [30]	43.7M	92.21	-
CIAO [46]	17.9M	-	94.50
PAtt-Lite	1.10M	95.05	95.55

The proposed PAtt-Lite outperformed all CNN-based existing work [14], [15], [42], [43] and transformer-based existing work [26] in terms of cross-validation mean accuracy for CK+ by achieving 100.00% mean accuracy across 10-fold cross-validation. To the best of our knowledge, this is the best performance reported for a method that is trained without facial landmarks.

4.4.2 Results on RAF-DB.

The performance comparison of the proposed PAtt-Lite with state-of-the-art methods on RAF-DB is shown in Table 4. An accuracy and parameter comparison between the proposed method and state-of-the-art methods is also presented in Fig. 7. VTFF [31], TransFER [27], Facial Chirality [45], APViT [29], POSTER [28], and POSTER++ [30] are transformer-based architecture, while RAN [5], SCAN-CCI [15], and ARM [44] are CNN-based architecture.

Based on the comparison, the transformer-based methods are generally outperforming and have a greater number of parameters than the CNN-based methods, with SCAN-CCI [15] being the exception as it has 70M parameters despite having a CNN-based architecture. Our proposed PAtt-Lite achieved better performance with a CNN backbone than all transformer-based state-of-the-art. Specifically, PAtt-Lite has 2.84% over transformer-based POSTER++ [30],

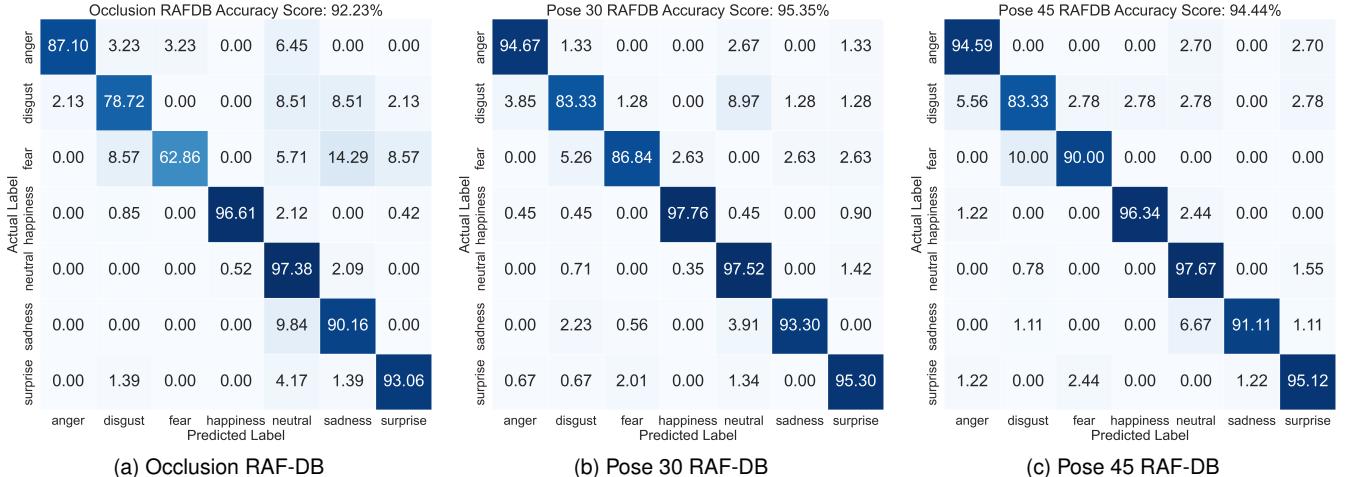


Fig. 2. Confusion matrices of patch extraction on challenging subsets of RAF-DB.

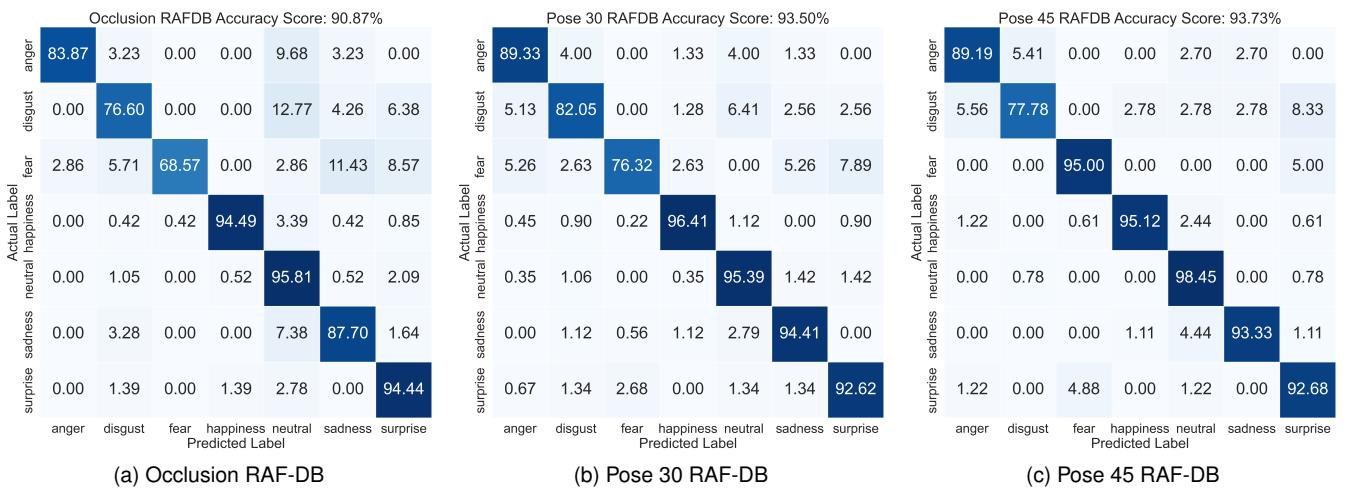


Fig. 3. Confusion matrices of patch attention on challenging subsets of RAF-DB.

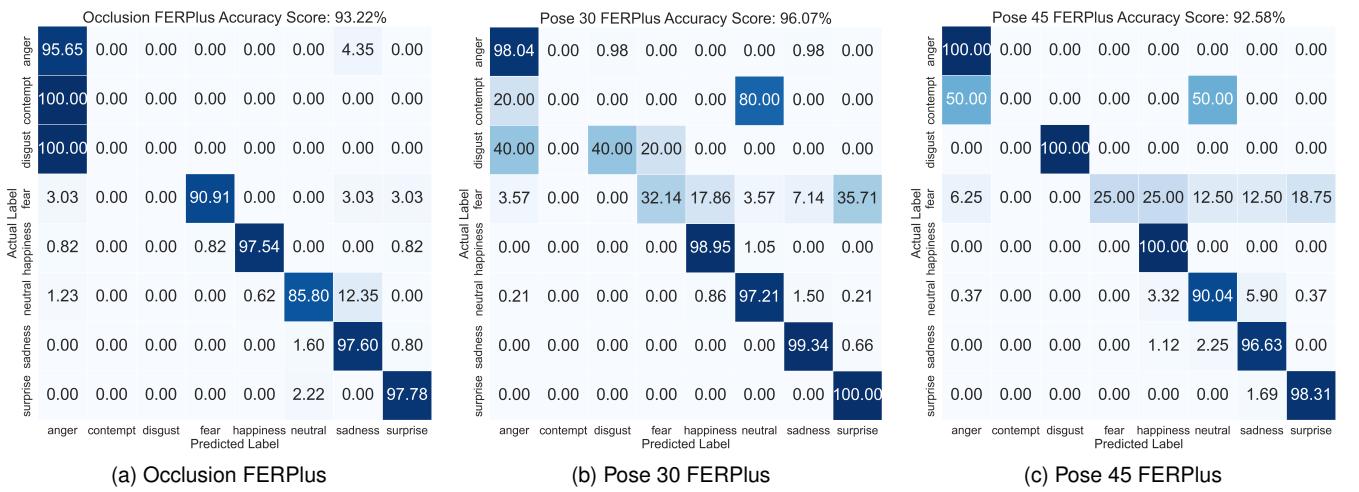


Fig. 4. Confusion matrices of patch extraction on challenging subsets of FERPlus.

which despite being the lightest transformer-based method, recorded the best result for RAF-DB. For comparison with CNN-based methods, PAtt-Lite achieved an improvement

of 4.63% overall accuracy over ARM [44]. In terms of performance comparison with lightweight methods, our proposed method outperformed RAN [5] and ARM [44] by 8.15% and

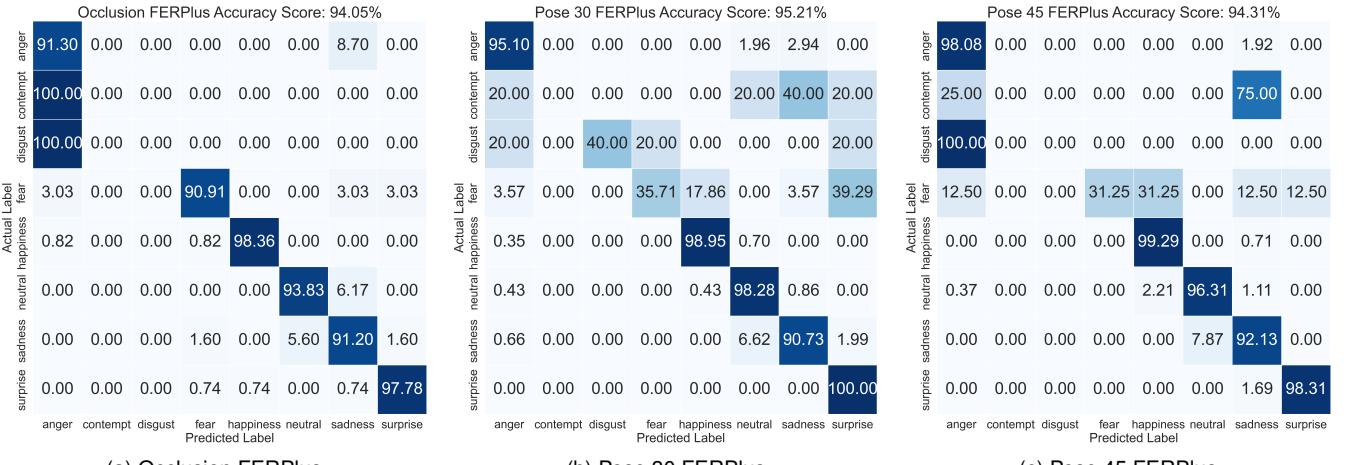


Fig. 5. Confusion matrices of patch extraction on challenging subsets of FERPlus.

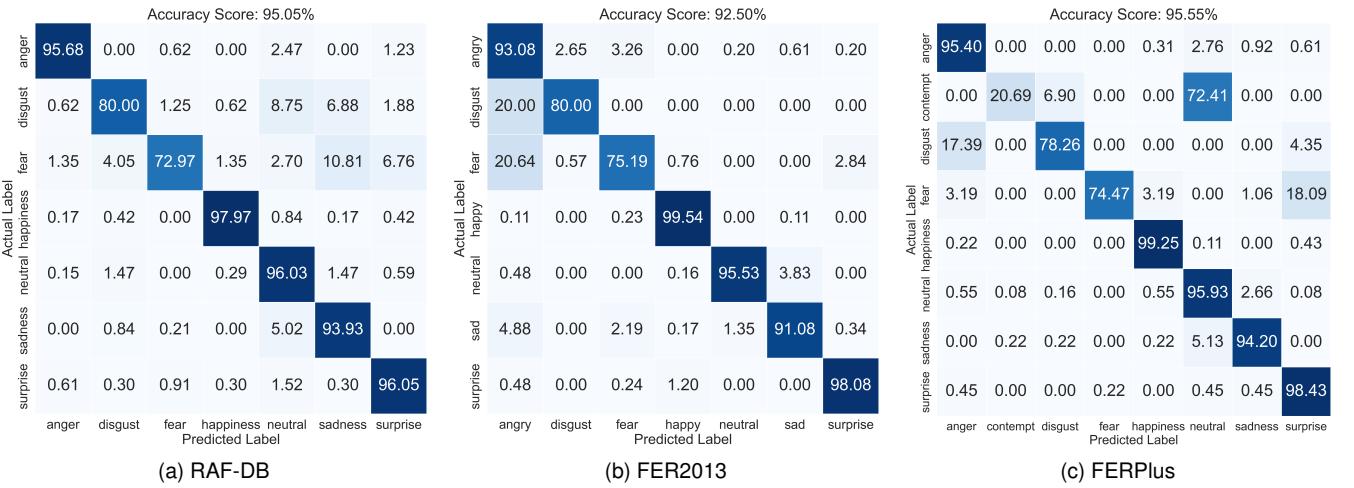


Fig. 6. Confusion matrices of PAtt-Lite on in-the-wild databases.

4.63% respectively, while having 10 times lesser parameters.

Furthermore, following [28], [29], [30], [31], a per-class performance comparison is added in Table 5 to compare the performance of PAtt-Lite on different classes of the database. For [15], [16], [27], [32], [44], which did not specifically report their per-class accuracy, their results from the confusion matrices are taken for comparison. Overall, a similar trend for the per-class performance is observed when PAtt-Lite is compared with existing work, where the per-class accuracy for the Disgust class and the Fear class is lacking behind the other classes. However, our method managed to slightly improve the per-class accuracy for the Disgust class over APViT [29] and perform on par with it for the Fear class while having 59 times fewer parameters. Overall, small improvements can be observed across all expression classes except class Anger and class Neutral, where significant improvement is recorded for the former and a small performance drop is noted for the latter. Noticeably, the proposed method managed to record a 6.79% improvement over all other state-of-the-art methods for class Anger, while recording a 1.87% performance drop over ARM [44]. Although the per-class accuracy for the Disgust class and the

Fear class can be further improved, PAtt-Lite demonstrated consistently strong performance by managing to achieve about 93% per-class accuracy or greater in the remaining classes, resulting in an average accuracy of 90.38%. This corresponds to an improvement in the average accuracy of 4.02% over APViT [29], which previously reported the best average accuracy. Moreover, when compared to Imponderous Net [16], which has the nearest number of parameters as our proposed method, PAtt-Lite achieved an improvement of 12.81% in terms of average accuracy, while performing stronger and more consistent on all expression classes.

4.4.3 Results on FER2013

The comparison of the proposed method with the state-of-the-art methods on FER2013 is depicted in Table 6. FER2013 has been a database that existing work struggled to perform well on until recently when FLEPNet [51] and NECM-PECM Ensemble [52] were proposed. Compared to the existing work, PAtt-Lite recorded a classification performance of 92.50%, which equates to an improvement of 4.5% over NECM-PECM Ensemble [52].

TABLE 5
Per-class performance comparison of the state-of-the-art results on RAF-DB. The best result is highlighted in bold.

Method	# Params	Anger	Disgust	Fear	Happiness	Neutral	Sadness	Surprise	Average
Imponderous Net [16]	1.45M	78.00	54.00	57.00	96.00	88.00	85.00	85.00	77.57
MViT [32]	33M	78.40	63.75	60.81	95.61	89.12	87.45	87.54	80.38
VTFF [31]	51.8M	85.80	68.12	64.86	94.09	87.50	87.24	85.41	81.86
SCAN-CCI [15]	70M	81.00	70.00	66.00	96.00	89.00	86.00	88.00	82.29
ARM [44]	11.2M	77.20	64.40	70.30	95.40	97.90	83.90	90.30	82.77
TransFER [27]	65.2M	88.89	79.37	68.92	95.95	90.15	88.70	89.06	85.86
POSTER++ [30]	43.7M	88.27	71.88	68.92	97.22	92.06	92.89	90.58	85.97
POSTER [28]	71.8M	88.89	75.00	67.57	96.96	92.35	91.21	90.27	86.04
APViT [29]	65.2M	86.42	73.75	72.97	97.30	92.06	88.70	93.31	86.36
PAtt-Lite	1.10M	95.68	80.00	72.97	97.97	96.03	93.93	96.05	90.38

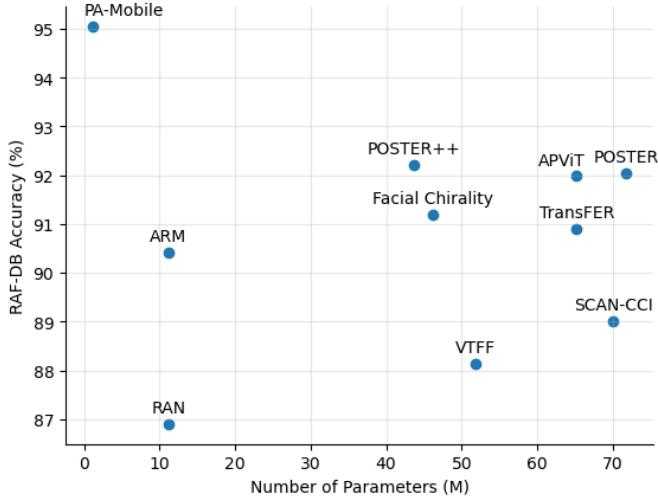


Fig. 7. Accuracy and parameter comparison on RAF-DB.

TABLE 6
Comparison of the state-of-the-art results on FER2013. The best result is highlighted in bold.

Methods	Accuracy
KDL [47]	71.28
MAFER [48]	73.45
PASM [49]	73.59
MoVE-CNNs [50]	77.70
FLEPNet [51]	80.72
NECM-PECM Ensemble [52]	88.00
PAtt-Lite	92.50

From the confusion matrix in Fig 6b, the per-class performance of our proposed method on FER2013 is observed to have a similar trend to RAF-DB and existing work [48], [51], where generally strong per-class accuracy can be achieved for all expression classes except the Disgust class and the Fear class, which are relatively weaker compared to the remaining classes. However, PAtt-Lite has managed to achieve significantly better performance across all expression classes. Specifically, the proposed method improves the per-class accuracy of all expressions except Disgust and Fear to more than 90% accuracy. This is unlike most existing work, which has struggled to achieve such performance on all expressions except the Happy class.

4.4.4 Results on FERPlus

FERPlus has an additional Contempt expression than the other in-the-wild databases in this research. The classification performance of our proposed method is compared with the state-of-the-art methods in the last column of Table 4. The accuracy and parameter of the proposed PAtt-Lite are also compared with state-of-the-art methods in Fig. 8. Like RAF-DB, transformer-based methods generally have better results than CNN-based methods, except for CIAO [46], which reported the best accuracy for FERPlus with a CNN backbone.

The proposed PAtt-Lite has achieved slight improvement on FERPlus when compared to the state-of-the-art methods. Specifically, PAtt-Lite achieved a classification accuracy of 95.55%, which corresponds to a performance improvement of 1.05% over CIAO [46], while having 16 times lesser parameters than CIAO [46]. When compared to transformer-based methods, our method achieved a 3.93% improvement in classification accuracy over POSTER [28] with 65 times lesser parameters. For comparison with lighter state-of-the-art methods, PAtt-Lite also outperformed RAN [5] by 6.39% in terms of overall accuracy with 10 times lesser parameters.

It is visible that our proposed method is still struggling with the Contempt class upon inspecting the confusion matrix in Fig. 6c, with the Disgust class and the Fear class following a general trend from RAF-DB and FER2013. However, this performance drop is expected due to the severe class imbalance between these classes and the Neutral class. Similar performance drops on some of these negative expression classes are also reported in existing work [15], [28], [29]. Despite this, our proposed method still achieved significant improvements in all other expression classes, with around or higher than 95% performance across these 5 classes.

4.4.5 Results on challenging subsets

The proposed method has shown great results in classifying samples under challenging conditions. Table 7 shows the difference in the performance of the proposed method compared to recent work on challenging subsets for RAF-DB and FERPlus. Specifically, PAtt-Lite achieved an improvement of more than 3% across all three challenging subsets for RAF-DB, outperforming Facial Chirality [45] which reported the best results on these subsets with 46.2M parameters. Meanwhile, to our best knowledge, our proposed method is the first to achieve more than 90% accuracy for all

TABLE 7

Comparison of the state-of-the-art results on challenging subsets for RAF-DB and FERPlus. The best result is highlighted in bold.

Methods	# Params	RAF-DB			FERPlus		
		Occlusion	Pose > 30	Pose > 45	Occlusion	Pose > 30	Pose > 45
RAN [5]	11.2M	82.72	86.74	85.20	83.63	82.23	80.40
Imponderous Net [16]	1.45M	83.40	86.12	84.41	83.47	86.84	84.83
VTFF [31]	51.8M	83.95	87.97	88.35	-	-	-
VTFF [31]	80.1M	-	-	-	84.79	88.29	87.20
OADN [17]	-	-	-	-	84.57	88.52	87.50
SCAN-CCI [15]	70M	85.03	89.82	89.07	86.12	88.89	88.15
MViT [32]	33M	85.17	87.99	88.40	-	-	-
Facial Chirality [45]	46.2M	88.16	91.50	90.86	-	-	-
PAtt-Lite	1.10M	92.23	95.35	94.44	93.22	96.07	92.58

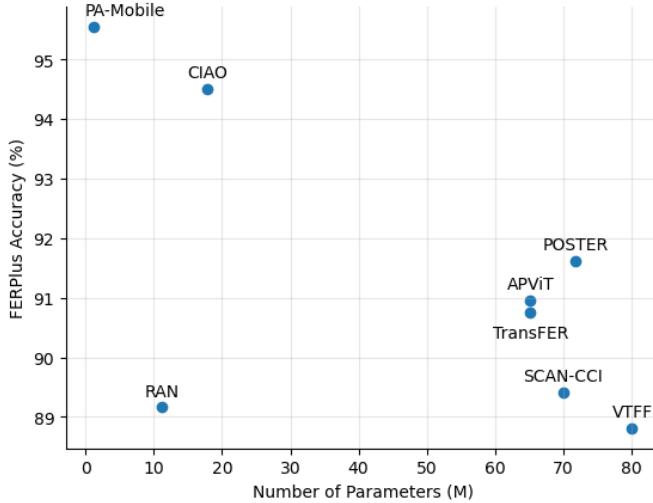


Fig. 8. Accuracy and parameter comparison on FERPlus.

challenging subsets for FERPlus. A performance improvement of around 4% to 7% is achieved across the three subsets compared to SCAN-CCI [15], which previously reported the best performance with 70M parameters. From the comparison with the best results on these subsets, the proposed PAtt-Lite achieved state-of-the-art performance across all subsets with significantly lesser parameters. Meanwhile, when compared with the lighter methods such as RAN [5] and Imponderous Net [16] on subsets of RAF-DB, PAtt-Lite managed to outperform these methods by around 9%. The proposed method also managed to outperform the lighter methods by more than 9% on the first two subsets and by 7.75% on the Pose 45 subset of FERPlus.

5 CONCLUSION

This work presents PAtt-Lite, a MobileNetV1-based solution, to improve the classification accuracy of FER under challenging conditions. The proposed PAtt-Lite achieves state-of-the-art performance in all benchmark databases and their subsets while being significantly lighter than other state-of-the-art methods at just 1.10M parameters. Specifically, the proposed patch extraction block improves the FER performance of PAtt-Lite under challenging conditions by enforcing the model to extract significant local facial features. On the other hand, the attention classifier is proposed

to learn the patched representation better and improve the overall performance of the proposed lightweight method.

This work provides valuable insight into potential future directions by highlighting the advantages and possible improvements of our proposed method. One such direction is robustness improvement to PAtt-Lite, especially towards low-resource expressions, which can be done by further refining the patch extraction blocks. We also suggest that there could be an effort for the development of a FER database with more reliable annotations. Existing FER methods have advanced to a point where they can potentially discover patterns that are difficult for human annotators to detect. Therefore, a crowdsourced database that incorporates insights from state-of-the-art FER methods could be a valuable resource for further research in this field.

Overall, this work has shown the potential to use MobileNetV1 as a baseline feature extractor in FER. The need for continued research and development is also highlighted to further improve the accuracy and reliability of automated facial expression recognition, especially under challenging conditions and low-resource expressions.

REFERENCES

- [1] Y. Wu and L. Shen, "An adaptive landmark-based attention network for students' facial expression recognition," in *2021 6th International Conference on Communication, Image and Signal Processing, CCISP 2021*, pp. 139–144, Institute of Electrical and Electronics Engineers Inc., 2021.
- [2] X. Li, R. Yue, W. Jia, H. Wang, and Y. Zheng, "Recognizing students' emotions based on facial expression analysis," in *11th International Conference on Information Technology in Medicine and Education, ITME 2021*, pp. 96–100, Institute of Electrical and Electronics Engineers Inc., 2021.
- [3] C. J. Meryl, K. Dharshini, D. S. Juliet, J. A. Rosy, and S. S. Jacob, "Deep learning based facial expression recognition for psychological health analysis," *Proceedings of the 2020 IEEE International Conference on Communication and Signal Processing, ICCSP 2020*, pp. 1155–1158, 7 2020.
- [4] J. Ye, Y. Yu, G. Fu, Y. Zheng, Y. Liu, Y. Zhu, and Q. Wang, "Analysis and recognition of voluntary facial expression mimicry based on depressed patients," *IEEE Journal of Biomedical and Health Informatics*, 2023.
- [5] K. Wang, X. Peng, J. Yang, D. Meng, and Y. Qiao, "Region attention networks for pose and occlusion robust facial expression recognition," *IEEE Transactions on Image Processing*, vol. 29, pp. 4057–4069, 2020.
- [6] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.

- [8] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [10] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- [11] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilennets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [12] S. Xie and H. Hu, "Facial expression recognition using hierarchical features with deep comprehensive multipatches aggregation convolutional neural networks," *IEEE Transactions on Multimedia*, vol. 21, no. 1, pp. 211–220, 2018.
- [13] S. Zhao, H. Cai, H. Liu, J. Zhang, and S. Chen, "Feature selection mechanism in cnns for facial expression recognition," in *BMVC*, p. 317, 2018.
- [14] Y. Li, J. Zeng, S. Shan, and X. Chen, "Occlusion aware facial expression recognition using cnn with attention mechanism," *IEEE Transactions on Image Processing*, vol. 28, no. 5, pp. 2439–2450, 2018.
- [15] D. Gera and S. Balasubramanian, "Landmark guidance independent spatio-channel attention and complementary context information based facial expression recognition," *Pattern Recognition Letters*, vol. 145, pp. 58–66, 2021.
- [16] D. Gera and S. Balasubramanian, "Imponderous net for facial expression recognition in the wild," *arXiv preprint arXiv:2103.15136*, 2021.
- [17] H. Ding, P. Zhou, and R. Chellappa, "Occlusion-adaptive deep network for robust facial expression recognition," in *2020 IEEE International Joint Conference on Biometrics (IJCB)*, pp. 1–9, IEEE, 2020.
- [18] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*, pp. 448–456, pmlr, 2015.
- [19] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [20] Y. Wu and K. He, "Group normalization," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018.
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [22] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [23] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training data-efficient image transformers & distillation through attention," in *International conference on machine learning*, pp. 10347–10357, PMLR, 2021.
- [24] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- [25] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong, et al., "Swin transformer v2: Scaling up capacity and resolution," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12009–12019, 2022.
- [26] M. Aouayeb, W. Hamidouche, C. Soladie, K. Kpalma, and R. Seguier, "Learning vision transformer with squeeze and excitation for facial expression recognition," *arXiv preprint arXiv:2107.03107*, 2021.
- [27] F. Xue, Q. Wang, and G. Guo, "Transfer: Learning relation-aware facial expression representations with transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3601–3610, 2021.
- [28] C. Zheng, M. Mendieta, and C. Chen, "Poster: A pyramid cross-fusion transformer network for facial expression recognition," *arXiv preprint arXiv:2204.04083*, 2022.
- [29] F. Xue, Q. Wang, Z. Tan, Z. Ma, and G. Guo, "Vision transformer with attentive pooling for robust facial expression recognition," *IEEE Transactions on Affective Computing*, 2022.
- [30] J. Mao, R. Xu, X. Yin, Y. Chang, B. Nie, and A. Huang, "Poster++: A simpler and stronger facial expression recognition network," *arXiv preprint arXiv:2301.12149*, 2023.
- [31] F. Ma, B. Sun, and S. Li, "Facial expression recognition with visual transformers and attentional selective fusion," *IEEE Transactions on Affective Computing*, pp. 1–1, 2021.
- [32] H. Li, M. Sui, F. Zhao, Z. Zha, and F. Wu, "Mvt: mask vision transformer for facial expression recognition in the wild," *arXiv preprint arXiv:2106.04520*, 2021.
- [33] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [34] M.-T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," *arXiv preprint arXiv:1508.04025*, 2015.
- [35] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.
- [36] J. Cheng, L. Dong, and M. Lapata, "Long short-term memory-networks for machine reading," *arXiv preprint arXiv:1601.06733*, 2016.
- [37] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, "The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression," in *2010 ieee computer society conference on computer vision and pattern recognition-workshops*, pp. 94–101, IEEE, 2010.
- [38] S. Li, W. Deng, and J. Du, "Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild," in *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pp. 2584–2593, IEEE, 2017.
- [39] I. J. Goodfellow, D. Erhan, P. L. Carrier, A. Courville, M. Mirza, B. Hammer, W. Cukierski, Y. Tang, D. Thaler, D.-H. Lee, et al., "Challenges in representation learning: A report on three machine learning contests," in *International conference on neural information processing*, pp. 117–124, Springer, 2013.
- [40] E. Barsoum, C. Zhang, C. C. Ferrer, and Z. Zhang, "Training deep networks for facial expression recognition with crowd-sourced label distribution," in *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pp. 279–283, 2016.
- [41] "Transfer learning and fine-tuning — TensorFlow Core."
- [42] J. Cai, Z. Meng, A. S. Khan, J. O'Reilly, Z. Li, S. Han, and Y. Tong, "Identity-free facial expression recognition using conditional generative adversarial network," in *2021 IEEE International Conference on Image Processing (ICIP)*, pp. 1344–1348, IEEE, 2021.
- [43] D. Ruan, Y. Yan, S. Lai, Z. Chai, C. Shen, and H. Wang, "Feature decomposition and reconstruction learning for effective facial expression recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7660–7669, 2021.
- [44] J. Shi, S. Zhu, and Z. Liang, "Learning to amend facial expression representation via de-albino and affinity," *arXiv preprint arXiv:2103.10189*, 2021.
- [45] L. Lo, H. Xie, H.-H. Shuai, and W.-H. Cheng, "Facial chirality: From visual self-reflection to robust facial feature learning," *IEEE Transactions on Multimedia*, vol. 24, pp. 4275–4284, 2022.
- [46] P. Barros and A. Sciutti, "Ciao! a contrastive adaptation mechanism for non-universal facial expression recognition," in *2022 10th International Conference on Affective Computing and Intelligent Interaction (ACII)*, pp. 1–8, IEEE, 2022.
- [47] M. A. Mahmoudi, A. Chetouani, F. Boufara, and H. Tabia, "Kernelized dense layers for facial expression recognition," in *2020 IEEE International Conference on Image Processing (ICIP)*, pp. 2226–2230, IEEE, 2020.
- [48] F. V. Massoli, D. Cafarelli, C. Gennaro, G. Amato, and F. Falchi, "Mafer: A multi-resolution approach to facial expression recognition," *arXiv preprint arXiv:2105.02481*, 2021.
- [49] P. Liu, Y. Lin, Z. Meng, L. Lu, W. Deng, J. T. Zhou, and Y. Yang, "Point adversarial self-mining: A simple method for facial expression recognition," *IEEE Transactions on Cybernetics*, 2021.
- [50] J. X. Yu, K. M. Lim, and C. P. Lee, "Move-cnns: Model averaging ensemble of convolutional neural networks for facial expression recognition," *IAENG International Journal of Computer Science*, vol. 48, no. 3, 2021.
- [51] M. Karnati, A. Seal, A. Yazidi, and O. Krejcar, "Fleynet: Feature level ensemble parallel network for facial expression recognition," *IEEE Transactions on Affective Computing*, vol. 13, no. 4, pp. 2058–2070, 2022.
- [52] P. Phattarasooksirot and A. Sento, "Facial emotional expression recognition using hybrid deep learning algorithm," in *2022 7th*

