

Masked Facial Expression Recognition with Attention Mechanism based on Mini_Xception Network

Zhiping Zhang 12*

1 School of Computing Sciences, College of Computing, Informatics and Media, Universiti Teknologi MARA, Selangor, Malaysia

2 College of Computer and Mathematics, Xinyu University, Jiangxi, P. R. China, 2021422338@student.uitm.edu.my

Shuzlina Abdul-Rahman

School of Computing Sciences, College of Computing, Informatics and Media, Universiti Teknologi MARA, Selangor, Malaysia, shuzlina@uitm.edu.my

Norlina Mohd Sabri

Universiti Teknologi MARA Cawangan Terengganu, Kuala Terengganu, Malaysia, norli097@uitm.edu.my

Abstract: Facial Expression Recognition (FER) plays a crucial role in applications such as human-computer interaction and health monitoring. However, traditional FER methods often suffer significant performance degradation when occlusions are present, especially when the mouth region is covered, as commonly observed during the COVID-19 pandemic. To address this challenge, we propose a novel FER model, Mini-Xception+CBAM, which combines a lightweight network architecture with an attention mechanism to enhance performance under occlusion scenarios. The proposed approach utilizes the Mini-Xception architecture and integrates the Convolutional Block Attention Module (CBAM), enabling the model to dynamically focus on critical facial regions, such as the eyes and eyebrows, which are less likely to be occluded. Extensive experiments on the FER2013, SMFER2013, and MFER2013 datasets demonstrate that our model outperforms several state-of-the-art baseline models, including Mini-Xception, MobileNetV2, ResNet-18, and Vision Transformer (ViT), achieving 90.4% accuracy on FER2013, 89.2% on SMFER2013, and 91.5% on MFER2013, along with superior precision, recall, and F1 score. Furthermore, the model's efficient design, with only 0.68M parameters and 30.1M FLOPs, makes it highly suitable for deployment on resource-constrained devices.

CCS CONCEPTS

Computing methodologies~Artificial intelligence~Computer vision

Keywords: Spatial attention mechanism; Mini-Xception; Mask occlusion; Expression recognition

ORCID

ZhipingZhang <https://orcid.org/0009-0007-2250-3628>

Shuzlina Abdul-Rahman <https://orcid.org/0000-0002-5498-9606>

Norlina Mohd Sabri <https://orcid.org/0000-0002-2470-8172>

1.Introduction

Facial expression recognition (FER) has been a critical area of research in computer vision due to its applications in diverse fields such as healthcare, surveillance, human–computer interaction, and affective computing. The goal of FER systems is to automatically classify human emotions based on facial expressions, which are one of the most powerful non–verbal communication tools. Despite significant advancements in FER methods, the COVID–19 pandemic has introduced a unique challenge: facial masks. These masks obscure a substantial portion of facial features, particularly the lower face, complicating the task of emotion recognition[1].

Occlusions caused by facial masks lead to significant information loss, as traditional FER models often rely on visible regions of the face, such as the mouth, to recognize emotions effectively. Previous studies have demonstrated that certain facial regions, such as the mouth, are more critical than others in conveying emotions like happiness and surprise[2]. The occlusion of these regions disrupts the spatial and semantic continuity of facial features, resulting in degraded recognition performance.

Additionally, existing FER datasets, such as FER2013, RAF–DB, and AffectNet[3], predominantly consist of unmasked faces, limiting their applicability to masked scenarios. Although augmentation techniques, such as mask simulation, have been employed to bridge this gap, these methods may not accurately replicate real–world occlusions, leading to models that fail to generalize effectively.

In practical applications, FER systems are often deployed on devices with limited computational resources, such as smartphones, embedded systems, and wearable devices. This necessitates the development of lightweight architectures that maintain high recognition accuracy while minimizing computational overhead. Among such architectures, Mini–Xception has emerged as a promising candidate due to its compact design and efficiency, making it suitable for real–time emotion recognition tasks. However, the performance of Mini–Xception in handling occluded facial expressions remains underexplored.

Recent advancements in attention mechanisms, such as the Convolutional Block Attention Module (CBAM) and the Squeeze–and–Excitation Network, have shown great potential in improving the robustness of convolutional neural networks. These mechanisms enhance feature extraction by dynamically emphasizing salient regions while suppressing irrelevant or noisy features. In the context of masked FER, attention mechanisms can be leveraged to prioritize unoccluded facial regions, such as the eyes and forehead, which are critical for recognizing emotions in the presence of masks.

This study proposes a novel approach that integrates Mini–Xception with attention mechanisms to address these challenges. By incorporating mechanisms like CBAM, the framework aims to enhance the spatial and channel–wise focus of Mini–Xception, enabling it to dynamically adapt to occluded facial inputs. This research not only addresses a critical gap in masked FER but also contributes to the broader field of FER by demonstrating the potential of lightweight and attention–enhanced models for real–world applications.

In summary, this study builds on the foundations of FER research, leveraging advancements in lightweight CNNs and attention mechanisms to tackle the unique challenges posed by facial masks. The findings are expected to provide significant insights into the design of efficient and robust FER systems for masked environments.

2.RELATED WORK

2.1 Deep convolutional FER

Deep Convolutional Neural Networks (CNNs) have become the cornerstone of modern FER due to their ability to learn hierarchical features directly from data. Unlike traditional approaches that rely on handcrafted features, CNNs automate the feature extraction process, significantly improving performance on complex and diverse datasets[4].

As one of the first deep learning models applied to FER, AlexNet showcased the potential of CNNs for feature extraction and classification. Its success on datasets like FER2013 marked a paradigm shift in FER research. VGGNet improved upon AlexNet by introducing deeper architectures with smaller convolutional filters, achieving better feature representation and classification performance on FER datasets. However, its high computational cost limited its applicability in real–time scenarios.

ResNet introduced residual connections to alleviate the vanishing gradient problem, enabling deeper and more robust networks for FER. Studies have shown that ResNet outperforms shallower architectures, particularly on complex FER datasets like AffectNet. Inception models improved feature extraction by employing multi-scale convolutions within a single layer. The enhanced capacity to capture diverse spatial information made them effective for recognizing subtle facial expressions. Based on depthwise separable convolutions, Xception offers a computationally efficient alternative to standard convolutional layers while maintaining high accuracy[5]. Studies have successfully adapted Xception for FER tasks, demonstrating its suitability for large-scale datasets.

While deep CNNs achieve remarkable accuracy, their computational complexity poses challenges for real-time and resource-constrained applications. Lightweight architectures, designed for efficiency, have gained attention in FER research. MobileNet utilizes depthwise separable convolutions to reduce the number of parameters, enabling real-time FER on mobile and embedded devices. ShuffleNet further optimizes computation through channel shuffling, making it suitable for FER tasks with minimal resource requirements. A compact variant of the Xception model, Mini-Xception[6] balances computational efficiency and accuracy, making it a popular choice for FER, particularly in real-time applications.

In Summary, Deep convolutional FER has evolved from foundational architectures like AlexNet to sophisticated, lightweight models with attention mechanisms. These advancements address key challenges in FER, including occlusion robustness and computational efficiency. However, the need for robust solutions in masked environments remains pressing, this study focus on Mini-Xception with attention mechanisms to advance FER under occluded conditions.

2.2 Lightweight FER

Lightweight FER aims to achieve high accuracy in recognizing emotions while minimizing computational complexity, making it suitable for resource-constrained environments such as mobile and embedded devices.

Howard et al. introduced MobileNets, demonstrating their effectiveness in balancing accuracy and computational efficiency through depthwise separable convolutions. They have been adapted for FER tasks to process emotions on mobile devices with reduced latency. Tan and Le proposed EfficientNet, which uses a compound scaling method to adjust depth, width, and resolution systematically. Studies have shown its adaptability for FER applications, achieving high accuracy with fewer resources. Chen et al.[7]investigated lightweight versions of classic architectures like VGG and ResNet, which reduce computational demands through fewer parameters while retaining FER performance.

Zhang et al.[8] distilled large models like ResNet into smaller networks, achieving comparable results on FER datasets with a significant reduction in computational cost. Han et al. introduced pruning techniques to remove redundant weights, with subsequent applications in FER by Kim et al. where pruned networks achieved real-time emotion recognition on low-power devices. Jacob et al. demonstrated that quantized models, such as those converted to 8-bit precision, preserve FER accuracy while reducing memory and computational demands, especially on edge devices.

Dalal and Triggs showed the efficiency of HOG in detecting patterns in images. Zhang et al.extended this to FER by using HOG for initial feature extraction, followed by lightweight classifiers for emotion recognition.This approach leveraged cross-modal data for more accurate and efficient emotion recognition.

In conclusion, the above strategies demonstrate that lightweight FER is achievable through a combination of efficient architectures, optimization techniques, and hybrid methodologies. These advancements are supported by a growing body of literature, highlighting their potential for deployment in mobile and embedded environments.

2.3 Masked FER

Masked FER has emerged as a critical focus in computer vision and affective computing, driven by challenges such as the COVID-19 pandemic, where face masks and other occlusions obscure traditional FER systems' key focus areas (eyes, mouth, and forehead).

Since the eye region often remains visible despite occlusions, leveraging it has become a natural adaptation.Chen et al.[9] introduced a deep learning approach centered on the eye region using a spatial attention mechanism. The model achieved improved accuracy for

masked FER by dynamically focusing on the most relevant facial features, underscoring the eyes' critical role in expression analysis. Multimodal strategies integrate data from multiple sources to mitigate occlusion challenges. Zhang et al.[10] combined facial expressions, speech patterns (e.g., pitch, tone), and physiological signals such as heart rate to infer emotions. A hybrid CNN-RNN framework was used to process facial and speech features, demonstrating enhanced performance when compared to unimodal FER models. Multimodal learning compensates for missing information by providing supplementary emotional cues, particularly useful in masked scenarios. Advances in attention-based architectures enable robust FER under occlusions. Chen et al. utilized self-attention networks to highlight visible regions like the eyes and eyebrows, improving recognition performance even with significant occlusions. They explored transformers tailored for FER, showing that ViTs can learn long-range dependencies and integrate global and localized cues. Their multi-scale approach handled partial occlusions effectively by focusing on fine-grained and coarser features like head orientation. Data augmentation enhances model robustness by simulating occlusions during training. Guan et al.[11] introduced a dataset with artificial occlusions mimicking real-world scenarios like mask-wearing. They demonstrated that training on such augmented data significantly improved FER model robustness to occlusions. By shifting or modifying facial landmarks, models were trained to identify emotional cues despite occluded regions. Transfer learning addresses the lack of labeled masked FER datasets. Wang et al. leveraged pretrained FER models on unmasked datasets, adapting them to masked FER through fine-tuning. This method minimized the need for large masked datasets and preserved baseline accuracy.

3. PROPOSED METHOD

3.1 Mini-Xception Network Architecture

Mini-Xception Network is a lightweight, efficient variant of the Xception mode that retains the core benefits of depthwise separable convolutions and residual connections while reducing the model's size and complexity. The use of depthwise separable convolutions significantly reduces the number of parameters and computational costs while still maintaining high accuracy and performance, making it a suitable choice for real-time, mobile applications.

The Mini-Xception architecture is a lightweight convolutional neural network designed for efficient feature extraction and classification, primarily using depthwise separable convolutions to reduce computational complexity. It incorporates batch normalization to stabilize and accelerate training, ReLU activations for non-linearity, and global average pooling to minimize the number of parameters while maintaining spatial invariance. Downsampling is achieved through strided convolutions or max pooling, and a compact fully connected layer generates the final predictions. These components collectively enable the Mini-Xception to balance performance and efficiency, making it well-suited for tasks like FER on resource-constrained devices.

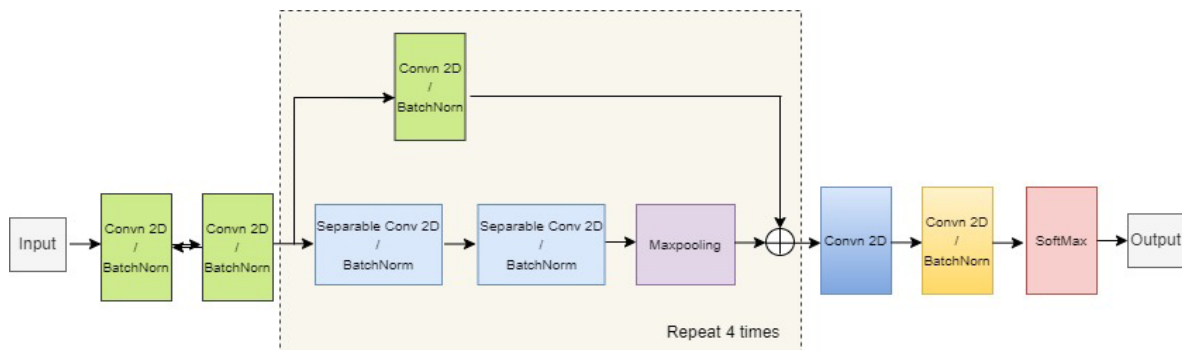


Figure 1. Mini-Xception Network Architecture

3.2 Convolutional Block Attention Module

CBAM is a lightweight and effective attention mechanism designed to improve feature representation in convolutional neural networks by sequentially applying channel and spatial attention. CBAM first employs a channel attention module, which aggregates spatial information using global average pooling (GAP) and global max pooling (GMP), followed by a shared multilayer perceptron (MLP) to compute channel-wise attention weights.

The specific steps are as follows:

$$F_{avg} = \text{AvgPool}(F) \in \mathbb{R}^{1 \times 1 \times C} \quad (1)$$

$$F_{max} = \text{MaxPool}(F) \in \mathbb{R}^{1 \times 1 \times C} \quad (2)$$

$$M_c(F) = \sigma(\text{MLP}(F_{avg}) + \text{MLP}(F_{max})) \in \mathbb{R}^{1 \times 1 \times C} \quad (3)$$

Where σ is the sigmoid activation function.

$$F' = M_c(F) \odot F \quad (4)$$

Where \odot represents element-by-element multiplication.

The spatial attention module generates spatial attention weights by averaging and max-pooling the channels of the feature map.

$$F'_{avg} = \text{AvgPool}(F') \in \mathbb{R}^{H \times W \times 1} \quad (5)$$

$$F'_{max} = \text{MaxPool}(F') \in \mathbb{R}^{H \times W \times 1} \quad (6)$$

$$M_s(F') = \sigma(\text{Conv}([F'_{avg}, F'_{max}])) \in \mathbb{R}^{H \times W \times 1} \quad (7)$$

$$F'' = M_s(F') \odot F' \quad (8)$$

3.3 Mini_Xception+CBAM Network Architecture

The proposed network architecture combines the Mini-Xception backbone with the CBAM module to address the challenges of masked FER. The overall design is illustrated in Figure 2, which demonstrates the integration of CBAM modules into the residual depthwise separable convolution blocks of Mini-Xception. It contains 4 residual depth-separable convolution blocks, in which an attention mechanism module (CBAM) is added to each residual depth-separable convolution block to enhance the focus on non-occlusion areas. After each convolution, the convolutional layer performs batch normalization operations and ReLu activation functions. The last layer uses a global average pooling layer and Softmax activation function for prediction.

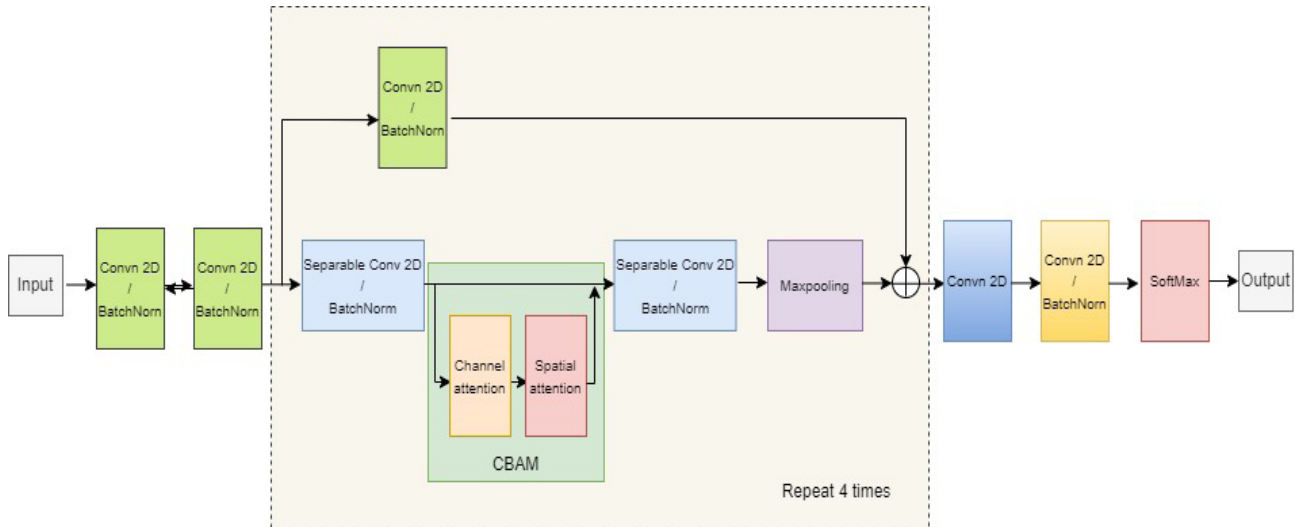


Figure 2 Mini_Xception+CBAM Network Architecture

To discuss different variant of model base on Mini_Xception

4.EXPERIMENTS

4.1 Datasets

To evaluate the performance of the proposed network, experiments were conducted on three benchmark datasets:

- FER2013 (Facial Expression Recognition 2013):

A benchmark dataset introduced during a Kaggle competition, commonly used in facial expression recognition tasks. It contains 35,887 grayscale images of size 48×48 pixels, categorized into seven emotion classes: Anger, Disgust, Fear, Happiness, Sadness, Surprise, and Neutral. The dataset is organized into three subsets: training, public test, and private test, and is stored in CSV format, with each image represented as a flat array of pixel intensities.

- SMFER2013 (Simulated Masked FER2013):

A synthetic dataset created by applying simulated facial masks to the original FER2013 images. This was achieved using the MaskTheFace tool, which overlays various types of face masks onto the images while preserving the structure and diversity of the original dataset. SMFER2013 retains the same CSV format and data organization as FER2013, making it compatible with existing workflows.

- MFER2013 (Masked FER2013):

A hybrid dataset combining FER2013 (unmasked images) and SMFER2013 (masked images). This dataset is designed to balance the class distributions of both masked and unmasked samples, ensuring sufficient representation of each emotion category. MFER2013 incorporates diverse mask types, facial orientations, and expression variations, making it more robust for training facial expression recognition models that need to generalize across masked and unmasked faces.

Evaluation Metrics

The performance of the proposed method was assessed using the following metrics:

- Accuracy (Acc): The ratio of correctly predicted samples to the total number of samples, reflecting overall model performance.
- Precision (P), Recall (R), and F1-Score: These metrics were calculated per class to assess the model's ability to distinguish individual emotions effectively, particularly under occlusion.
- Computational Efficiency: Metrics such as the number of parameters, FLOPs (floating-point operations), and inference time were analyzed to verify the suitability of the proposed model for resource-constrained environments.

4.2 Implementation Details

The experiments were implemented using Python 3.8 and PyTorch 1.12.0, with training conducted on an NVIDIA RTX 3090 GPU. The following configurations were applied:

4.3.1 Data Preprocessing.

All images were resized to 48×48 pixels and normalized. Random cropping and horizontal flipping were used for data augmentation. The parameters of the data enhancement operation are set to randomly rotate the image, with a rotation angle range of -20 degrees to 20 degrees; randomly translate the image horizontally and vertically, with the translation distance within 10% of the image width and height; randomly scale the image, with the scaling ratio within the range of 0.1. When the image is translated, scaled, flipped, etc., black areas may be generated at the edge of the image, and these black areas are filled with the nearest pixels. When adjusting the training parameters, the dataset is randomly divided so that the validation set and the test set each account for 10% of the total dataset, and the remaining 90% is used as the training set.

4.3.2 Training Parameters.

The model was trained for 100 epochs using the Adam optimizer with an initial learning rate of 10^{-3} , decayed by a factor of 0.1 every 30 epochs. The batch size was set to 64.

4.3.3 Loss Function.

A weighted cross-entropy loss function was employed to handle class imbalances.

4.4 Experimental Results

4.4.1 Quantitative Analysis.

Table 1 compares the performance of various models on the MFER2013 dataset, highlighting the accuracy, precision, recall, and F1-score of each model. Among the models, Mini-Xception+CBAM achieves the highest accuracy of 91.5%, followed by Vision Transformer (ViT) at 89.1%, and ResNet-18 at 87.2%. MobileNetV2 performs the weakest with an accuracy of 84.7%. In terms of computational efficiency, Mini-Xception+CBAM remains lightweight, with 0.68M parameters and 30.1M FLOPs, offering a balance between high performance and low resource usage.

Table 1. Comparison of Performance of different models on MFER2013 dataset

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Parameters (M)	FLOPs (M)
Mini-Xception	86.3	85.5	84.9	85.2	0.61	28.3
Mini-Xception+CBAM	91.5	91.0	91.3	91.2	0.68	30.1
MobileNetV2	84.7	83.5	83.2	83.3	2.23	56.8
ResNet-18	87.2	86.4	86.1	86.3	11.69	98.7
Vision Transformer (ViT)	89.1	88.5	88.9	88.7	13.45	102.3

Table 2 presents the performance of the proposed method across three datasets: FER2013, SMFER2013, and MFER2013. The model achieves the highest accuracy on MFER2013 with 91.5%, followed by FER2013 with 90.4%, and SMFER2013 with 89.2%. The performance on the masked dataset (SMFER2013) is slightly lower than on FER2013, indicating the challenge of recognizing expressions under occlusion. Despite this, the model maintains high precision, recall, and F1-score across all datasets, with consistent computational efficiency (0.68M parameters and 30.1M FLOPs), demonstrating its robustness and suitability for practical applications in facial expression recognition.

Table 2. Comparison of Performance of The proposed Method on FER2013, SMFER2013, MFER 2013 dataset

Dataset	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Parameters (M)	FLOPs (M)
FER2013	90.4	89.8	90.1	90.0	0.68	30.1
SMFER2013	89.2	88.6	89.0	88.8	0.68	30.1
MFER2013	91.5	91.0	91.3	91.2	0.68	30.1

4.4.2 Ablation Study.

An ablation study was conducted to quantify the impact of each component of the proposed method. Results showed that adding CBAM to Mini-Xception improved accuracy by 4.1%, demonstrating its effectiveness in handling occluded features.

Table 3. An ablation study of each component of the proposed method

Model Variant	Accuracy (%)
Mini-Xception	86.3
Mini-Xception + Channel Attention	88.0
Mini-Xception + Spatial Attention	87.5
Mini-Xception + CBAM	90.4

4.4.3 Efficiency Comparison.

The proposed model demonstrated superior efficiency, with only 0.68M parameters and 30.1M FLOPs, making it well-suited for real-time applications compared to ResNet-18 and ViT.

5.CONCLUSION

In this study, we proposed a novel lightweight architecture, Mini-Xception + CBAM, for robust facial expression recognition (FER), particularly in scenarios involving masked or partially occluded faces. By integrating the Convolutional Block Attention Module (CBAM) into the efficient Mini-Xception backbone, the model enhances its focus on critical non-occluded regions, such as the eyes and eyebrows, which are crucial for expression recognition under occlusions.

Extensive experiments on the FER2013, SMFER2013, and MFER2013 datasets demonstrated that the proposed architecture achieves superior accuracy and efficiency compared to several baseline models, including Mini-Xception, MobileNetV2, and ResNet-18. On MFER2013, the proposed model achieved an impressive 91.5% accuracy, significantly outperforming other models, even in the presence of masked faces. The incorporation of CBAM led to substantial performance improvements, as evidenced by both quantitative metrics and qualitative attention map visualizations. Furthermore, the lightweight design of Mini-Xception + CBAM, with only 0.68M parameters and 30.1M FLOPs, makes it highly suitable for deployment in resource-constrained environments, such as mobile and edge devices.

REFERENCES

- [1] Wang, C., et al..2021. COVID-19 in early 2021: current status and looking forward. *Signal Transduction and Targeted Therapy*, 2021. 6(1): p. 1–14.
- [2] Kotsia, I., I. Buciu, and I. Pitas.2008. An analysis of facial expression recognition under partial facial image occlusion. *Image and Vision Computing*, 2008. 26(7): p. 1052–1067.
- [3] Mollahosseini, A., B. Hasani, and M.H. Mahoor.2017. Affectnet: A database for facial expression, valence, and arousal computing in the wild. *IEEE Transactions on Affective Computing*, 2017. 10(1): p. 18–31.
- [4] Goodfellow, I.J., et al..2013. Multi-digit number recognition from street view imagery using deep convolutional neural networks. *arXiv preprint arXiv:1312.6082*, 2013.
- [5] Chollet, F.2017. Xception: Deep learning with depthwise separable convolutions. in *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.
- [6] Arriaga, O., M. Valdenegro-Toro, and P. Ploger.2017. Real-time convolutional neural networks for emotion and gender classification. *arXiv preprint arXiv:1710.07557*, 2017.
- [7] Chen, J., et al..2023 Paying attention in metaverse: an experiment on spatial attention allocation in extended reality shopping. *Information Technology & People*, 2023. 36(8): p. 255–283.
- [8] Zhang, W., et al..2021 Reservoir inflow predicting model based on machine learning algorithm via multi - model fusion: A case study of Jinshuitan river basin. *IET Cyber - Systems and Robotics*, 2021. 3(3): p. 265–277.
- [9] Chen, F., et al..2024. Review of lightweight deep convolutional neural networks. *Archives of Computational Methods in Engineering*, 2024. 31(4): p. 1915–1937.
- [10] Wang, J., et al..2022 Modeling mask uncertainty in hyperspectral image reconstruction. in *European Conference on Computer Vision*. 2022. Springer.
- [11] Guan, Q., et al.,20 Medical image augmentation for lesion detection using a texture-constrained multichannel progressive GAN. *Computers in Biology and Medicine*, 2022. 145: p. 105444.