# Gender Recognizer based on Human Face using CNN and Bottleneck Transformer Encoder

Adri Priadana
*Department of Electrical, Electronic and Computer Engineering*
*University of Ulsan*
Ulsan, Korea
priadana3202@mail.ulsan.ac.kr

Muhamad Dwisnanto Putro
*Department of Electrical Engineering*
*Universitas Sam Ratulangi*
Manado, Indonesia
dwisnantoputro@unsrat.ac.id

Jinsu An
*Department of Electrical, Electronic and Computer Engineering*
*University of Ulsan*
Ulsan, Korea
jinsu5023@islab.ulsan.ac.kr

Duy-Linh Nguyen
*Department of Electrical, Electronic and Computer Engineering*
*University of Ulsan*
Ulsan, Korea
ndlinh301@mail.ulsan.ac.kr

Xuan-Thuy Vo
*Department of Electrical, Electronic and Computer Engineering*
*University of Ulsan*
Ulsan, Korea
xthuy@islab.ulsan.ac.kr

Kang-Hyun Jo
*Department of Electrical, Electronic and Computer Engineering*
*University of Ulsan*
Ulsan, Korea
acejo@ulsan.ac.kr

*Abstract*—Several applications, such as Human-Robot interactions and offline advertising platforms, perform gender recognition based on a human face to profile their audience. These applications demand gender recognition that can operate in real-time on a low-cost or CPU device. This work proposes a gender recognizer based on a human face using a Convolutional Neural Network (CNN) and Bottleneck Transformer Encoder (BTE) that renders low parameters and operation. BTE is offered to support the primary CNN feature extractor in learning the global representation of the feature maps efficiently. This work uses three face gender datasets benchmark, namely UTKFace, Labeled Faces in the Wild (LFW), and Adience, to train and validate the proposed network. The CNN network consisted of the BTE achieves competitive accuracy compared to the state-of-the-art network. The recognizer can operate in real-time on a CPU with 150 frames per second.

*Index Terms*—Bottleneck Transformer, Convolutional Neural Network (CNN), Efficient Architecture, Gender Recognition, Self-Attention

## I. INTRODUCTION

As a unique soft biometric property, gender is an essential part that can be used in profiling users. Several applications utilize gender by recognizing it based on a human face to support their purpose [1]. In Human-Robot interactions [2], [3], robots can use gender property to characterize their interlocutors to personalize and make conversations more pleasant. Foggia et al. [2] demonstrated that people prefer conversations with robots that can personalize their interlocutors according to gender by using appropriate pronouns and honorifics. In advertising applications such as digital signage, gender recognition is used as a basis to provide relevant advertising content in real-time [4]. This technology can analyze the gender of individuals facing the platform and adapting the displayed advertisements accordingly. By tailoring relevant advertising messages to the audience, businesses can enhance the effectiveness of their marketing efforts and provide a more personalized experience to potential customers.

Gender recognition through a human face is initiated by detecting the audience's face via a camera and then classifying it to the gender class. In robotics or digital signage applications utilizing a low-cost device, there is a need for the gender recognition operation to function effectively in real time when individuals are facing it. Therefore, an efficient gender recognition mechanism with low computation is required. Hence, robotics or digital signage platform that uses a low-cost device, or at least a CPU device, can operate it properly.

Convolutional Neural Networks (CNNs) have recently demonstrated substantial effectiveness in recognition tasks [5]. To construct a recognition system, particularly for gender recognition based on human faces, many researchers have created various CNN designs. Savchenko [6] proposed a CNN architecture modified from MobileNet to predict facial gender, gaining 91.95% based on accuracy on the UTKFace [7] dataset. HyperFace-ResNet [8], improved ResNet-101 [9] architecture by fusing the lower and deeper layers using an element-wise addition operation, is offered to predict facial gender. This presented architecture gained 94% based on accuracy on the LFW [10] datasets. Greco et al. [4] applied MobileNetV2 architecture to develop a gender recognizer system that supports a digital signage platform. Other works [11]–[13] designed more efficient CNN architectures for gender recognition based on human face generating fewer parameters and achieved great accuracy on Adience [14], UTKFace [7], and LFW [10] datasets. These architectures can perform fast in a CPU device. Employing CNN architecture with fewer parameters will improve efficiency and lead the recognizer to operate more quickly, especially on low-cost or CPU devices.

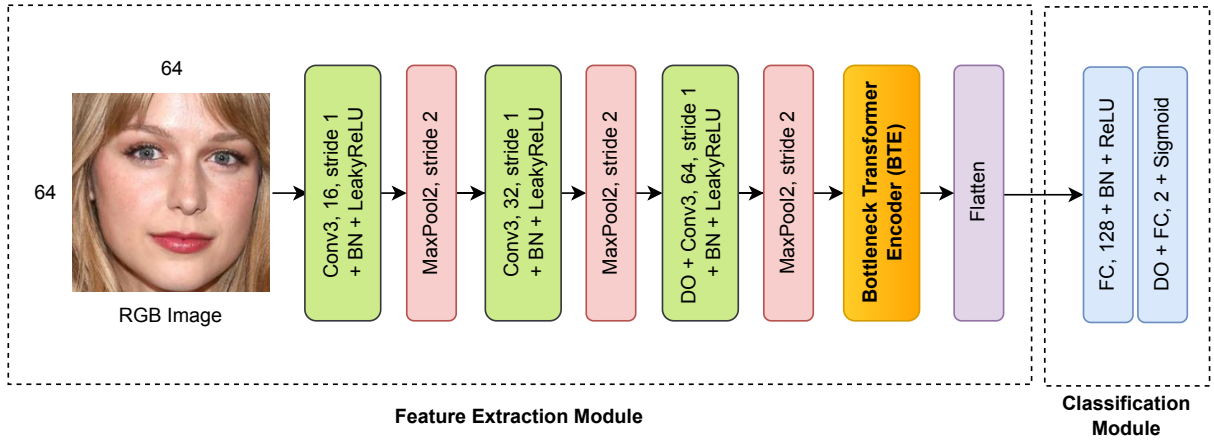Considering the efficiency and computation, a lightweight

Fig. 1. Overall view of the proposed CNN network with Bottleneck Transformer Encoder (BTE) for facial gender recognition based on human faces.

CNN network is proposed incorporated with Bottleneck Transformer Encoder (BTE) to perform gender recognition based on a human face. This architecture generates a few parameters and low operation, appropriate to implement on a CPU device. In this work, the contribution is outlined as follows:

1) A gender recognition based on a human face with only 555,770 parameters and 24.8 MFLOPs achieves very competitive accuracy and efficiency compared to other networks on three datasets, UTKFace [7], Adience [14], and LFW [10].
2) A Bottleneck Transformer Encoder (BTE) leads the architecture not only able to capture the local representation but also the global representation by efficiently calculating the relationship between each spatial information.
3) The proposed gender recognizer based on a human face can operate in real-time on a CPU with a speed of 131.85 frames per second, outperforming the other gender recognizer based on a human face.

## II. PROPOSED ARCHITECTURE

The proposed gender recognition architecture has a feature extraction and classification modules shown in Fig. 1. It produces only 555,770 parameters.

### A. Feature Extraction Module

The feature extraction module serves as an instrument to grasp facial features from the input face area image, containing three of $3 \times 3$ convolution layers with a growing number of channels from 16 to 64. This mechanism is designed to capture additional details from the higher-level facial image features. This module utilizes Batch Normalization (BN) as well as Leaky Rectified Linear Unit (LeakyReLU) activation in each convolutional layer to address gradient-related issues [15]. Before the final convolutional layer, a dropout (DO) operation [16] is applied to mitigate the risk of overfitting. After each convolutional layer, a downsampling mechanism is employed, consisting of three times $2 \times 2$ max-pooling

operations with a stride of two. Involving several convolution layers captures only the local representation of the facial feature. Therefore, we propose a Bottleneck Transformer Encoder (BTE) to grasp the global representation of the feature maps. In this architecture, BTE is positioned following the final max-pooling operation and just before the flattening operation.

### B. Bottleneck Transformer Encoder (BTE)

Modern networks concentrate on improving the self-attention mechanism oriented to the success of the Vision Transformer (ViT) [17] technique. However, it brings high operations (multiplication and addition), leading to inadequate implementation for vision tasks on low-cost devices. This work proposes a Bottleneck Transformer Encoder (BTE) to alleviate the mentioned above problem. Similar to the general Transformer Encoder, this function will transform a tensor input $\mathbf{X}$ of shape $H \times W \times C$ into a query ($\mathbf{Q}$), key ($\mathbf{K}$), and value ($\mathbf{V}$) by applying a $1 \times 1$ convolution operation with Rectified Linear Unit (ReLU). However, a reduction channel $r$ and multi-head mechanism with $NumHead$ are applied in this work to produce a slimmer tensor $H \times W \times ((\frac{C}{NumHead})/r)$, followed by reshaping operation to produce a $\mathbf{Q}$, $\mathbf{K}$, and $\mathbf{V}$ matrix with shape $HW \times ((\frac{C}{NumHead})/r)$. After that, we compute these matrices using scaled dot-product attention shown in Fig. 2, which is described as follows:

$$SDPA\left(\mathbf{Q}, \mathbf{K}, \mathbf{V}\right) = softmax\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}, \qquad (1)$$

where $d_k$ is a scaling factor to control the softmax temperature and $T$ is a transpose matrix operation. To reform the output matrix shape into $H \times W \times ((\frac{C}{NumHead})/r)$, the reshaping operation is also performed. Then, a concatenation operation is used to merge all output from all heads, followed by $1 \times 1$ convolution operation with a dropout (DO) and Rectified Linear Unit (ReLU) activation layers to restore the number of channels according to the input tensor $\mathbf{X}$ with shape $H \times W \times C$ that represents a bottleneck mechanism. Once again, a linear projection using $1 \times 1$ convolution operation
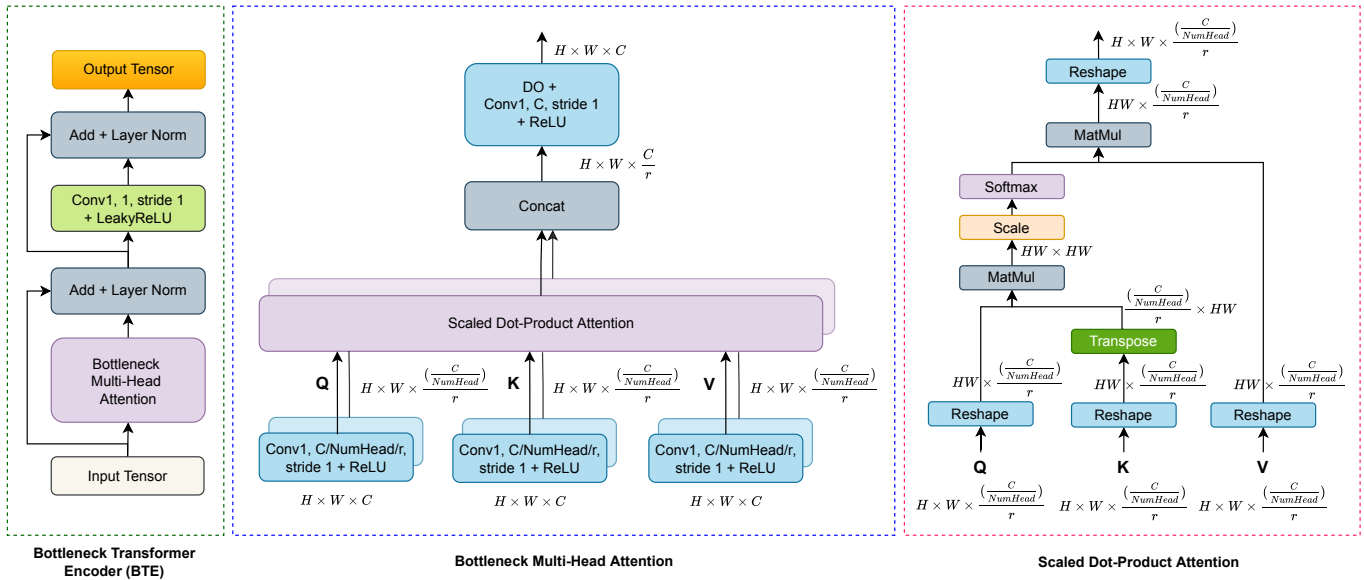
Fig. 2. The proposed Bottleneck Transformer Encoder (BTE).

with a LeakyReLU activation is performed, integrating with a residual connection [9] and Layer Normalization (LN) [18] around the multi-head attention and last convolution layer.

### C. Classification Module

The second phase of the gender recognition architecture proposed here involves a classification module comprising two fully connected (FC) layers, one with 128 units and the other with 2 units. The first FC layer incorporates BN and ReLU (Rectified Linear Unit) activation to mitigate gradient-related challenges. The second FC layer employs a dropout mechanism [16] to prevent overfitting and utilizes Sigmoid activation to produce the input for the final prediction decision. This classification module is utilized to compute the probability associated with each gender class, aiding in the discrimination between male and female based on facial features.

### III. IMPLEMENTATION CONFIGURATION

In this experiment, UTKFace, LFW, and Adience datasets are used to train the proposed architecture. The training process utilizes an initial learning rate of $10^{-2}$, 512 batch size, and 300 epochs. Tensorflow and Keras are used as a framework, while an Nvidia GTX 1080Ti with 11GB of GPU is employed as the accelerator. The Binary Cross-Entropy loss and Adam optimizer are used for the training. To further enhance training, a learning rate reduction strategy is implemented, reducing the rate by 75% when accuracy fails to improve over 20 consecutive epochs. The computational speed of the proposed architecture is assessed on an Intel Core i7-9750H CPU clocked at 2.6 GHz, with 20GB of RAM.

TABLE I
ASSESSMENT OUTCOMES ON THREE DATASETS USED IN THIS WORK.

| Networks | Number of Parameters | Accuracy on Validation (%) |
|---|---|---|
| **UTKFace** | | |
| EfficientNetV2B1 [19] | 77,54,679 | 90.31 |
| VGG16 [20] | 138,357,544 | 91.90 |
| Shahzeb et al. [21] | 14,715,201 | 91.94 |
| Savchenko (Modified MobileNet) [6] | 3,491,521 | 91.95 |
| SufiaNet [12] | 226,574 | 92.05 |
| MPConvNet [11] | 659,650 | 92.32 |
| MudaNet [13] | 674,760 | 92.66 |
| **Proposed** | **555,770** | **92.78** |
| **LFW** | | |
| Rouhsedaghat et al. [22] | 16,900 | 94.63 |
| SufiaNet [12] | 226,574 | 95.66 |
| MudaNet [13] | 674,760 | 96.22 |
| MPConvNet [11] | 659,650 | 96.30 |
| Greco et al. [23] | 3,538,984 | **98.73** |
| **Proposed** | **555,770** | **96.50** |
| **Adience** | | |
| SufiaNet [12] | 226,574 | 84.60 |
| MudaNet [13] | 674,760 | 84.85 |
| Saha et al. [24] | 3,190,913 | 84.94 |
| MPConvNet [11] | 659,650 | 85.67 |
| Opu et al. [25] | 210,050 | 85.77 |
| **Proposed** | **555,770** | **85.86** |

### IV. EXPERIMENTAL RESULTS

#### A. Assessment on Datasets for Face Gender Classification

*1) UTKFace:* There are more than 23,000 facial images on this dataset labeled in gender, age, and ethnicity, covering many variations such as expression, pose, age, illumination, etc. The proposed architecture is trained on 70% of this dataset, while the remains for testing sets. The proposed architecture attains a validation accuracy of 92.78% on this

TABLE II
MODEL ANALYSIS OF THE PROPOSED MODEL BASED ON UTKFACE.

| Settings | Number of Parameters | MFLOPs | Validation Accuracy (%) |
|---|---|---|---|
| The Backbone without BTE | 549,218 | 23.92 | 92.46 |
| The Backbone with BTE | 555,770 | 24.80 | 92.78 |

dataset while utilizing only 555,770 parameters. Table I shows that the accuracy surpasses the other State-of-The-Art (SOTA) architectures, which differed by 0.12% from the second-best. The proposed architecture also generates fewer parameters compared to the second-best.

*2) LFW:* This dataset includes more than 13,000 face images, divided into testing (30%) and training (70%). This dataset has two labels, females and males, with a substantial disparity, with females accounting for 23% and males constituting 77%. Table I describes that the proposed architecture achieves competitive validation accuracy with 96,50% on this dataset. It becomes second-best compared to the other SOTA architectures. However, the proposed architecture generates far fewer parameters.

*3) Adience:* This dataset comprises over 26,000 facial images annotated with gender and age labels, encompassing a range of variations including age, pose, noise, lighting, appearance, and more. Some pre-processing is established on this dataset in this experiment, such as eradicating data with missing values and generating 17,492 face images. Following the two other datasets mentioned before, this data set is also split into 30% and 70% as a testing and training set, respectively. As a result, the proposed architecture gains a validation accuracy of 85.86%, shown in Table I. This result surpasses the other SOTA architectures, which differed by 0.09% from the second-best.

*B. Ablation Study*

*1) Model Analysis:* This analysis is established on the UTKFace dataset to investigate the impact of the proposed BTE on recognition accuracy in this work. Table II reveals that the proposed BTE can improve the validation accuracy by 0.32%. With this increased performance, the proposed BTE only generates 6,552 and 0.88 more additional parameters and MFLOPs, respectively.

*2) Channel Reduction Analysis:* This work arranges this analysis to examine the optimal channel reduction value on the proposed BTE concerning the recognition performance. Table III shows that using a channel reduction value of one does not enhance performance. The proposed BTE with a channel reduction value of eight produces the highest validation accuracy in this work.

*C. Runtime Efficiency*

The proposed architecture, with only 555,770 parameters and 24.8 MFLOPs, is principally developed to perform on a CPU device in real-time scenarios. Due to this, the proposed architecture can operate 345.69 FPS for gender recognition

TABLE III
CHANNEL REDUCTION ANALYSIS OF THE PROPOSED BTE
ON UTKFACE DATASET

| Value of Reduction $r$ | Number of Parameters | MFLOPs | Validation Accuracy (%) |
|---|---|---|---|
| - | 570,274 | 27.18 | 92.49 |
| 2 | 561,986 | 25.72 | 92.53 |
| 4 | 557,842 | 25.09 | 92.57 |
| 6 | 556,288 | 24.87 | 92.63 |
| **8** | **555,770** | **24.80** | **92.78** |
| 10 | 555,252 | 24.73 | 92.60 |

- indicates that the BTE is performed without channel reduction $r$

TABLE IV
EFFICIENCY IN TERMS OF RUNTIME ON A CPU SETUP

| Networks | Number of Parameters | MFLOPs | GR (FPS) | FD + GR (FPS) |
|---|---|---|---|---|
| VGG16 with BN | 39,782,722 | 2,290 | 42.42 | 36.35 |
| ResNet50V2 | 23,568,898 | 571 | 57.15 | 46.67 |
| MudaNet [13] | 674,760 | 69 | 58.57 | 48.07 |
| MobileNetV2 | 2,260,546 | 50 | 118.56 | 80.43 |
| SufiaNet [12] | 226,574 | 22 | 127.50 | 84.90 |
| MPConvNet [11] | 659,650 | 67 | 265.07 | 131.85 |
| **Proposed** | **555,770** | **25** | **345.69** | **150.49** |

GR indicates the Gender Recognition
FD + GR indicates the Gender Recognition integrated with Face Detection

and 150.49 FPS when integrated with an efficient face detector, dubbed LWFCPU [26], as shown in Table IV. The proposed gender recognizer based on a human face outperforms the other public and light architecture used in this case. Fig. 3 demonstrates the gender recognition result of the proposed gender recognizer based on a human face on a CPU. The male and female faces are indicated with fuchsia and yellow bounding box, respectively.

*D. Limitations*

The UTKFace dataset used to train the proposed architecture for the gender recognizer in this work encompasses a range of pose variations. Nevertheless, it lacks face instances with extreme yaw and pitch pose, pushing to inaccurate predictions in certain cases involving faces in extreme yaw and pitch poses, as shown in the last row of Fig. 3. In this case, a female face is recognized as a male face.

V. CONCLUSION

This work proposes a gender recognizer based on a human face. It offers a light CNN as a feature extractor supported by a Bottleneck Transformer Encoder (BTE) to learn the global representation of the feature maps efficiently. This encoder will enhance the feature map quality generated by the feature extractor. The proposed network achieves rivalrous accuracy compared to the SOTA network on three face gender dataset benchmarks (UTKFace, LFW, and Adience). The proposed gender recognizer based on a human face consisting of the proposed gender recognition integrated with a face detector can operate in real-time on the CPU with 150.49 frames per

Fig. 3. The accurate recognition outcome (the three first rows) and the inaccurate recognition outcomes (the last row) generated by the proposed gender recognizer based on human faces.

second. In future research, we will investigate alternative approaches to address the limitations identified in this study and to enhance both the accuracy and speed of gender recognition through a human face.

## REFERENCES

[1] P. Foggia, A. Greco, A. Saggese, and M. Vento, "Multi-task learning on the edge for effective gender, age, ethnicity and emotion recognition," *Engineering Applications of Artificial Intelligence*, vol. 118, p. 105651, 2023.

[2] P. Foggia, A. Greco, G. Percannella, M. Vento, and V. Vigilante, "A system for gender recognition on mobile robots," in *Proceedings of the 2nd international conference on applications of intelligent systems*, 2019, pp. 1–6.

[3] M. A. Uddin, M. S. Hossain, R. K. Pathan, and M. Biswas, "Gender recognition from human voice using multi-layer architecture," in *2020 International Conference on INnovations in Intelligent SysTems and Applications (INISTA)*, 2020, pp. 1–7.

[4] A. Greco, A. Saggese, and M. Vento, "Digital signage by real-time gender recognition from face images," in *2020 IEEE International Workshop on Metrology for Industry 4.0 & IoT*. IEEE, 2020, pp. 309–313.

[5] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 11 966–11 976.

[6] A. V. Savchenko, "Efficient facial representations for age, gender and identity recognition in organizing photo albums using multi-output convnet," *PeerJ Computer Science*, vol. 5, p. e197, 2019.

[7] Z. Zhang, Y. Song, and H. Qi, "Age progression/regression by condi-

tional adversarial autoencoder," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5810–5818.

[8] R. Ranjan, V. M. Patel, and R. Chellappa, "Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 1, pp. 121–135, 2017.

[9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016, pp. 770–778.

[10] G. B. Huang, M. Mattar, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database forstudying face recognition in unconstrained environments," in *Workshop on faces in'Real-Life'Images: detection, alignment, and recognition*, 2008.

[11] A. Priadana, M. D. Putro, and K.-H. Jo, "An efficient face gender detector on a cpu with multi-perspective convolution," in *2022 13th Asian Control Conference (ASCC)*. IEEE, 2022, pp. 453–458.

[12] A. Priadana, M. D. Putro, C. Jeong, and K.-H. Jo, "A fast real-time face gender detector on cpu using superficial network with attention modules," in *2022 International Workshop on Intelligent Systems (IWIS)*. IEEE, 2022, pp. 1–6.

[13] A. Priadana, M. D. Putro, X.-T. Vo, and K.-H. Jo, "A facial gender detector on cpu using multi-dilated convolution with attention modules," in *2022 22nd International Conference on Control, Automation and Systems (ICCAS)*. IEEE, 2022, pp. 190–195.

[14] E. Eidinger, R. Enbar, and T. Hassner, "Age and gender estimation of unfiltered faces," *IEEE Transactions on information forensics and security*, vol. 9, no. 12, pp. 2170–2179, 2014.

[15] A. Priadana, M. D. Putro, D.-L. Nguyen, X.-T. Vo, and K.-H. Jo, "Human face detector with gender identification by split-based inception block and regulated attention module," in *International Workshop on Frontiers of Computer Vision*. Springer, 2023, pp. 163–177.

[16] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.

[17] 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.

[18] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.

[19] P. Vidyarthi, S. Dhavale, and S. Kumar, "Gender and age estimation using transfer learning with multi-tasking approach," in *2022 2nd Asian Conference on Innovation in Technology (ASIANCON)*. IEEE, 2022, pp. 1–5.

[20] A. Krishnan, A. Almadan, and A. Rattani, "Understanding fairness of gender classification algorithms across gender-race groups," in *2020 19th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 2020, pp. 1028–1035.

[21] M. Shahzeb, S. Dhavale, D. Srikanth, and S. Kumar, "Dcnn-based transfer learning approaches for gender recognition," in *International Conference on Data Management, Analytics & Innovation*. Springer, 2023, pp. 357–365.

[22] M. Rouhsedaghat, Y. Wang, X. Ge, S. Hu, S. You, and C.-C. J. Kuo, "Facehop: A light-weight low-resolution face gender classification method," in *International Conference on Pattern Recognition*. Springer, 2021, pp. 169–183.

[23] A. Greco, A. Saggese, M. Vento, and V. Vigilante, "A convolutional neural network for gender recognition optimizing the accuracy/speed tradeoff," *IEEE Access*, vol. 8, pp. 130771–130781, 2020.

[24] A. Saha, S. N. Kumar, and P. Nithyakani, "Age and gender prediction using adaptive gamma correction and convolutional neural network," in *2023 International Conference on Computer Communication and Informatics (ICCCI)*. IEEE, 2023, pp. 1–5.

[25] M. N. I. Opu, T. K. Koly, A. Das, and A. Dey, "A lightweight deep convolutional neural network model for real-time age and gender prediction," in *2020 Third International Conference on Advances in Electronics, Computers and Communications (ICAECC)*. IEEE, 2020, pp. 1–6.

[26] M. D. Putro, D.-L. Nguyen, and K.-H. Jo, "Lightweight convolutional neural network for real-time face detector on cpu supporting interaction of service robot," in *2020 13th International Conference on Human System Interaction (HSI)*. IEEE, 2020, pp. 94–99.