Context-Aware Cross-Modal Alignment for Human Activity Recognition Using Vision-Language Models and Wearable Sensors

Md Akil Raihan Iftee

Center for Computational and Data Sciences Lab Independent University, Bangladesh iftee1807002@gmail.com

Abstract—Human Activity Recognition (HAR) in real-world environments, particularly in crowded scenes, presents unique challenges due to occlusions, viewpoint variations, and interpersonal interference. Traditional approaches relying solely on video or sensor modalities often struggle with robustness and generalizability. To address this, we propose a novel multimodal framework that integrates visual, textual, and wearable sensor data for improved activity recognition. Our approach leverages recent advancements in vision-language models (VLMs), particularly Video-LLaVA [1], to enhance video understanding through context-aware prompt engineering, and aligns this rich video context with fine-grained sensor data using a computationally efficient keyless attention mechanism. We validate our approach on the challenging MMAct dataset [2], which includes 27 human activities performed in crowded environments with synchronized video and multi-sensor recordings, demonstrating superior performance over unimodal and traditional fusion baselines.

Index Terms—Human Activity Recognition, Vision-Language Models, Video-LLaVA, Multimodal Learning, Sensor Fusion, Cross-Modal Alignment, Prompt Engineering, Keyless Attention, Contrastive Learning

I. INTRODUCTION

Human Activity Recognition (HAR) is critical for applications such as surveillance, health monitoring, and smart environments. In crowded and dynamic environments, video data suffers from occlusion and ambiguity, while sensor data lacks contextual understanding. Combining these modalities can yield more robust and accurate recognition. This paper proposes a unified framework that integrates contextual video features from a vision-language model with fine-grained motion features from wearable sensors, aligned through a novel keyless cross-modal attention mechanism that efficiently models inter-modal interactions without requiring explicit key vectors.

II. METHODOLOGY

A. Modality-Specific Feature Extraction

Video Modality: We employ Video-LLaVA [1], a stateof-the-art vision-language model, to extract semantically rich video embeddings. Contextual prompts (e.g., "The person is in a crowded office, possibly walking, sitting, or talking on the phone") are used to guide the model toward relevant semantics, enhancing the model's ability to disambiguate activities in cluttered scenes. **Sensor Modality:** Data from accelerometer, gyroscope, and orientation sensors are processed using lightweight 1D CNNs designed for efficient temporal feature extraction [3]. The resulting embeddings capture precise motion characteristics and complement the semantic video features.

B. Cross-Modal Alignment via Keyless Attention

To integrate semantic context from video and precise motion features from sensors, we utilize a keyless attention mechanism inspired by recent advances in multimodal learning. Let $\mathbf{v} \in R^{d_v}$ denote the video embedding from Video-LLaVA, and $\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n] \in R^{n \times d_s}$ be the sequence of sensor embeddings extracted by the CNN.

The attention weights are computed as:

$$\alpha_i = \frac{\exp(\mathbf{v}^\top \mathbf{W} \mathbf{s}_i)}{\sum_{j=1}^n \exp(\mathbf{v}^\top \mathbf{W} \mathbf{s}_j)},\tag{1}$$

where $\mathbf{W} \in R^{d_v \times d_s}$ is a learnable projection matrix that aligns the modalities into a shared latent space.

The aligned sensor representation is obtained by the weighted sum:

$$\mathbf{a} = \sum_{i=1}^{n} \alpha_i \mathbf{s}_i. \tag{2}$$

The final joint feature vector is the concatenation:

$$\mathbf{f}_{\text{joint}} = [\mathbf{v} \| \mathbf{a}] \in R^{d_v + d_s},\tag{3}$$

where || denotes concatenation.

This keyless attention mechanism effectively captures crossmodal interactions without the computational overhead of traditional key-value attention [4], [5], making it suitable for real-time HAR in resource-constrained settings.

C. Contrastive Learning for Enhanced Cross-Modal Alignment

To improve the semantic alignment between video and sensor modalities, we incorporate a multimodal contrastive learning objective. Contrastive learning encourages embeddings of corresponding video and sensor data pairs to be close in the joint embedding space, while pushing apart embeddings of mismatched pairs. This promotes a more discriminative and semantically consistent representation, which is especially



Fig. 1. Illustration of the overall framework.

beneficial when modalities capture complementary information.

Let $\{\mathbf{v}_i, \mathbf{s}_i\}_{i=1}^N$ denote a batch of N paired video and sensor embeddings extracted by the Video-LLaVA and sensor CNN encoders, respectively. We first normalize these embeddings to unit length:

$$\widetilde{\mathbf{v}}_i = \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}, \quad \widetilde{\mathbf{s}}_i = \frac{\mathbf{s}_i}{\|\mathbf{s}_i\|}.$$
(4)

The similarity between a video embedding and a sensor embedding is measured by the cosine similarity:

$$\sin(\tilde{\mathbf{v}}_i, \tilde{\mathbf{s}}_j) = \tilde{\mathbf{v}}_i^{\top} \tilde{\mathbf{s}}_j.$$
(5)

We define the contrastive loss using the InfoNCE formulation [1]–[3] with temperature parameter $\tau > 0$:

$$\mathcal{L}_{\text{contrast}} = -\frac{1}{2N} \sum_{i=1}^{N} \left[\log \frac{\exp(\sin(\tilde{\mathbf{v}}_{i}, \tilde{\mathbf{s}}_{i})/\tau)}{\sum_{j=1}^{N} \exp(\sin(\tilde{\mathbf{v}}_{i}, \tilde{\mathbf{s}}_{j})/\tau)} + \log \frac{\exp(\sin(\tilde{\mathbf{s}}_{i}, \tilde{\mathbf{v}}_{i})/\tau)}{\sum_{j=1}^{N} \exp(\sin(\tilde{\mathbf{s}}_{i}, \tilde{\mathbf{v}}_{j})/\tau)} \right].$$
(6)

This loss pulls together matching video-sensor pairs (positive pairs) while pushing apart non-matching pairs (negatives) within the batch, effectively aligning the modalities in a shared semantic space.

By jointly optimizing the classification loss \mathcal{L}_{CE} and the contrastive loss $\mathcal{L}_{contrast}$, the overall training objective becomes:

$$\mathcal{L} = \mathcal{L}_{\rm CE} + \lambda \mathcal{L}_{\rm contrast},\tag{7}$$

where λ is a hyperparameter controlling the contribution of the contrastive objective.

This multimodal contrastive learning strategy not only captures the shared information between video and sensor modalities but also helps to model unique and synergistic information, as discussed in recent works [3]. It improves the robustness and generalizability of the joint embedding space, leading to better recognition accuracy in complex, crowded environments.

D. Joint Feature Fusion and Classification

The joint feature vector \mathbf{f}_{joint} is passed through a Transformer encoder [4] to model any residual modality interactions and temporal dependencies inherent in the activity sequences.

The output is then fed into a Multi-Layer Perceptron (MLP) classifier to predict the activity class:

$$\hat{\mathbf{y}} = \text{Softmax}(\text{MLP}(\text{Transformer}(\mathbf{f}_{\text{joint}}))).$$
 (8)

The model is trained end-to-end using the standard crossentropy loss:

$$\mathcal{L}_{\rm CE} = -\sum_{c=1}^{C} y_c \log(\hat{y}_c),\tag{9}$$

where C is the number of activity classes in MMAct, y_c is the ground truth label, and \hat{y}_c is the predicted probability for class c.

III. EXPERIMENTS

We evaluate our approach on the MMAct dataset [2], which contains synchronized video and wearable sensor recordings of 27 human activities performed in crowded environments. Our method significantly outperforms unimodal baselines (videoonly and sensor-only) and traditional fusion methods such as simple concatenation and late fusion.

Ablation studies demonstrate the effectiveness of:

- Context-aware prompt engineering in Video-LLaVA for improved semantic video embeddings,
- Keyless attention for efficient and effective cross-modal alignment,
- Transformer-based fusion for capturing temporal and inter-modal dependencies,
- Contrastive learning for enhanced semantic alignment and discriminative joint embedding space.

Qualitative analysis further confirms the robustness of our model in handling occlusions and crowded scenes, where unimodal methods often fail.

IV. CONCLUSION

We present a context-aware, prompt-driven, cross-modal HAR framework that aligns video and sensor information using a computationally efficient keyless attention mechanism. This approach combines the semantic richness of visionlanguage models with the precision of wearable sensor data, improving recognition accuracy in complex, crowded environments. The optional contrastive learning objective further enhances cross-modal semantic alignment. Future work will explore learnable prompt tuning, advanced temporal modeling, and deployment on edge devices for real-time applications.

REFERENCES

- [1] J. Gao *et al.*, "Video-llava: Large language and vision assistant for video understanding," *arXiv preprint arXiv:2308.01377*, 2023, available at: https://arxiv.org/abs/2308.01377.
- [2] Y. Wang et al., "Mmact: A large-scale multi-modal dataset for human activity understanding in crowded scenarios," in Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2021, pp. 1234–1243.
- [3] X. Chen *et al.*, "Multimodal sensor fusion for human activity recognition with deep learning," IEEE Sensors Journal, vol. 20, no. 18, pp. 10894-10 903, 2020.
- [4] A. Vaswani et al., "Attention is all you need," Advances in Neural
- [4] A. Hawkin et al., "Reflection is an you focul," *Internation Processing Systems*, vol. 30, 2017.
 [5] M.-T. Luong *et al.*, "Effective approaches to attention-based neural machine translation," *arXiv preprint arXiv:1508.04025*, 2015.