

MAiVAR-T: Multimodal Audio-image and Video Action Recognizer using Transformers

Muhammad Bilal Shaikh*, Douglas Chai†
School of Engineering
Edith Cowan University
*mbshaikh@our.ecu.edu.au, †d.chai@ecu.edu.au

Syed Mohammed Shamsul Islam
School of Science
Edith Cowan University
syed.islam@ecu.edu.au

Naveed Akhtar

Department of Computer Science & Software Engineering
The University of Western Australia
naveed.akhtar@uwa.edu.au

Abstract—In line with the human capacity to perceive the world by simultaneously processing and integrating high-dimensional inputs from multiple modalities like vision and audio, we propose a novel model, MAiVAR-T (Multimodal Audio-Image to Video Action Recognition Transformer). This model employs an intuitive approach for the combination of audio-image and video modalities, with a primary aim to escalate the effectiveness of multimodal human action recognition (MHAR). At the core of MAiVAR-T lies the significance of distilling substantial representations from the audio modality and transmuting these into the image domain. Subsequently, this audio-image depiction is fused with the video modality to formulate a unified representation. This concerted approach strives to exploit the contextual richness inherent in both audio and video modalities, thereby promoting action recognition. In contrast to existing state-of-the-art strategies that focus solely on audio or video modalities, MAiVAR-T demonstrates superior performance. Our extensive empirical evaluations conducted on a benchmark action recognition dataset corroborate the model's remarkable performance. This underscores the potential enhancements derived from integrating audio and video modalities for action recognition purposes.

Index Terms—Multimodal Fusion, Transformers, Human Action Recognition, Deep Learning.

I. INTRODUCTION

Human action recognition has become a critical task in various fields such as surveillance [1], robotics [2], interactive gaming [3], and health care [4]. Traditionally, most approaches have focused on visual cues [5]. However, human actions are not limited to visual manifestations; they also consist of rich auditory information [6]. Accordingly, Multimodal human action recognition (MHAR) that incorporates both visual and audio cues can provide more comprehensive and accurate recognition results [7].

Despite these promising prospects, the performance of current MHAR models is hampered by challenges of multimodal data fusion. Existing methods, including Convolutional Neural Networks (CNNs) [8]–[10] require significantly more computation than their image counterparts, some architectures factorise convolutions across spatiotemporal dimensions. Contrastingly, Recurrent Neural Networks (RNNs) and Long

Short-Term Memory (LSTMs) [11] have demonstrated constraints in processing large sequences, memory efficiency and parallelism.

In this paper, we propose a novel transformer-based model, Multimodal Audio-image and Video Action Recognizer using Transformers (MAiVAR-T). Our approach capitalizes on the self-attention mechanism inherent in transformers [12] to extract relevant features from both modalities and fuse them effectively. The proposed MAiVAR-T model outperforms state-of-the-art MHAR models on benchmark datasets [13], demonstrating the potential of transformer-based architectures in improving multimodal fusion and recognition accuracy.

To summarize, the contributions made in this paper are:

- A new feature representation strategy is proposed to select the most informative candidate representations for audio-visual fusion;
- Collection of effective audio-image-based representations that complement video modality for better action recognition are included;
- We apply a novel MAiVAR-T framework (see Fig. 1) for audio-visual fusion that supports different audio-image representations and can be applied to different tasks; and
- State-of-the-art results for action recognition on the audio-visual dataset have been reported.

The remainder of the paper is organized as follows: we begin with a review of related works on MHAR (Section II), followed by a detailed discussion of the proposed methodology (Section III). We then present the experimental setup (Section IV) and report the results (Section V). Finally, we conclude the paper with future directions (Section VI).

II. RELATED WORK

A. Deep Learning for MHAR

Recently, deep learning models have shown remarkable results in MHAR [14]. They are capable of automatically learning a hierarchy of intricate features from raw multimodal data, which are beneficial for action recognition tasks.

CNNs have been widely adopted for MHAR to automatically extract spatial features from input data [15], and LSTMs

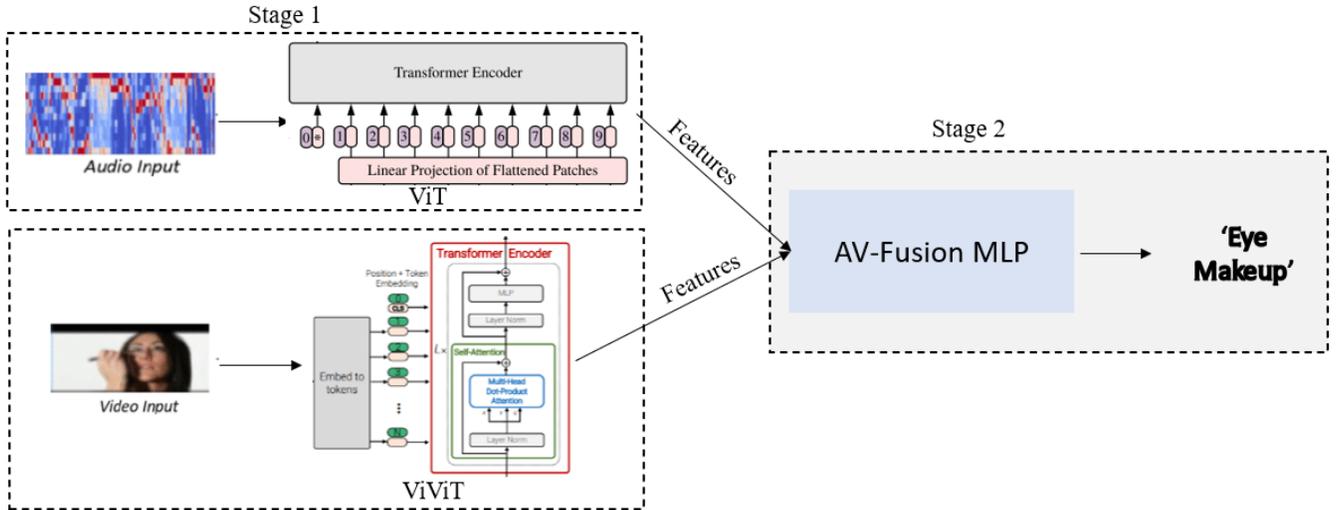


Fig. 1: The proposed framework contains two stages. The first stage extracts the features influencing the recognition while the second stage performs classification on the fused features. The input sequence consists of image and audio-image patches. These are then projected into tokens and appended to special CLS (classification). Our transformer encoder then uses self attention to model unimodal information, and send cross-modal information flow through to fusion network.

are typically used for modelling the temporal dynamics of actions [11]. However, the traditional combination of CNNs and LSTMs for MHAR faces challenges such as ineffective multimodal fusion and difficulty handling long temporal sequences.

Transformers, introduced by Vaswani et al. [12], have demonstrated their superiority in many fields like natural language processing [16], image classification [17], and video understanding [18]. The self-attention mechanism through its optimal complexity (see Table I) in transformers could potentially enhance the capability of feature extraction and multimodal fusion in MHAR tasks. However, the utilization of transformers in MHAR is relatively unexplored and demands further investigation.

B. Audiovisual Learning and Fusion

The field of audiovisual multimodal learning has a long and diverse history, both preceding and during the deep learning era [19]. Early research focused on simpler approaches, utilizing hand-designed features and late-stage processing, due to limitations in available data and computational resources [20]. However, with the advent of deep learning, more sophisticated strategies have emerged, enabling the implicit learning of modality-specific or joint latents to facilitate fusion. As a result, significant advancements have been achieved in various supervised audiovisual tasks [21].

It is common to jointly train multiple modality-specific convolution networks, where the intermediate activations are combined either through summation [22]. On the other hand, in transformer-based architectures, the incorporation of Vision Transformers (ViT) [17] and Video Vision Transformers (ViViT) [18] has brought about significant advancements in multimodal human action recognition. Initially, ViT proved

TABLE I: Complexity comparison for different types of layer. Notations: n : sequence length, d : representation dimension, k kernel size.

Layer Type	Complexity per layer	Sequential Operations	Maximum Path Length
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$

TABLE II: Hyper-parameters of the network.

Parameter	Value
Batch size	256
Initial learning rate	0.001
lr decay (every 4 epochs)	0.10
Learning rate patience	10
Epochs	100

instrumental in dissecting images into smaller segments, to interpret these patches as a sequence for more accurate image understanding. This ability greatly improved the recognition and classification of human actions within still images. The introduction of ViViT further extended this capacity, applying transformer techniques to analyze video data. By processing sequences of video frames, ViViT effectively interprets the spatio-temporal dynamics involved in human movements. Together, the use of Vision Transformers and Video Vision Transformers can produce a shift in multimodal human action recognition, enhancing the capability of systems to accurately classify and understand complex human activities across visual and audio domains.

III. PROPOSED METHODOLOGY

Data Collection: We collected human actions from a benchmark dataset called UCF101 [13], with each instance containing video clips and their corresponding audio streams. UCF-101 contains an average length of 180 frames per video. We observed that half of the videos in the dataset contained no audio. Thus, in order to focus on the effect of audio features, we used only those videos that contained audio. This resulted in 6837 videos across 51 categories. Whilst this led the dataset to be significantly reduced, the distribution of the audio dataset was similar to the video dataset. We used the first train-test split setting provided with this dataset, which resulted in 4893 training and 1944 testing samples. We reported the top 1 accuracies obtained by training on split 1.

Data Preprocessing: The video and audio data were preprocessed separately, as described in the following subsections. The video data was transformed into frames, while the audio data was converted into six audio-image representations following [14], [23]. Standard normalization techniques were applied to both modalities.

Audio image representations: Following are some of the key characteristics of audio-image representations (shown in Figure 3).

- Audio image representations provide a significant reduction in dimensionality. For example, spectral centroid images represent the frequency content of the audio signal over time, which is a lower-dimensional representation of the original video dataset. This can make it easier and faster to process the data and extract meaningful features.
- Audio images are based on the audio signal, which is less affected by visual changes, such as changes in lighting conditions or camera angles. This makes these representations more robust to visual changes and can improve the accuracy of human action analysis.
- Standardization as audio images can be standardized to a fixed size and format, which can make it easier to compare and combine data from diverse sources. This can be useful for tasks such as cross-dataset validation and transfer learning. Hence, this dataset can serve as a standard benchmark for evaluating the performance of different machine-learning algorithms for human action analysis based on audio signals.
- Suitable for privacy-oriented applications such as surveillance or healthcare monitoring, which may require the analysis of human actions without capturing the original visual information.

Architecture: The MAiVAR-T model comprises an audio transformer, a video transformer, and a cross-modal attention layer. The transformers process the audio and video inputs separately, after which the cross-modal attention layer fuses the outputs. Finally, a classification layer predicts the action present in the input data.

Audio Stream: The audio stream uses Vision Transformer (ViT) [24] to process 2D images with minimal changes. In particular, ViT extracts N non-overlapping image patches,

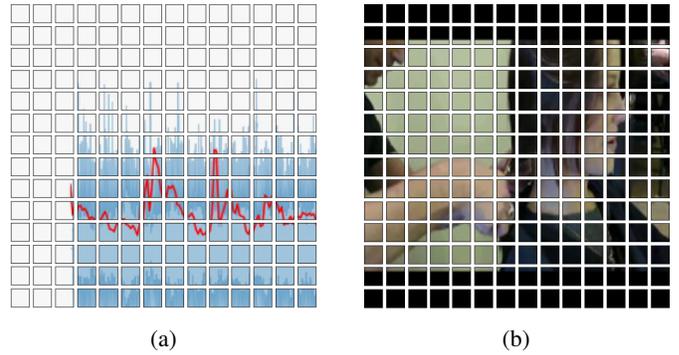


Fig. 2: Image patches (a) Audio-image representation, (b) RGB video frame.

$x_i \in \mathbb{R}^{h \times w}$, performs a linear projection and then rasterises them into 1D tokens $z_i \in \mathbb{R}^d$. The sequence of tokens input to the following transformer encoder is

$$\mathbf{z} = [z_{cls}, \mathbf{E}x_1, \mathbf{E}x_2, \dots, \mathbf{E}x_N] + \mathbf{p}, \quad (1)$$

where the projection by \mathbf{E} is equivalent to a 2D convolution. In addition, a learned positional embedding, $p \in \mathbb{R}^{N \times d}$, is added to the tokens to retain positional information, as the subsequent self-attention operations in the transformer are permutation invariant. The tokens are then passed through an encoder consisting of a sequence of L transformer layers. The MLP consists of two linear projections separated by a GELU non-linearity and the token-dimensionality, d , remains fixed throughout all layers. Finally, a linear classifier is used to classify the encoded input based on $z_{cls}^L \in \mathbb{R}^d$, if it was prepended to the input, or a global average pooling of all the tokens, z^L , otherwise. As the transformer [12], which forms the basis of ViT [17], is a flexible architecture that can operate on any sequence of input tokens $z \in \mathbb{R}^{N \times d}$, we describe strategies for tokenising videos next.

Video Feature Stream: We consider mapping a video $\mathbb{V} \in \mathbb{R}^{T \times H \times W \times C}$ to a sequence of tokens $z' \in \mathbb{R}^{n_t \times n_h \times n_w \times d}$. We then add the positional embedding and reshape into $\mathbb{R}^{N \times d}$ to obtain z , the input to the transformer.

IV. EXPERIMENTS

A. Audio preprocessing

Each audio image representation was broken into patches as illustrated in the examples shown in Figure 2. For spatial context, positional embeddings for each input were projected into the architecture (see Figure 4). An internal schematic of the transformer model has been illustrated in Figure 5. Training data was batched into mini-batches of 16 instances each. Augmentation techniques like random cropping and time-stretching were applied to increase model robustness.

B. Video preprocessing

Following [18], the features extracted are then fed to the multimodal fusion module (AV-Fusion MLP) which later performs the classification for each action class.

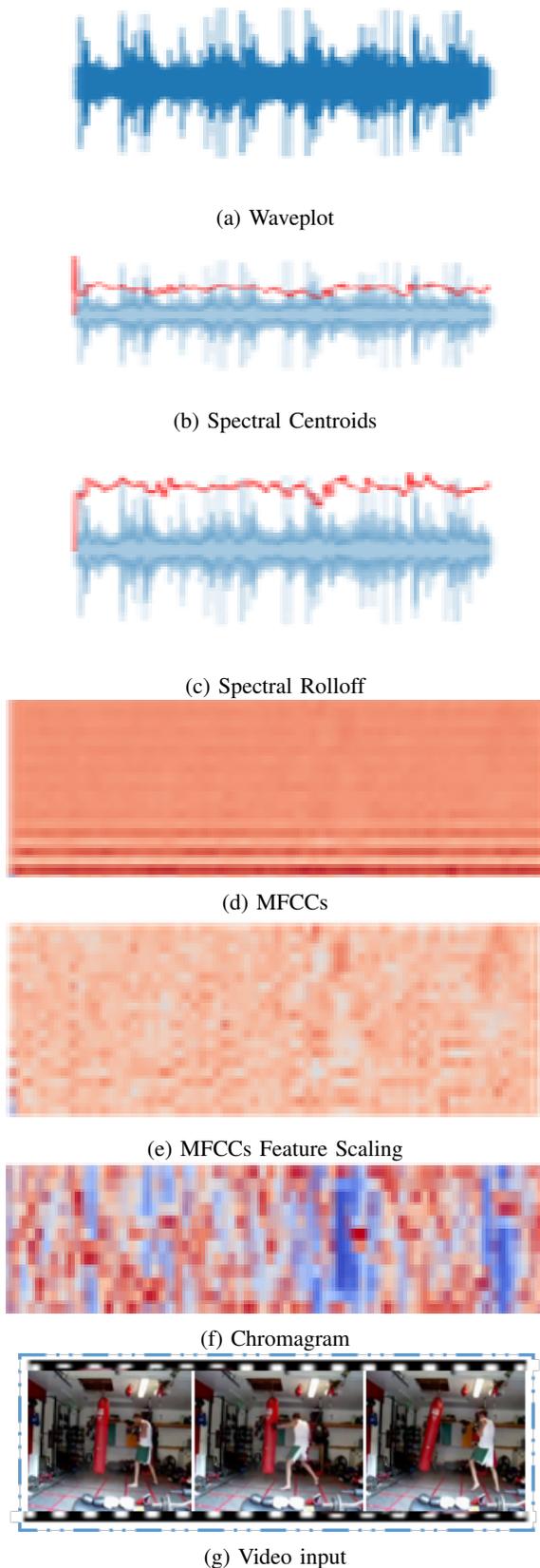


Fig. 3: Segmented video input and six different audio-image representations of the same action.

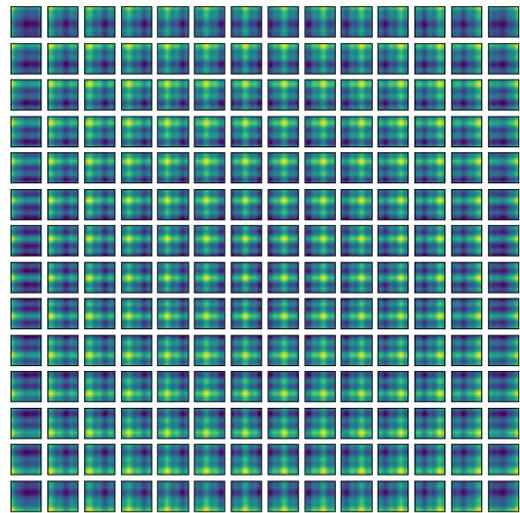


Fig. 4: Positional embeddings.

C. Training

We utilized a multimodal cross-entropy loss function for training, balancing both audio and video modalities. The network hyperparameters are reported in Table II.

Hardware and Schedule: The training was performed on a high-performance computing cluster, equipped with GeForce GTX 1080 Ti GPUs. We trained the transformer-based model for 100 epochs, with a learning rate (α) schedule that decreased the rate by 10% every 4 epochs. **Optimizer:** The Adam optimizer [25] was used due to its effectiveness in training deep networks. **Regularization:** Dropout techniques [26] were applied to prevent overfitting during training.

V. RESULTS

To assess the contribution of each component in our model, we performed an ablation study. Results demonstrate that both the audio and video transformers, as well as the cross-modal attention layer, contribute significantly to the final action recognition performance. The process of attention mechanism in the extraction of features through robust audio-image representations could be visualized in Figures 6 and 7. We have used an accuracy metric that measures the proportion of correct predictions made by the model out of all the predictions and defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad (2)$$

where TP are the correctly predicted positive values. TN are the correctly predicted negative values. FP , also known as Type I errors, are the negative values incorrectly predicted as positive. FN , also known as Type II errors, are the positive values incorrectly predicted as negative.

Table III compares the performance of transformer-based feature extractors with CNN-based counterparts. Proposed MAiVAR-T outperforms prior methods by a +3% as presented in Table IV.

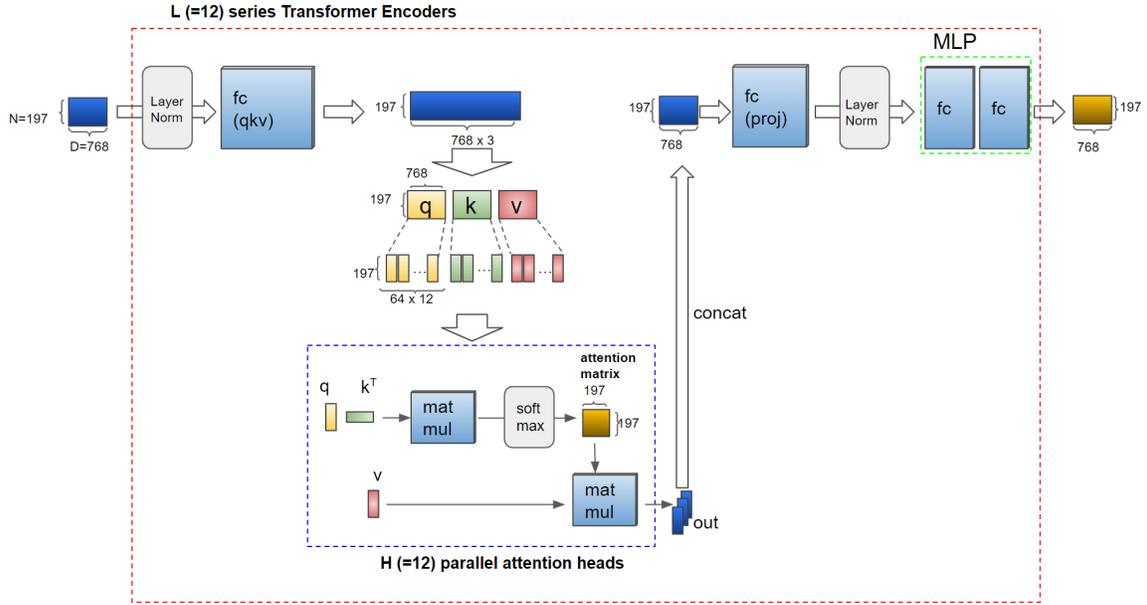


Fig. 5: Schematic of Vision Transformer Encoder.

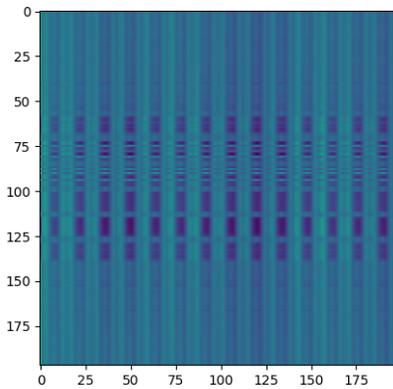


Fig. 6: Attention matrix for an audio-image representation.

TABLE III: Test accuracy of different audio representations with CNN and transformer-based backbones (InceptionResNet-v4(IRV4) and Vision Transformer (ViT) respectively)

Representation	IRV4	ViT
Waveplot	12.08	19.7 (+7)
Spectral Centroids	13.22	28.65 (+15)
Spectral Rolloff	16.46	26.85 (+10)
MFCCs	12.96	18.26 (+6)
MFCCs Feature Scaling	17.43	17.44 (+0.01)
Chromagram	15.48	19.08 (+3)

VI. CONCLUSION

Over the past decade, Convolutional Neural Networks (CNNs) with video-based modalities have been a staple in the

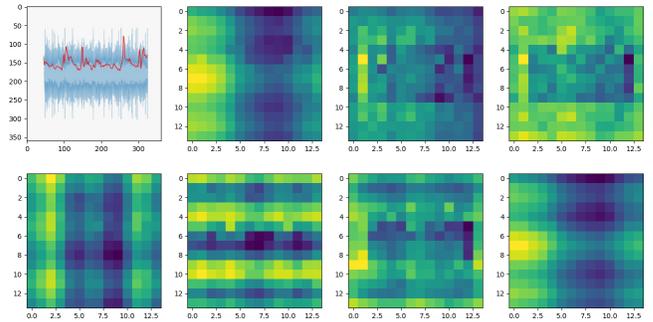


Fig. 7: Visualization of attention.

field of action video classification. However, in this paper, we challenge the indispensability of video modalities and propose a transformer-based multi-modal audio-image to video action recognition framework called Multi-modal Audioimage-Video Action Recognizer using Transformers (MAiVAR-T). This fusion-based, end-to-end model for audio-video classification features a transformer-based architecture that not only simplifies the model but also enhances its performance.

Experimental results demonstrate that our transformer-based audio-image to video fusion methods hold their own against traditional image-only methods, as corroborated by previous research. Given the significant improvements observed with pre-training on larger video datasets, there is considerable potential for further enhancing our model's performance. In future work, we aim to validate the efficacy of integrating text modality with audio and visual modalities. Furthermore, the scalability of MAiVAR-T on large-scale audio-video action recognition datasets, such as Kinetics 400/600/700 will be explored. Additionally, we plan to explore better architectural

TABLE IV: Classification accuracy of MAiVAR compared to the state-of-the-art methods on UCF51 dataset after fusion of audio and video features.

YEAR	METHOD	ACCURACY [%]
2015	C3D [27]	82.23
2016	TSN (RGB) [28]	60.77
2017	C3D+AENet [29]	85.33
2018	DMRN [30]	81.04
2018	DMRN [30] + [31] features	82.93
2020	Attention Cluster [32]	84.79
2020	IMGAUD2VID [6]	81.10
2022	STA-TSN (RGB) [33]	82.1
2022	MAFnet [31]	86.72
2022	MAiVAR-WP [14]	86.21
2022	MAiVAR-SC [14]	86.26
2022	MAiVAR-SR [14]	86.00
2022	MAiVAR-MFCC [14]	83.95
2022	MAiVAR-MFS [14]	86.11
2022	MAiVAR-CH [14]	87.91
Ours	MAiVAR-T	91.2

designs to integrate our proposed approach with more innovative ideas, such as integrating generative AI-based transformer architectures, into our network could provide valuable insights into the impact of transformers on MHAR.

ACKNOWLEDGMENT

This work is jointly supported by Edith Cowan University (ECU) and the Higher Education Commission (HEC) of Pakistan under Project #PM/HRDI-UESTPs/UETs-I/Phase-1/Batch-VI/2018. Dr. Akhtar is a recipient of Office of National Intelligence National Intelligence Postdoctoral Grant # NIPG-2021-001 funded by the Australian Government.

REFERENCES

[1] H. Park, Z. J. Wang, N. Das, A. S. Paul, P. Perumalla, Z. Zhou, and D. H. Chau, "SkeletonVis: Interactive visualization for understanding adversarial attacks on human action recognition models," in *Proc. of AAAI*, vol. 35, no. 18, 2021, pp. 16 094–16 096. 1

[2] G.-Z. Yang, J. Bellingham, P. E. Dupont, P. Fischer, L. Floridi, R. Full, N. Jacobstein, V. Kumar, M. McNutt, R. Merrifield, B. J. Nelson, B. Scassellati, M. Taddeo, R. Taylor, M. Veloso, Z. L. Wang, and R. Wood, "The grand challenges of science robotics," *Science Robotics*, vol. 3, no. 14, p. eaar7650, 2018. 1

[3] H. Oinas-Kukkonen and M. Harjumaa, "Persuasive systems design: key issues, process model and system features 1," in *Routledge handbook of policy design*. Routledge, 2018, pp. 87–105. 1

[4] R. Liu, A. A. Ramli, H. Zhang, E. Henricson, and X. Liu, "An overview of human activity recognition using wearable sensors: Healthcare and artificial intelligence," in *Proc. of Internet of Things-ICIOT*. Springer, 2022, pp. 1–14. 1

[5] C. Li, Q. Zhong, D. Xie, and S. Pu, "Collaborative spatiotemporal feature learning for video action recognition," in *Proc. of CVPR*, 2019, pp. 7872–7881. 1

[6] R. Gao *et al.*, "Listen to look: Action recognition by previewing audio," in *Proc. of CVPR*. IEEE, 2020, pp. 10 457–10 467. 1, 6

[7] M. B. Shaikh and D. Chai, "RGB-D data-based action recognition: a review," *Sensors*, vol. 21, no. 12, p. 4246, 2021. 1

[8] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989. 1

[9] 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. 1

[10] K. He, X. Zhang, S. Ren *et al.*, "Deep residual learning for image recognition," in *Proc. of CVPR*. IEEE, 2016, pp. 770–778. 1

[11] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997. 1, 2

[12] 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. 1, 2, 3

[13] K. Soomro, A. R. Zamir, and M. Shah, "UCF101: A dataset of 101 human actions classes from videos in the wild," *arXiv preprint arXiv:1212.0402*, 2012, doi:10.48550/arXiv.1212.0402. 1, 3

[14] M. B. Shaikh, D. Chai, S. M. S. Islam, and N. Akhtar, "MAiVAR: Multimodal audio-image and video action recognizer," in *Proc. of VCIP*. IEEE, 2022, pp. 1–5, doi:10.1109/VCIP56404.2022.10008833. 1, 3, 6

[15] J. Lin, C. Gan, and S. Han, "TSM: Temporal shift module for efficient video understanding," in *Proc. of ICCV*, 2019, pp. 7083–7093. 1

[16] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018. 2

[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. 2, 3

[18] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid, "ViViT: A video vision transformer," in *Proc. of ICCV*, October 2021, pp. 6836–6846. 2, 3

[19] D. Ramachandram and G. W. Taylor, "Deep multimodal learning: A survey on recent advances and trends," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 96–108, 2017. 2

[20] T. Chen and R. Rao, "Audio-visual integration in multimodal communication," *Proc. of the IEEE*, vol. 86, no. 5, pp. 837–852, 1998. 2

[21] Y. Kim, H. Lee, and E. M. Provost, "Deep learning for robust feature generation in audiovisual emotion recognition," in *Proc. of ICASSP*. IEEE, 2013, pp. 3687–3691. 2

[22] E. Kazakos *et al.*, "EPIC-Fusion: Audio-visual temporal binding for egocentric action recognition," in *Proc. of ICCV*, 2019, pp. 5492–5501. 2

[23] M. B. Shaikh, D. Chai, S. M. S. Islam, and N. Akhtar, "PyMAiVAR: An open-source python suite for audio-image representation in human action recognition," *Software Impacts*, p. 100544, 2023. 3

[24] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021. 3

[25] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014, doi:10.48550/arXiv.1412.6980. 4

[26] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *JMLR*, vol. 15, no. 1, pp. 1929–1958, 2014. 4

[27] D. Tran, L. Bourdev, R. Fergus *et al.*, "Learning spatiotemporal features with 3d convolutional networks," in *Proc. of ICCV*, 2015, pp. 4489–4497, doi:10.1109/ICCV.2015.510. 6

[28] L. Wang *et al.*, "Temporal segment networks: Towards good practices for deep action recognition," in *Proc. of the ECCV*, 2016, pp. 20–36, doi:10.1007/978-3-319-46484-8_2. 6

[29] N. Takahashi, M. Gygli, and L. Van Gool, "AENet: Learning deep audio features for video analysis," *IEEE TMM*, vol. 20, no. 3, pp. 513–524, 2017. 6

[30] Y. Tian, J. Shi, B. Li, Z. Duan, and C. Xu, "Audio-visual event localization in unconstrained videos," in *Proc. of ECCV*, 2018, pp. 247–263. 6

[31] M. Brousmiche, J. Rouat, and S. Dupont, "Multimodal attentive fusion network for audio-visual event recognition," *Information Fusion*, vol. 85, pp. 52–59, 2022. 6

[32] X. Long, G. De Melo, D. He, F. Li, Z. Chi, S. Wen, and C. Gan, "Purely attention based local feature integration for video classification," *IEEE TPAMI*, pp. 2140 – 2154, 2020. 6

[33] G. Yang *et al.*, "STA-TSN: Spatial-temporal attention temporal segment network for action recognition in video," *PLoS one*, vol. 17, no. 3, pp. 1–19, 2022. 6