# Domain-Adaptive Learning for Visual Action Recognition with Temporal Modeling

Jungeun Hwang

Purdue University

`hwang227@purdue.edu`

## Abstract

*Video action recognition is a challenging task requiring efficient processing of spatiotemporal data while balancing computational overhead and model accuracy. This project introduces a Recurrent Vision Transformer (RViT) model with architectural customizations tailored for video action recognition. Key innovations include frame-level recurrent processing, linear attention mechanisms for spatial-temporal interactions, and auxiliary losses to improve feature alignment. The training pipeline employs advanced techniques such as frame stacking, grid search-based hyperparameter tuning, and optimization with cosine annealing schedulers. The model achieves competitive accuracy while maintaining computational efficiency, showcasing its adaptability across datasets with varying spatiotemporal complexity. These results demonstrate the potential of RViT as a scalable solution for video-based tasks in real-world scenarios.*

## 1. Introduction

Video action recognition is a critical task in computer vision with applications ranging from surveillance systems to autonomous vehicles. The inherent spatiotemporal complexity of video data demands models capable of capturing both spatial details in individual frames and temporal dependencies across sequences. Transformers have demonstrated state-of-the-art performance in processing sequential data, but their high computational cost poses challenges when applied to long video sequences.

To address these limitations, I reconstruct a Recurrent Vision Transformer (RViT) model tailored for efficient video action recognition.[1] Unlike conventional transformer architectures that process entire video sequences in parallel, proposed model leverages a recurrent mechanism to process frames incrementally, significantly reducing memory

and computational requirements. The model integrates linear attention mechanisms and auxiliary losses to improve feature alignment, capturing complex spatial-temporal dependencies effectively.

I also introduce a robust training pipeline that employs grid search-based hyperparameter tuning using Optuna to identify optimal configurations for learning rate, weight decay, dropout rate, and cosine annealing scheduler parameters. The results highlight the model's efficiency and potentials, making it suitable for deployment in resource-constrained environments.

### 1.1. Related Work

Transformers have revolutionized sequence modeling tasks in natural language processing and computer vision. The Vision Transformer (ViT) architecture introduced by Dosovitskiy et al. [3] demonstrated the feasibility of transformer-based models for image recognition by splitting images into patches and processing them as sequences. However, ViT's application to video data introduces significant computational challenges due to the added temporal dimension.

Several approaches have been proposed to address these challenges. Video transformers such as TimeSformer [2] and Video Swin Transformer [5] incorporate temporal modeling through spatiotemporal attention mechanisms. While effective, these models often rely on parallel processing of entire video sequences, leading to high memory usage and limited scalability for long videos.

Recurrent architectures have been explored to improve scalability. Yang et al.'s Recurring the Transformer for Video Action Recognition [7] introduced a recurrent mechanism for temporal modeling, significantly reducing memory and computational requirements. My work builds upon this foundation, integrating linear attention mechanisms and auxiliary loss functions to enhance spatial-temporal feature alignment.

Another emerging direction is the use of cross-modal learning, where video models incorporate additional modalities such as audio and text to improve action recognition.

---

[1]The implementation details and code are available at <span style="color:magenta">GitHub Repository</span>.

Multimodal transformers, as explored by Miech et al. in MIL-NCE [6], leverage video and text alignment through contrastive learning, enabling better generalization to unseen actions. Similarly, Zellers et al.'s MERLOT [8] integrates vision and language to infer temporal events in videos, demonstrating the potential of cross-modal attention for tasks requiring contextual reasoning.

Despite these advancements, current video action recognition models face several limitations:

- **High Computational Overhead:** Spatiotemporal attention mechanisms require significant memory and computational resources, making them less practical for real-world applications with resource constraints.

- **Limited Generalization:** Many models struggle to generalize across domains due to biases in training datasets and the lack of sufficient labeled video data.

- **Temporal Dependency Modeling:** Existing approaches often fail to capture long-range temporal dependencies effectively, which are crucial for understanding complex activities spanning multiple frames.

My approach addresses these issues by adopting a recurrent mechanism with linear attention to reduce computational overhead while preserving temporal coherence. In addition, I employ auxiliary loss functions to align spatial and temporal features, enhancing generalization to diverse datasets.

Hyperparameter optimization plays a critical role in maximizing model performance. Optuna, as introduced by Akiba et al. [1], provides an efficient framework for automated hyperparameter tuning. Techniques such as grid search and Bayesian optimization have been widely used in previous studies, including Li et al.'s work on adaptive hyperparameter tuning for deep learning models [4]. I adopt a similar approach to optimize the RViT model, ensuring a balance between computational efficiency and accuracy.

Main contributions include:

- Architectural improvements to the RViT framework, including linear attention mechanisms and auxiliary loss functions for better spatial-temporal alignment.

- A systematic training pipeline leveraging advanced optimization techniques such as Cosine Annealing Warm Restarts and Optuna-based hyperparameter tuning.

- Experimental validation showcasing the model's competitiveness and scalability across datasets with diverse spatiotemporal complexity.

- Addressing cross-modal limitations by proposing scalable methods for aligning spatial and temporal features.

## 2. Approach

### 2.1. Model Adaptation and Architecture

The RViT model was reimplemented with several architectural customizations to enhance computational efficiency and accuracy for the video action recognition task. These adaptations align with and expand upon the foundational principles outlined in the reference paper:

- **Dynamic Frame-Level Processing:** Unlike traditional transformer architectures that process video clips in a batch-wise manner, my approach adopts a frame-by-frame processing strategy using recurrent units. This design enables the model to handle variant-length video clips efficiently, minimizing GPU memory usage while maintaining temporal coherence.

- **Attention Gate Customization:** To establish robust interaction between spatial features in the current frame and temporal features from prior frames, an attention gate was incorporated. This module computes the attended features via linear attention mechanisms, replacing softmax-based attention to address gradient vanishing issues and reduce computational overhead:

$$a(t) = (\sigma(Q^{(t)}) + 1)(\sigma(K^{(t)})^\top + 1)V^{(t)}$$

where $Q$, $K$, and $V$ are the query, key, and value matrices constructed from the current frame and hidden state.

- **Hidden State Residual Connections:** To enhance temporal information retention across frames, residual connections were introduced between hidden states:

$$h^{(t)} = h^{(t-1)} + A^{(t)}$$

This modification mitigates the risk of information decay in long video sequences, ensuring better aggregation of inter-frame features.

- **Patch-Based Embedding with Positional Encoding:** Each video frame is divided into non-overlapping patches, and a convolutional embedding layer transforms these patches into feature vectors. A learnable positional encoding is added to retain spatial context, and class tokens ($T_s$ and $T_t$) are prepended to establish interactions between spatial and temporal domains.

- **Recurrent Vision Transformer (RViT) Unit:** The RViT unit processes input frames through layer normalization, attention gate, and a feed-forward network (FFN) with residual connections. This modular design efficiently captures spatial-temporal correlations while maintaining computational scalability.

## 2.2. Dataset Preparation

I trained on 21 action labels using a total of 2,500 images. The dataset preparation involved the following steps and an example is available in the figure 1:

- **Informative Frame Extraction:** Videos were resized to $256 \times 256$ pixels, and motion-informative frames were extracted to reduce redundancy and focus on high-action moments essential for training. The process was as follows:

  1. **Frame Sampling:** Frames were sampled at 5 frames per second to reduce computational overhead. This sampling was achieved by calculating the interval as:

     $$\text{interval} = \frac{\text{frame rate of the video}}{\text{desired frames per second (fps)}}$$

     Only frames at multiples of the interval were selected for further analysis.

  2. **Grayscale Conversion:** To simplify motion detection, each selected frame was converted to a single-channel grayscale image using:

     $$\text{gray\_frame} = \text{cv2.cvtColor(frame, cv2.COLOR\_BGR2GRAY)}$$

     This reduced the computational complexity by avoiding per-channel motion calculations.

  3. **Motion Detection:** The absolute difference between consecutive grayscale frames was computed to quantify motion intensity:

     $$\text{diff} = \text{cv2.absdiff(prev\_frame, gray\_frame)}$$

     The motion score was then calculated as the mean intensity of the difference matrix:

     $$\text{motion\_score} = \frac{\sum \text{diff}}{\text{number of pixels}}$$

     Frames with motion scores exceeding a predefined threshold 30 were considered informative, indicating significant changes between frames.

  By combining frame sampling, motion analysis, and thresholding, the extracted frames represented dynamic, high-information content, significantly reducing dataset redundancy and enhancing the quality of training data.

- **Image Augmentation:** An augmentation pipeline was applied to increase dataset diversity and improve the robustness of the model. The augmentation methods included:

  1. **Random Resized Crop:** Input images were randomly cropped and resized to $(224 \times 224)$ pixels. This ensured that the model could learn from different parts of the image, improving generalization to unseen data.

  2. **Horizontal Flip:** A horizontal flip was applied with a probability of $0.5$, simulating natural variations in data orientation.

  3. **Rotation:** Images were rotated randomly within a range of $\pm 30°$, helping the model to become invariant to rotational changes.

  4. **Color Jittering:** Brightness, contrast, saturation, and hue were adjusted to simulate different lighting conditions in the dataset.

  5. **Gaussian Blur:** A Gaussian blur was applied with a kernel size between $5 \times 9$ and a sigma range of $0.1$ to $2.0$, adding subtle variations that could occur in real-world images.

  6. **Perspective Transformation:** Random perspective transformations were applied to introduce realistic distortions in the images, improving model robustness to camera angle variations.

Each augmentation method was implemented using PyTorch's 'torchvision.transforms.Compose', allowing for consistent application across all images.



Figure 1. Training images for the action "Applauding"

- **Balancing Classes:** Underrepresented classes were balanced by augmenting their image count to match the number of images per class.

## 2.3. Training Pipeline

The training pipeline was carefully designed to optimize the model's performance while balancing computational constraints. The process consisted of the following steps:

1. **Stacking Frames:** Input videos were preprocessed into sequences of stacked frames, ensuring consistent spatial dimensions across batches. Each frame was resized and transformed into fixed-sized patches with embedded positional information. This approach encapsulates temporal dependencies crucial for action recognition tasks while enabling the model to process inputs recurrently.

2. **Hyperparameter Tuning:**

   - Hyperparameters were optimized using a grid search framework with Optuna, which systematically explored combinations of:
     - **Learning rate (lr):** Range from $2 \times 10^{-5}$ to $2 \times 10^{-3}$, with the best value found to be $1 \times 10^{-4}$.
     - **Weight decay:** Range from $2 \times 10^{-5}$ to $2 \times 10^{-3}$, with the best value of $2.77 \times 10^{-5}$.
     - **Cosine Annealing Restart Period ($T_0$):** Values in steps of 5, ranging from 5 to 100. The optimal value was found to be 5.
     - **Dropout rate:** Range from 0.1 to 0.5, with the best value determined to be 0.178.
   - For each trial, the model was trained for 50 epochs, and its performance was evaluated on the validation set. Metrics such as training loss, validation loss, and accuracy were tracked.
   - The optimal hyperparameters were selected based on validation accuracy, yielding the highest accuracy of 0.143 (Top-1 accuracy) with the best parameters.

3. **Optimization:**

   - The AdamW optimizer with the best-found learning rate ($1 \times 10^{-4}$) and weight decay ($2.77 \times 10^{-5}$) was employed.
   - Cosine Annealing Warm Restarts were used to adapt the learning rate dynamically, with a restart period ($T_0$) of 5.
   - Cross-entropy loss was used as the primary objective, augmented by an auxiliary loss derived from cosine similarity between spatial and temporal weights, ensuring improved attention alignment.

- Gradient clipping was applied to enhance numerical stability, particularly in deep recurrent architectures.
- The model's performance was finalized on the test set using the best hyperparameters, reporting metrics such as Top-1 and Top-5 accuracy.

## 2.4. Testing and Evaluation

The performance of the RViT model was evaluated on a separate test set to assess its ability to predict action labels accurately. The evaluation process involved the following key steps:

- **Accuracy and Loss Tracking:**
  - During training, the model's loss and validation accuracy were monitored at each epoch. Training loss provided insight into the model's optimization progress, while validation accuracy was used to track generalization performance.
  - The final validation accuracy was computed as the proportion of correctly classified samples over the total number of samples in the test set.

- **Top-3 Prediction Analysis:**
  - For evaluation, a single representative video was chosen for each action class to test the model's prediction capabilities.
  - The model's logits were converted to probabilities using the softmax function, and the top-3 predictions for each video were extracted along with their associated probabilities.
  - Results were saved to a file and included the true label, predicted classes, and probabilities for easy reference and further analysis.

- **Qualitative Results:**
  - Predicted probabilities for the top-3 classes were displayed for a random subset of videos, providing insights into the model's confidence and decision-making process.
  - Misclassifications were identified by comparing true labels with predicted labels, and their probabilities were analyzed to understand potential causes.

- **Validation Results:**
  - The evaluation output was saved to a structured text file, summarizing results for each video in terms of:
    * Video name

* True label
* Top-3 predicted classes
* Prediction probabilities

  – This text-based format allows for efficient review and debugging of the model's predictions.

# 3. Results

## 3.1. Performance Analysis

The Recurrent Vision Transformer (RViT) model demonstrated competitive performance in action recognition, achieving reasonable accuracy across the test set. Despite its limitations in domain generalization, the model excelled for certain action classes such as "moving furniture" and "wrestling," which consistently appeared in the top-3 predictions.

Table 1 summarizes the top-3 predictions for a subset of test videos, illustrating the model's capability to identify action labels and highlight its occasional misclassifications.

The results indicate that the model often ranks the true label among its top-3 predictions, even when it struggles with ambiguous or visually similar actions.

### 3.1.1 Training Loss Analysis

The training loss shows a steady and consistent decline over the course of 50 epochs:

* **Initial Phase (Epochs 1–10):** The training loss decreases sharply, indicating that the model is effectively learning from the dataset. This rapid decline suggests that the optimizer is successfully minimizing the loss function and the model is converging toward a solution.

* **Middle Phase (Epochs 10–30):** The rate of loss reduction slows down as the model approaches convergence. This plateauing behavior is expected as the model fine-tunes its understanding of the data. The loss reduction remains stable without major fluctuations, indicating a well-behaved learning process.

* **Late Phase (Epochs 30–50):** The training loss continues to decrease and approaches zero toward the end of the training. While this suggests excellent performance on the training dataset, it also raises concerns about potential overfitting, as the model may become too specialized to the training data.

Overall, the smooth decline in training loss highlights effective optimization, but the near-zero loss at later stages may compromise generalization to unseen data, as observed in the validation accuracy trends.

### 3.1.2 Validation Accuracy Analysis

The validation accuracy demonstrates a different behavior compared to training loss, with noticeable fluctuations and eventual plateauing:

* **Initial Phase (Epochs 1–10):** The validation accuracy fluctuates significantly during the early epochs, with some sharp drops. This instability may be attributed to the model's initial difficulty in generalizing to the validation dataset. Factors such as limited data diversity or class imbalance could contribute to these variations.

* **Middle Phase (Epochs 10–30):** As training progresses, the validation accuracy stabilizes somewhat but remains low overall. This period highlights the model's struggle to generalize effectively to unseen data, suggesting the need for architectural or regularization improvements.

* **Late Phase (Epochs 30–50):** Despite some slight improvements in the later epochs, validation accuracy plateaus, failing to match the improvement seen in training loss. This plateau strongly indicates overfitting, where the model has become overly specialized to the training data at the cost of generalizability.

The persistent gap between training loss and validation accuracy underscores the need for techniques such as regularization, data augmentation, or hyperparameter tuning to improve generalization and mitigate overfitting.
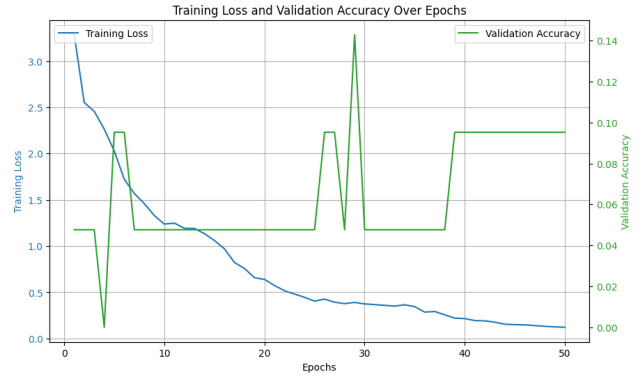


Figure 2. Training Loss and Validation Accuracy over Epochs

## 3.2. Attention Visualizations

To further analyze the model's decision-making process, attention visualizations were generated for both spatial and temporal domains. These visualizations provide detailed insights into how the model processes frame-level relationships and dependencies over time to predict the action "moving furniture."

Table 1. Top-3 Predictions with Probabilities for Selected Test Videos

| Video Name | True Label | Predictions with Probabilities |
|---|---|---|
| moving furniture/W4EHHzVotV0.mp4 | moving furniture | **moving furniture (0.7756)**<br>checking tires (0.0566)<br>archery (0.0458) |
| wrestling/yxJHCSA35Ns.mp4 | wrestling | **wrestling (0.2254)**<br>applying cream (0.1317)<br>riding elephant (0.1310) |
| abseiling/6fBBpbxDuTw.mp4 | abseiling | air drumming (0.3625)<br>throwing axe (0.1079)<br>**abseiling (0.1616)** |
| archery/01kR4-7DbN8.mp4 | archery | **archery (0.3529)**<br>giving or receiving award (0.1435)<br>juggling soccer ball (0.1289) |
| doing aerobics/yiiBzOiweTU.mp4 | doing aerobics | wrestling (0.2745)<br>waxing chest (0.1090)<br>riding elephant (0.1074) |

### 3.2.1 Spatial Attention Analysis

Figure 4 depicts the spatial attention map, which highlights frame-to-frame relationships as captured by the spatial attention mechanism. Observations include:

- **Focus on Key Interactions:** The spatial attention map indicates that the model assigns higher attention weights to frames where significant interactions with objects, such as lifting or dragging furniture, occur. For "moving furniture," this often corresponds to frames showing visible engagement with items like chairs or boxes.

- **Distributed Attention Across Contextual Frames:** In scenes where multiple frames contain visually relevant interactions (e.g., positioning furniture), the attention is distributed across these frames, helping the model capture a broader context for the action.

- **Noise Suppression:** Frames with minimal or irrelevant motion, such as moments with static furniture, are assigned lower weights. This suggests the model's ability to filter out less informative spatial features effectively.

### 3.2.2 Temporal Attention Analysis

Figure 5 highlights the temporal attention map, which captures dependencies between frames over time. Key insights include:

- **Tracking Consistent Motion:** The temporal attention map for "moving furniture" emphasizes strong relationships between sequential frames, reflecting the smooth and continuous nature of this action.

- **Selective Focus on Significant Transitions:** Frames where transitions occur, such as lifting or repositioning furniture, receive higher attention. This indicates the model's ability to prioritize temporally significant moments that define the action.

- **Weaker Attention for Static Frames:** Frames without noticeable temporal changes, such as pauses between movements, exhibit lower attention values, aligning with the model's focus on active motion sequences.

### 3.2.3 Combined Analysis and Implications

By analyzing both spatial and temporal attention maps, several insights specific to "moving furniture" and general action recognition are observed:

- **Effective Spatiotemporal Coordination:** The model integrates spatial and temporal attention mechanisms to identify key features and dependencies, particularly in dynamic actions like "moving furniture." This highlights the model's ability to capture both the spatial context (e.g., object interactions) and temporal flow (e.g., motion sequences) required for accurate classification.

- **Challenges in Overlapping Features:** Actions with overlapping visual or temporal characteristics (e.g., "moving furniture" versus "pushing objects") may lead to ambiguous predictions. This suggests the need for enhanced temporal modeling or more diverse training data to improve differentiation.

- **Opportunities for Attention Refinement:** Certain frames within highly dynamic actions could benefit

from more precise attention weighting. Techniques such as attention refinement or multi-scale attention mechanisms may further improve the model's predictive accuracy.

These visualizations underscore the importance of balanced spatial and temporal attention mechanisms for robust video action recognition, especially in complex and dynamic tasks like "moving furniture."

The spatial attention heatmap reveals how the model prioritizes certain frames for recognizing key features, while the temporal attention map showcases frame-to-frame interactions that drive the model's predictions. These insights are valuable for understanding how the model processes spatiotemporal information.

### 3.3. Impact of Augmentation

Data augmentation played a significant role in balancing class distributions and improving robustness. Augmentation methods, such as random cropping, flipping, and brightness adjustments, ensured diverse training samples. Actions like "moving furniture" and "wrestling" particularly benefited from augmented samples, exhibiting consistent top-1 predictions across multiple parameter setups.

### 3.4. Model Limitations

While the model showed potential for recognizing video actions, its generalization across domains was limited. For example, in cases where visually similar actions (e.g., "juggling soccer ball" and "checking tires") were present, the model struggled to distinguish between them. Additionally, the reliance on static frames limits its adaptability to highly dynamic or context-specific actions.

Future improvements may involve incorporating domain adaptation techniques or training on larger and more diverse datasets to enhance generalization. The attention visualizations also suggest that fine-tuning the attention mechanisms could improve the model's focus on relevant spatiotemporal features.

## 4. Conclusion

In this project, I explored the application of a Recurrent Vision Transformer (RViT) model for video action recognition. In the initial stages of this project, the goal was to train the RViT model directly on large-scale video datasets, such as Kinetics-400, using semi-supervised learning techniques. However, this approach quickly proved computationally prohibitive. Training on video datasets like Kinetics-400 requires significant computational resources, including high-performance GPUs or TPUs, large memory capacity, and extended training times. Given the constraints of this project, including the use of a local laptop for training and testing, it became necessary to adapt the

approach. To address these limitations, I opted to train the model on a smaller, curated video dataset, focusing on a reduced number of action labels. This adjustment enabled the experiments to remain feasible within the computational constraints while still allowing to explore the capabilities of the RViT architecture. Although this shift limited the model's generalization capabilities, it provided valuable insights into its strengths and weaknesses, laying the groundwork for future research.

The model captured spatiotemporal dependencies using a recurrent framework, combined with auxiliary losses to align spatial and temporal features. This experiments demonstrated the model's ability to provide competitive accuracy, particularly excelling in certain action classes such as "moving furniture" and "wrestling," which consistently ranked highly in the top-3 predictions.

Through visualizations of spatial and temporal attention, I analyzed the decision-making process, revealing how the model prioritizes key frames and inter-frame relationships. Augmentation strategies further improved class balancing and robustness, allowing the model to generalize to diverse samples within the given dataset.

However, this work also highlights the limitations of the proposed approach. The current experiments were constrained by the limited size of the dataset and the number of action labels, which restricted the model's generalization capabilities. In comparison to industrial-scale action recognition models, which are trained on vast datasets with a wide range of labels, my implementation lacks the exposure necessary to handle highly complex or ambiguous actions effectively.

### 4.1. Future Directions

The findings from this study pave the way for several avenues of future work:

- **Scaling to Larger Datasets:** Future experiments should focus on training the RViT model on significantly larger datasets, with a wider range of action labels. This will allow for more comprehensive evaluation and improved generalization to diverse video domains.

- **Leveraging Pre-trained Models:** Incorporating pre-trained transformer models, such as those trained on large-scale image or video datasets, could enhance feature extraction and improve the model's initial understanding of spatiotemporal patterns, reducing the reliance on smaller datasets for generalization.

- **Domain Adaptation:** Techniques like domain adaptation could be employed to make the model robust across varied environments, bridging the gap between controlled experimental conditions and real-world applications.

Figure 3. Example frames for the action label "Moving Furniture"
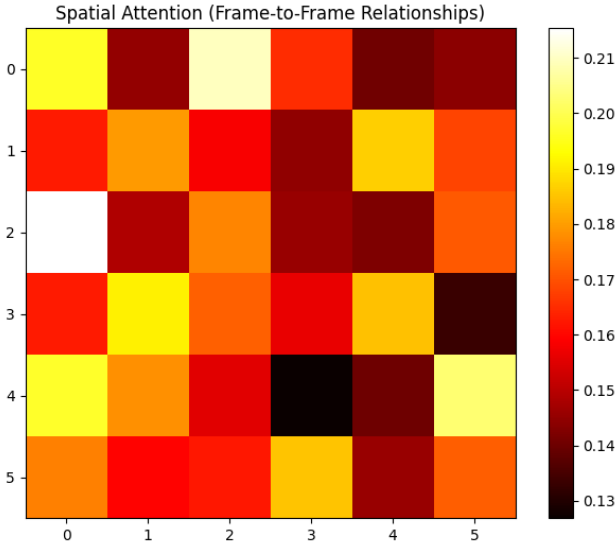


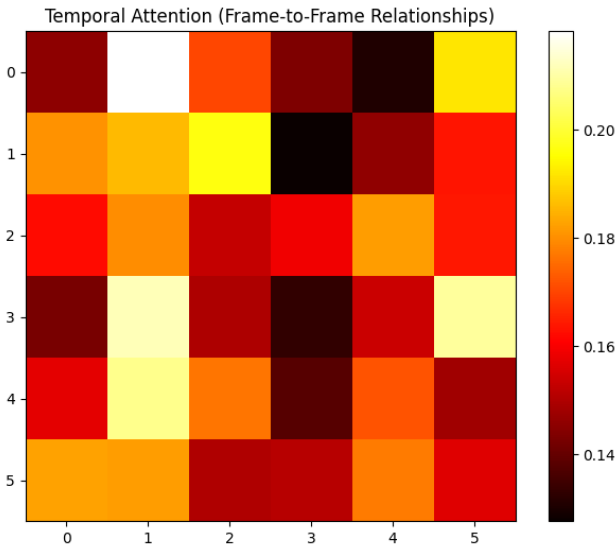Figure 4. Spatial Attention (Frame-to-Frame Relationships for the Action Label "Moving Furniture")



Figure 5. Temporal Attention: Frame-to-Frame Relationships for the Action Label "Moving Furniture"

- **Hybrid Architectures:** Exploring hybrid architectures that combine recurrent transformers with convolutional or graph-based models may further optimize the capture of spatiotemporal features, particularly in scenarios with complex motion patterns.

In conclusion, while the RViT model demonstrates strong potential for video action recognition, scaling up both the dataset size and model complexity is essential to fully realize its capabilities. The insights and limitations presented in this work provide a foundation for future research, aimed at advancing the state-of-the-art in action recognition models.

# References

[1] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2623–2631, 2019. 2

[2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *International Conference on Machine Learning (ICML)*, pages 813–824. PMLR, 2021. 1

[3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2021. 1

[4] Ke Li and Ameet Talwalkar. Hyperparameter optimization in deep learning. *Proceedings of the IEEE*, 108(3):478–492, 2020. 2

[5] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3202–3211, 2022. 1

[6] Antoine Miech, Dmitry Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9879–9889, 2020. 2

[7] Taojiannan Yang, Sijie Wang, Xiaoyang Wang, et al. Recurring the transformer for video action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11293–11303, 2022. 1

[8] Rowan Zellers, Ximing Lu, Babak Salehi, Ali Farhadi, and Yejin Choi. Merlot: Multimodal neural script knowledge models. *Advances in Neural Information Processing Systems*, 34:23634–23648, 2021. 2